

A Grouping Genetic Algorithm for Joint Stratification and Sample Allocation Designs

Mervyn O’Luing,¹ Steven Prestwich,¹ S. Armagan Tarim²

Abstract

Finding the optimal stratification and sample size in univariate and multivariate sample design is hard when the population frame is large. There are alternative ways of modelling and solving this problem, and one of the most natural uses genetic algorithms (GA) combined with the Bethel-Chromy evaluation algorithm. The GA iteratively searches for the minimum sample size necessary to meet precision constraints in partitionings of atomic strata created by the Cartesian product of auxiliary variables. We point out a drawback with classical GAs when applied to the grouping problem, and propose a new GA approach using “grouping” genetic operators instead of traditional operators. Experiments show a significant improvement in solution quality for similar computational effort.

Keywords: Grouping genetic algorithm; Optimal stratification; Sample allocation; R software.

¹ *Insight Centre for Data Analytics, Department of Computer Science, University College Cork, Ireland.* Email: mervyn.oluing@insight-centre.org, steven.prestwich@insight-centre.org

² *Cork University Business School, University College Cork, Ireland.* Email: armagan.tarim@ucc.ie

1 Introduction

In this paper we address the optimization problem of jointly determining stratification and sample allocation for univariate and multivariate scenarios. To serve this purpose, we refer to (Ballin and Barcaroli, 2013). In principle the optimal stratification (i.e. that which yields the smallest sample size) can be found by testing all possible partitionings of *atomic strata*, but the number of possible partitionings grows exponentially with the number of atomic strata.

An efficient search algorithm is necessary to avoid evaluating each possible partitioning. Genetic algorithms (GAs) often converge quickly to optimal or near optimal solutions, and are particularly good at navigating rugged search spaces containing many local minima. The Bethel-Chromy algorithm combines similar algorithms from (Bethel, 1985, 1989) and (Chromy, 1987) and is suitable for univariate and multivariate cases. It uses lagrangian multipliers to find the minimum sample size that meets precision constraints for a given stratification. (Ballin and Barcaroli, 2013) combine a GA with this algorithm to search for the minimum sample size. It is used to evaluate each partitioning created by the GA. A full description of the methodology and problem statement is found in (Ballin and Barcaroli, 2013). However, they use a classical GA which is known to be unsuitable for partitioning problems.

In this paper we propose to apply genetic operators to the GA that are better suited to this application. It is an example of the class of evolutionary algorithms called *Grouping Genetic Algorithms* (GGAs). The GA has been updated following this work (Barcaroli, 2019). Section 2 motivates the work and introduces GGAs. Section 2.3 describes our GGA for the problem. Section 3 compares the original GA with our GGA on publicly-available test data. Section 4 describes a version of our GGA with enhanced performance, using a fast C++ implementation of the *bethel.r* function which we integrated into R using the Rcpp package. Section 5 concludes the paper.

44 2 Classical vs grouping genetic algorithms

45 In this section we discuss “classical” and “grouping” GAs, and explain why the latter are
46 more appropriate for our problem.

47 2.1 Classical genetic algorithms

48 GAs are a nature-inspired class of optimisation algorithms, modelled on the ability of
49 organisms to solve the complex problem of adaptation to life on Earth. The variables of
50 an optimisation problem are called *genes* and their values *alleles*. A candidate solution is
51 a list of alleles called a *chromosome*. A set of chromosomes is usually called a *population*,
52 so to avoid confusion with the target population we shall use *chromosome population* when
53 referring to GAs. The objective function (which is maximised by convention) is called the
54 chromosome’s *fitness*. The search for fit chromosomes (solutions with high objective) uses
55 two *genetic operators*: small random changes called *mutation*, equivalent to small local
56 moves in a hill-climbing algorithm; and large changes called *crossover* in which the genes
57 of two *parent chromosomes* are *recombined*. One well-known recombination operator is
58 *single-point crossover*: choose two *parent* chromosomes with alleles

$$a_1, \dots, a_N \quad b_1, \dots, b_N$$

59 select a random integer i (the *crossover point*) such that $1 \leq i < N$, and generate two
60 new *offspring* chromosomes

$$a_1, \dots, a_i, b_{i+1}, \dots, b_N \quad b_1, \dots, b_i, a_{i+1}, \dots, a_N$$

61 These might be further subjected to random *mutation*, in which a few alleles are changed,
62 before placing them back into the chromosome population. There are a variety of methods
63 for selecting parents and replacing existing chromosomes. In *generational* GAs the entire
64 chromosome population is replaced by offspring, and parents are often selected randomly

65 but with a bias toward fitter chromosomes; while in *steady-state* GAs only one offspring
66 is generated in each GA iteration, and usually replaces the least-fit chromosome in the
67 chromosome population. GAs often give more robust results than search algorithms
68 based on hill-climbing, because of their use of recombination. They have found many
69 applications since their introduction in 1975 by John Holland.

70 The original GA which is represented in the *R* (R Core Team, 2015) package *Sam-*
71 *plingStrata* (Barcaroli et al., 2014), is an elitist generational GA in which the atomic
72 strata L are considered to be elements of a set (or genes) for a standard crossover strat-
73 egy. In each iteration the best solutions (the *elite*) are carried over to the next generation.
74 Each gene represents a variable in the problem. We refer to this as a classical GA because
75 a classical problem representation and genetic operators are used, as described below.

76 Dividing atomic strata into disjoint groups is an example of a *grouping* problem,
77 related to *cutting*, *packing* and *partitioning* problems. The motivation for our work
78 is that classical GAs are known to perform poorly on grouping problems. The reason
79 is that the chromosomal representation of a grouping contains a great deal of *symmetry*
80 (or *redundancy*): permuting the group names yields an equivalent grouping, so each
81 grouping has multiple representations. Symmetry has a damaging effect on GAs because
82 recombining similar parent groupings might yield a very different offspring grouping,
83 violating the basic GA principle that parents should tend to produce offspring with similar
84 fitness. In extreme cases, a classical GA might perform even worse than a completely
85 random search. We provide two examples to illustrate the problem.

86 To illustrate the problem with symmetry in our first example the parents represent
87 the same grouping in different ways. Note that to increase readability, letters A - F are
88 used as alleles instead of integers in the presentation here. Consider the following two
89 chromosomes:

chromosome	groups represented					
	A	B	C	D	E	F
ABCDEF	{1}	{2}	{3}	{4}	{5}	{6}
FEDCBA	{6}	{5}	{4}	{3}	{2}	{1}

90 which both represent the grouping $\{\{1\}, \{2\}, \{3\}, \{4\}, \{5\}, \{6\}\}$. Now suppose we ap-
 91 ply single-point crossover to obtain two new offspring chromosomes from these parents.
 92 Arbitrarily choosing the center of the chromosomes as the crossover point, we obtain
 93 offspring:

chromosome	groups represented					
	A	B	C	D	E	F
ABCCBA	{1, 6}	{2, 5}	{3, 4}	\emptyset	\emptyset	\emptyset
FEDDEF	\emptyset	\emptyset	\emptyset	{3, 4}	{2, 5}	{1, 6}

94 which both represent the completely unrelated grouping $\{\{1, 6\}, \{2, 5\}, \{3, 4\}\}$: no groups
 95 at all are passed from the parents to the offspring. Hence the offspring and parent fitnesses
 96 can be completely unrelated to each other, which reduces the GA to near-random search.
 97 As another example, consider the following two classical chromosomes:

chromosome	groups represented					
	A	B	C	D	E	F
AECFEC	{1}	\emptyset	{3, 6}	\emptyset	{2, 5}	{4}
DFFDAA	{5, 6}	\emptyset	\emptyset	{1, 4}	\emptyset	{2, 3}

98 which in turn represent the different groupings $\{\{1\}, \{3, 6\}, \{2, 5\}, \{4\}\}$ and $\{\{5, 6\}, \{1, 4\}, \{2, 3\}\}$.
 99 Using the same crossover strategy we obtain offspring:

chromosome	groups represented					
	A	B	C	D	E	F
AECDAA	{1, 5, 6}	\emptyset	{3}	{4}	{2}	\emptyset
DFFDEC	\emptyset	\emptyset	{6}	{1}	{5}	{2, 3, 4}

104 representing the groupings $\{\{1, 5, 6\}, \{3\}, \{4\}, \{2\}\}$ and $\{\{6\}, \{1\}, \{5\}, \{2, 3, 4\}\}$. Note
105 that these offspring have very little in common with their parents, as the only preserved
106 groups are $\{1\}$ and $\{4\}$.

107 **2.2 Grouping genetic algorithms**

108 The symmetry problem can be tackled by designing more complex genetic representa-
109 tions and operators (Galinier and Hao, 1999) or by clustering techniques (Pelikan and
110 Goldberg, 2000). The risk of clustering is that genetic diversity may be lost if the clus-
111 ters are too tight, leading to search stagnation (Prügel-Bennett, 2004). Instead we follow
112 the former approach by designing a GGA (Falkenauer, 1998), which have been shown to
113 perform far better than classical GAs on grouping problems.

114 GGAs are designed specifically to solve grouping problems and have found many appli-
115 cations, including WiFi network deployment (Agustín-Blas et al., 2011), wireless network
116 design (Brown and Vroblefski, 2004), steel plate cutting (Hung, Sumichrast, and Brown,
117 2003), production plant layout (De Lit, Falkenauer, and Delchambre, 2000) and social
118 network analysis (James et al., 2010). They may use the same heuristics as other GAs
119 (parent selection, offspring replacement, etc) but they use different genetic encoding and
120 operators: that is, how they map a problem to chromosomes and how they perform
121 recombination and mutation. We shall illustrate these differences on the above examples.

122 GGAs represent a grouping as an ordered list of subsets, omitting empty sets. The
123 parents in the second example of Section 2.1 might be represented in this way:

$$\langle \{1\}, \{3, 6\}, \{2, 5\}, \{4\} \rangle \quad \langle \{5, 6\}, \{1, 4\}, \{2, 3\} \rangle$$

124 GGA mutation is simple: an item is moved from one group to another. However, the
125 GGA recombination operator is more complicated. Choose a *crossing section* in each
126 parent, for example $\langle \{1\}, \{3, 6\} \rangle$ from the 1st parent and $\langle \{1, 4\} \rangle$ from the 2nd parent.

127 Then *inject* the 1st crossing section into the 2nd parent at a random point, and vice-versa:

$$\langle \{1\}, \{3, 6\}, \underline{\{1, 4\}}, \{2, 5\}, \{4\} \rangle \quad \langle \{5, 6\}, \{1, 4\}, \underline{\{1\}}, \{3, 6\}, \{2, 3\} \rangle$$

128 Next remove any repeated objects that were already in the receiving parent:

$$\langle \emptyset, \{3, 6\}, \{1, 4\}, \{2, 5\}, \emptyset \rangle \quad \langle \{5\}, \{4\}, \{1\}, \{3, 6\}, \{2\} \rangle$$

129 Finally remove any empty sets:

$$\langle \{3, 6\}, \{1, 4\}, \{2, 5\} \rangle \quad \langle \{5\}, \{4\}, \{1\}, \{3, 6\}, \{2\} \rangle$$

130 These are the offspring. Clearly, both offspring have much in common with both parents,
131 as 5 of the 7 parent groups survive in the offspring: $\{1\}$, $\{4\}$, $\{1, 4\}$, $\{2, 5\}$ and $\{3, 6\}$. In
132 the first example of Section 2.1 it is easily verified that both offspring represent the same
133 grouping as the parents, as one would expect. This property of the GGA injection-based
134 recombination makes it much more likely that offspring have similar fitness to parents,
135 which in turn helps the GGA to iteratively improve the chromosome population.

136 It might be noticed that the GGA problem representation still contains symmetry:
137 any grouping still has multiple representations, obtained by permuting the subsets in the
138 ordered list. But the genetic operators are almost independent of this ordering so it is
139 almost irrelevant. The only effect of the ordering is to limit the set of possible injections:
140 in the second example of Section 2.1 we cannot inject a non-existent crossing section for
141 example such as $\langle \{1\}, \{4\} \rangle$ from parent 1 because those two groups are not adjacent.
142 This limit is removed by an additional genetic operator called *inversion* which selects a
143 section of the chromosome and reverses it. For example

$$\langle \{1\}, \{2\}, \underline{\{3, 6\}}, \{4\}, \{5\} \rangle \longrightarrow \langle \{1\}, \{2\}, \underline{\{5\}}, \{4\}, \{3, 6\} \rangle$$

144 This does not change the grouping represented by the chromosome, but reordering the
145 groups in the chromosome makes all injections possible.

146 Injection, mutation and inversion are the common operators used in GGAs, but there
147 is no canonical algorithm. Instead GGAs tend to be tailored for specific applications, and
148 in principle any GA can be adapted to grouping problems by using grouping operators.
149 In Section 2.3 we design a GGA for our problem.

150 **2.2.1 Note on implementation**

151 For the sake of clarity the descriptions in Section 2.2 omit implementation details, for
152 example the fact that GGA chromosomes are usually implemented in two parts (or
153 sometimes more). The first part uses a classical representation as above, while the second
154 part lists the nonempty groups as a permutation. Injection occurs on the second parts of
155 parent chromosomes and some renaming of groups is necessary.

156 Typically we decide in advance the number of iterations which we wish to run the
157 algorithm for. This should be enough to give the GGA a chance to converge on the
158 optimum solution after the mutation and inversion probabilities have been applied. If,
159 however, the optimum solution is known beforehand the algorithm can be set to stop at
160 this point.

161 The number of iterations is usually decided with experience of using the GGA on similar
162 target and auxiliary variables for similar datasets, or with the existing dataset and target
163 and auxiliary variables. It may require a number of experiments using the GGA (or
164 GA) before the number of iterations needed to reach convergence can be estimated. In
165 fact there is a possibility that either the GGA or GA would appear to have reached
166 convergence after a set number of iterations, but instead have become trapped in a local
167 minimum. It may be useful to increase the number of iterations and try alternative
168 mutation probabilities in order to be certain that it has converged on a global minimum.

169 This implies a number of trial runs before finally deciding the parameters under which
170 to run the algorithms. Therefore the fact the GGA has been shown to attain convergence

171 quicker than the GA is likely to compound the improvement in total processing time.
172 In the experiments described below we keep the number of iterations small as we want
173 to demonstrate the ability of the GGA to converge on a solution within that number of
174 iterations.

175 We use either the mutation settings specified in the examples provided by (Ballin and
176 Barcaroli, 2013) or the default mutation settings in (Barcaroli et al., 2014). We apply
177 grouping genetic operators and inversion to the GA designed by (Ballin and Barcaroli,
178 2013): it is the grouping genetic operators that make it a GGA. Thus we compare the
179 performance between the different GA and GGA genetic operators rather than experiment
180 with parameters such as varying the number of iterations, chromosome population size,
181 mutation probability, or elitism rate.

182 The mutation probability can be selected in advance by the user. Typically, the
183 probability of mutation should be such that it increases the chance of the GGA leaving a
184 local minimum, but not disrupt the natural evolution of chromosomes from one generation
185 to the next. On the other hand we have fixed the inversion probability at 0.01, because
186 this is enough to maintain diversity.

187 The size of the chromosome population can be decided by trial and error. It is advisable
188 to consider the evaluation time of each chromosome when setting the size: if there are
189 too many chromosomes in the set, it might take an extra long time to move from one
190 iteration to the next, and we found that the *bethel.r* algorithm (i.e. the Bethel-Chromy
191 evaluation algorithm in (Barcaroli et al., 2014)) takes several seconds to evaluate even
192 one chromosome for the larger datasets we used in this paper (we discuss this further in
193 Section 4).

194 For further details on the implementation of GGAs (e.g. elitism rate) we refer the
195 reader to papers such as (Falkenauer, 1998).

196 **2.3 Application to the joint stratification and sample allocation**
197 **problem**

198 As mentioned above our GGA is based on the GA described in (Ballin and Barcaroli,
199 2013) and represented in *R* in the *SamplingStrata* package (Barcaroli et al., 2014), but
200 with grouping operators and chromosomes instead of the classical versions. This change
201 is the only novelty of our algorithm (except for the optimisation described in Section 4)
202 but its effect on performance is large. We inserted the GGA into a modified version of the
203 function called *rbga.r* from the *genalg* R package (Willighagen, 2005). It is designed to
204 work with the other functions in *SamplingStrata*, and is applied to the joint stratification
205 and optimum sample size problem. The GGA is summarised in Figure 2.1.

Grouping Genetic Algorithm (GGA)

Step 1: Initialization

(a) Randomly generate a chromosome population of size N_p

Step 2: Selection part 1

- (a) Rank chromosomes based on sample size
- (b) Save best E chromosomes for the next generation

Step 3: Inversion

With probability 0.01 invert groups in the N_p chromosomes

Step 4: Selection part 2

For each of the remaining $N_p - E$ chromosomes in the new generation:

- (a) Draw parents 1 and 2 from the aforementioned N_p chromosomes
(higher ranked chromosomes have a higher probability of being selected)
- (b) Perform crossover as explained in Section 2.2
- (c) Remove empty groups
- (d) Renumber groups

Step 5: Mutation

Mutate integers in $N_p - E$ chromosomes at a selected probability

Step 6: if #iterations < maximum

(optional: and sample size > desired value) go to step 2

Figure 2.1: Pseudocode for our GGA

206 Following the problem statement in (Ballin and Barcaroli, 2013) we summarise the cost

207 function as follows:

$$C(n_1, \dots, n_H) = C_0 + \sum_{h=1}^H C_h n_h$$

208 where C_0 is the fixed cost and C_h is the average cost of interviewing one unit in stratum
 209 h and n_h is the number of units, or sample, allocated to stratum h . In our analysis C_0
 210 is set to 0, and C_h is set to 1. The expectation of the estimator of the “g-th” population
 211 total is:

$$E(\hat{T}_g) = \sum_{h=1}^H N_h \bar{Y}_{h,g} \quad (g = 1, \dots, G)$$

212 where $\bar{Y}_{h,g}$ is the mean of the G different target variables Y in each stratum h . The
 213 variance of the estimator is given by:

$$\text{VAR}(\hat{T}_g) = \sum_{h=1}^H N_h^2 \left(1 - \frac{n_h}{N_h}\right) \frac{S_{h,g}^2}{n_h} \quad (g = 1, \dots, G) \quad (1)$$

214 The upper limit of variance or precision U_g is expressed as a coefficient of variation CV
 215 for each \hat{T}_g :

$$CV(\hat{T}_g) = \frac{\sqrt{\text{VAR}(\hat{T}_g)}}{E(\hat{T}_g)} \leq U_g \quad (2)$$

216 The problem can be summarised as follows:

$$\begin{aligned} \min n &= \sum_{h=1}^H n_h \\ CV(\hat{T}_g) &\leq U_g \end{aligned}$$

217 **3 Comparing the Genetic Algorithms**

218 We now run a number of comparisons between the original GA and our GGA using
 219 publicly available datasets. Unless otherwise stated, for all the cases presented below, we
 220 adopt the following parameter setting for both genetic algorithms, where $N_P = 20$, U_g
 221 $\equiv 0.05$, the elitism rate is 0.2, and the mutation probability is 0.05.

222 3.1 A comparison for the iris dataset

223 (Ballin and Barcaroli, 2013) use the iris dataset (Anderson, 1935; Fisher, 1936; R Core
 224 Team, 2015) to demonstrate that the GA they propose can find the optimum stratification
 225 i.e. the stratification or grouping of atomic strata which supplies the minimum sample
 226 size. The iris dataset is small and is widely available. It has 150 observations for 5
 227 variables Sepal Length, Sepal Width, Petal Length, Petal Width and Species.

228 Species is a categorical variable which has three levels, setosa, versicolor and virginica,
 229 each of which have 50 observations. The remaining four variables are continuous mea-
 230 surements for length and width in centimetres. (Ballin and Barcaroli, 2013) select Petal
 231 Length and Petal Width as variables of interest, i.e. target variables. They select Sepal
 232 Length and Species as two auxiliary variables.

233 They convert Sepal Length to a categorical variable using a k-means algorithm (Har-
 234 tigan and Wong, 1979) to define three clusters (i.e. 4.3 to less than 5.5, 5.5 to less than
 235 6.5, 6.5 to 7.9). The cross product of the categorical version of Sepal Length with Species
 236 creates 9 atomic strata. However, one atomic atomic stratum is empty because there are
 237 no corresponding values in Petal Length and Petal Width. Therefore there are 8 usable
 238 atomic strata for this example.

Table 3.1: Reproduction of table of atomic strata for estimating the minimum sample size for the target variables of iris dataset as found in (Ballin and Barcaroli, 2013), P.379

Stratum	N	M1	M2	S1	S2	X1	X2	DOMAIN
[4.3; 5.5] (1)*setosa	45	1.466667	0.244444	0.17127	0.106574	[4.3; 5.5] (1)	setosa	1
[4.3; 5.5] (1)*versicolor	6	3.583333	1.166667	0.491313	0.205481	[4.3; 5.5] (1)	versicolor	1
[4.3; 5.5] (1)*virginica	1	4.5	1.7	0	0	[4.3; 5.5] (1)	virginica	1
[5.5; 6.5] (2)*setosa	5	1.42	0.26	0.172047	0.08	[5.5; 6.5] (2)	setosa	1
[5.5; 6.5] (2)*versicolor	35	4.268571	1.32	0.367051	0.189435	[5.5; 6.5] (2)	versicolor	1
[5.5; 6.5] (2)*virginica	23	5.230435	1.947826	0.318194	0.28873	[5.5; 6.5] (2)	virginica	1
[6.5; 7.9] (3)*versicolor	9	4.677778	1.455556	0.193091	0.106574	[6.5; 7.9] (3)	versicolor	1
[6.5; 7.9] (3)*virginica	26	5.876923	2.107692	0.494825	0.228579	[6.5; 7.9] (3)	virginica	1

239 The initial atomic strata are reproduced in Table 3.1 where M_g refers to the means

240 for the corresponding Y_g values in each atomic stratum l_k ; S_g refers to the corresponding
241 stratum population standard deviations. There are 4,140 possible partitionings of the 8
242 atomic strata. Consequently, it is possible to test within a reasonable amount of time
243 the sample size for the entire search space using the *bethel.r* function. This has already
244 been done (Ballin and Barcaroli, 2013) and the minimum sample size is known to be 11.

245 This test can be used to determine whether the new GA correctly finds the minimum
246 sample size without exploring the entire search space. We use $N_p = 10$ in this case.
247 For this test the *bethel.r* function will search for the minimum sample size, in integers
248 rather than real numbers. The chromosomes will then be ranked by sample size in
249 ascending order. Accordingly the elite chromosomes are taken into the next iteration
250 and the remaining chromosomes are generated using the recombination method for each
251 algorithm.

252 We will compare the number of chromosomes generated to find the optimal stratifica-
253 tion in the two algorithms as well as the number of iterations. Our anticipation is that
254 the GGA should be more efficient, and thus typically find the optimal solution in fewer
255 iterations than the GA.

256 The maximum number of iterations is set to 200, because using (Ballin and Barcaroli,
257 2013) as a guide we anticipate that both algorithms will find the correct solution in fewer
258 iterations than this. Thus we have added a piece of code to both algorithms such that
259 they stop when the optimal sample size, $n=11$, has been reached and supply the number
260 of iterations taken to reach that point. This approach is different to that of (Ballin and
261 Barcaroli, 2013) who report the number of times in 10 experiments the GA finds the
262 correct solution for a given number of iterations ranging incrementally from 25 to 200.
263 However, we feel this approach would better demonstrate that the GGA can find the
264 correct solution in less iterations even on the small iris dataset experiment.

265

Table 3.2: Iris dataset experiment results for GA and GGA

(a) GA			(b) GGA		
Number of			Number of		
Experiment	Iterations	Chromosomes	Experiment	Iterations	Chromosomes
1	14	228	1	11	180
2	8	132	2	7	116
3	17	276	3	6	100
4	40	644	4	22	356
5	31	500	5	9	148
6	13	212	6	11	180
7	15	244	7	8	132
8	9	148	8	7	116
9	15	244	9	9	148
10	15	244	10	11	180
11	14	228	11	3	52
12	8	132	12	9	148
13	17	276	13	27	436
14	40	644	14	12	196
15	31	500	15	16	260
16	13	212	16	6	100
17	15	244	17	20	324
18	9	148	18	6	100
19	15	244	19	7	116
20	15	244	20	6	100
21	16	260	21	11	180
22	67	1076	22	7	116
23	19	308	23	8	132
24	9	148	24	5	84
25	11	180	25	7	116
26	20	324	26	5	84
27	32	516	27	6	100
28	10	164	28	6	100
29	37	596	29	9	148
30	9	148	30	6	100

266 Table 3.2 provides the number of iterations (and chromosomes generated) taken to find
267 n=11 over 30 experiments for both GAs.

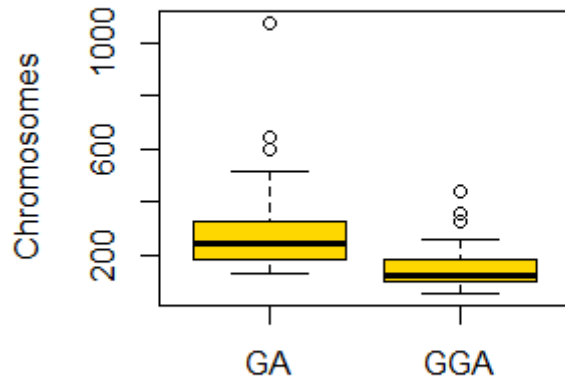


Figure 3.1: Boxplot distribution of number of Chromosomes generated to find $n=11$ for GA and GGA after 30 experiments

268 Figure 3.1 provides the distribution of the number of chromosomes generated to find
 269 the optimal solution for the GA and the GGA. The boxplots indicate that the GGA
 270 typically needs to generate fewer chromosomes to find the optimum solution.

Table 3.3: Example stratifications for the GA and GGA on the iris dataset for $n=11$

		Y1			Y2		
	Stratum	N	Mean	SD	Mean	SD	Sample Size
GA	1	50	1.462	0.1685	0.246	0.1026	2
	2	50	4.26	0.4562	1.326	0.1911	3
	3	1	4.5	0	1.7	0	1
	4	23	5.2304	0.3112	1.9478	0.2824	3
	5	26	5.8769	0.4852	2.1077	0.2241	2
Total		150					11
	Stratum	N	Mean	SD	Mean	SD	Sample Size
GGA	1	23	5.2304	0.3112	1.9478	0.2824	3
	2	50	1.462	0.1685	0.246	0.1026	2
	3	26	5.8769	0.4852	2.1077	0.2241	2
	4	51	4.2647	0.4529	1.3333	0.1962	4
Total		150					11

271 Table 3.3 provides example stratifications for the GA and GGA that both provide the
 272 optimal sample size necessary to meet precision constraints. (Ballin and Barcaroli, 2013)

273 indicate that a number of partitionings from the total of 4,140 possible partitionings
274 provide the minimum sample size. These range in size from 3 to 5 strata. It is seen
275 that the GGA results in fewer, less fragmented design strata. The same tendency can be
276 observed in the latter cases.

277 **3.2 Swiss municipality dataset**

278 The swissmunicipalities dataset provided by (Barcaroli et al., 2014) refers to the Swiss
279 municipalities in 2003. Each municipality belongs to one of seven regions which are
280 at the NUTS-2 level, i.e. equivalent to provinces. Each region contains a number of
281 cantons, which are administrative subdivisions. There are 26 cantons in Switzerland.
282 The data, which was sourced from the Swiss Federal Statistical Office and is included in
283 the *sampling* and *SamplingStrata* packages, contains 2,896 observations (each observation
284 refers to a Swiss municipality in 2003). They comprise 22 variables, details of which can
285 be examined in (Barcaroli et al., 2014).

286 The target estimates are the totals of the population by age class in each Swiss region.
287 In this case, the G target variables will be:

288 Y1: number of men and women aged between 0 and 19,

289 Y2: number of men and women aged between 20 and 39,

290 Y3: number of men and women aged between 40 and 64,

291 Y4: number of men and women aged 65 and over.

292 We consider 6 auxiliary variables, formed using the same k- means clustering method as
293 the iris dataset example:

294 X1: Classes of total population in the municipality. 18 categories.

295 X2: Classes of wood area in the municipality. 3 categories.

296 X3: Classes of area under cultivation in the municipality. 3 categories.

297 X4: Classes of mountain pasture area in the municipality. 3 categories.

298 X5: Classes of area with buildings in the municipality. 3 categories.

299 X6: Classes of industrial area in the municipality. 3 categories.

300 There are 7 regions, which we treat as population domains of design to distinguish them
301 from the design strata, replicating the experiment outlined in (Barcaroli et al., 2014).

302 The number of non-empty atomic strata is 641 in the population. We set the minimum
303 population size of stratum to be 2, and the maximum number of iterations to be 400.

304 The results for Sample Size and Strata after 30 experiments each with 400 iterations are
305 summarised in figure 3.2 below.

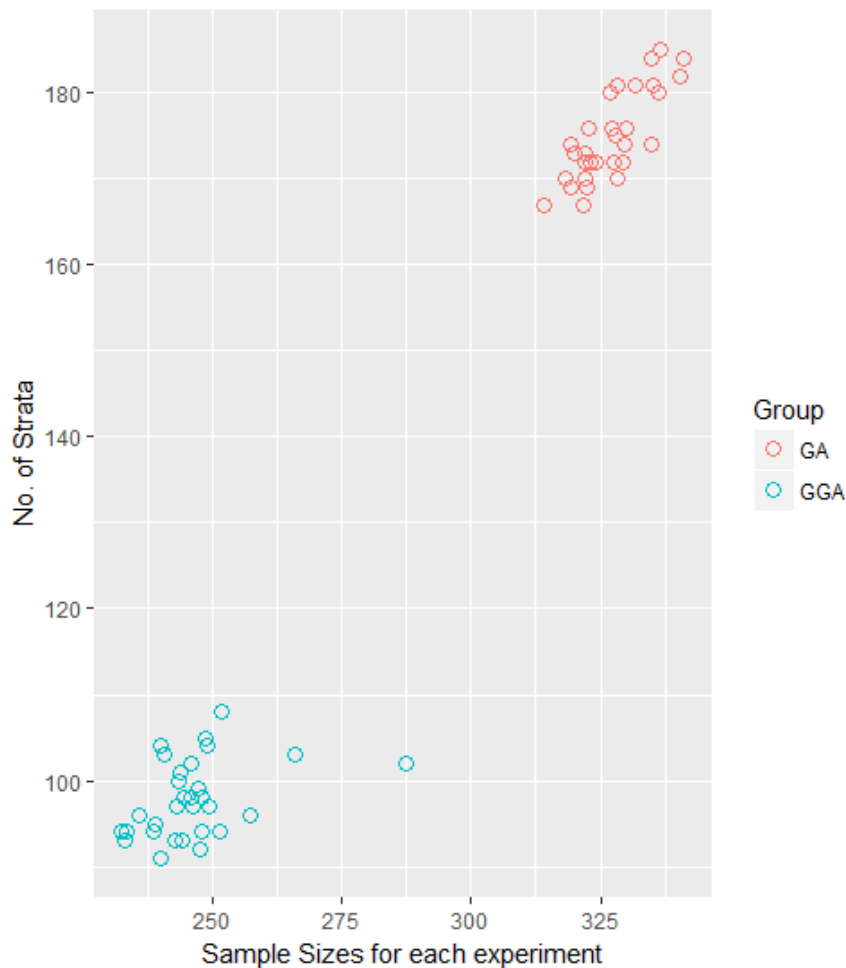


Figure 3.2: Scatterplot of Results for Strata v Sample size for GA and GGA after 30 experiments

306 Figure 3.2 clearly shows that the GGA returns a smaller sample size to the GA for
307 these settings. The median for the GGA, 246, is 25% lower than that for the GA, 328.

308 **3.3 2015 American Community Survey Public Use Microdata**

309 The United States has been conducting a decennial census since 1790. In the 20th
310 century censuses were split into long and short form versions. A subset of the population
311 was required to answer the longer version of the census, with the remainder answering
312 the shorter version. After the 2000 census the longer questionnaire became the annual
313 American Community Survey(ACS) (US Census Bureau, 2013). The 2015 ACS Public
314 Use Microdata Sample (PUMS) file (US Census Bureau, 2016) is a sample of actual
315 responses to the ACS representing 1% of the US population. The PUMS file contains
316 1,496,678 records each of which represents a unique housing unit or group quarters. There
317 are 235 variables. The full data dictionary is available in (US Census Bureau, 2016). We
318 selected the following to be target variables:

- 319 1. Household income (past 12 months)
- 320 2. Property value
- 321 3. Selected monthly owner costs
- 322 4. Fire/hazard/flood insurance (yearly amount).

323 and the following auxiliary variables:

- 324 1. Units in structure
- 325 2. Tenure
- 326 3. Work experience of householder and spouse
- 327 4. Work status of householder or spouse in family households
- 328 5. House heating fuel

329 6. When structure first built.

330 The PUMS data for which all the values are present contains 619,747 records. We use
331 the 51 states (based on census definitions) as domains.

332 In the convergence plots of Figure 3.3, the black line represents the best or lowest sample
333 size for the chromosome population in each iteration, whereas the red line represents the
334 mean sample size for the chromosome population in each iteration.

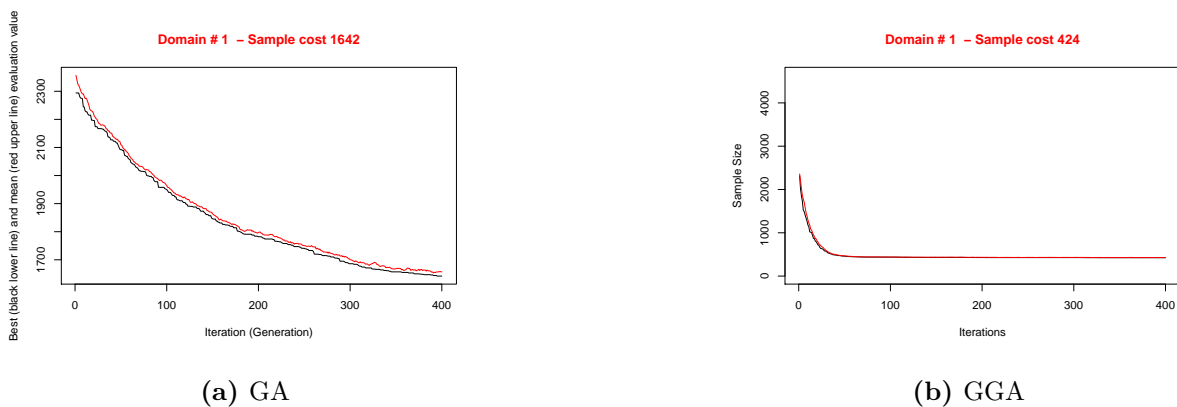


Figure 3.3: Convergence plots for Sample Size after the 1st experiment for GA and GGA. Note the different scales on the vertical axes

335 The GA appears to be reducing the sample size steadily but does not appear to have
336 reached a local minimum after 400 iterations. The GGA appears to have reached a local
337 or global minimum very quickly.

338 3.4 Kaggle Data Science for Good challenge Kiva Loans data

339 The online crowdfunding platform kiva.org provided a dataset of loans issued to people
340 living in poor and financially excluded circumstances around the world over a two year
341 period for a Kaggle Data Science for Good challenge. The dataset has 671,205 unique
342 records. We selected these target variables:

343 1. term in months

344 2. lender count

345 3. loan amount

346 and the following auxiliary variables:

347 1. sector

348 2. currency

349 3. activity

350 4. region

351 5. partner id

352 to create atomic strata. For these variables we removed any records with missing values.
353 We then proceeded to remove any countries with less than 10 records from the sampling
354 frame. This resulted in a sampling frame with 614,361 records. The variable country-code
355 defines the 73 design domains in this experiment.

Table 3.4: Sample size and strata for the Kiva Loans data from the GA and the GGA after 100 iterations

GA		GGA		Reduction	
Sample size	Strata	Sample size	Strata	Sample size	strata
78,018	43,030	11,963	1,793	84.67%	95.83%

356 Table 3.4 shows an 84.67% reduction in sample size and a 95.83% reduction in the
357 number of strata after 100 iterations. Figure 3.4 shows that for the same starting chro-
358 mosome population size for Domain 1 of the Kiva Loans dataset, the GGA attained a
359 good sample size in less than 100 iterations, but after 10,000 iterations the GA had not
360 converged and the sample size was still much higher than the GGA.

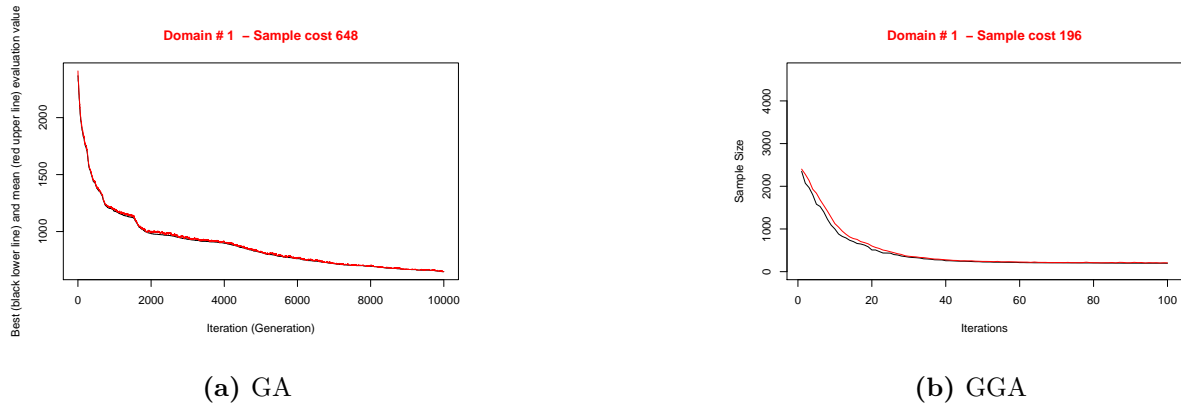


Figure 3.4: Convergence plots for Sample Size for the 1st Domain for GA (10,000 iterations) and GGA (100 iterations) in the Kiva Loans dataset experiment. Note the different scales on the vertical and horizontal axes

361 **3.5 UN Commodity Trade Statistics data**

362 Kaggle also hosts a copy of the UN Statistical Division Commodity Trade Statistics data.
 363 Trade records are available from 1962. We took a subset of data for the year 2011 and
 364 removed records with missing observations. This resulted in a data set with 351,057
 365 records. We selected the following target variable:

- 366 1. trade_usd

367 which refers to the value of trade in USD (US dollars), and the following auxiliary
 368 variables:

- 369 1. commodity
- 370 2. flow
- 371 3. category

372 The variable commodity is a categorical description of the type of commodity, e.g. Horses,
 373 live except pure-bred breeding. The variable flow describes whether the commodity was
 374 an import, export, re-import or re-export. The variable category describes the category

375 of commodity, e.g. silk or fertilisers. The 171 categories of country or area were selected
376 as domains.

Table 3.5: Sample size and strata for the UN Commodity Trade Statistics data from the GA and the GGA after 100 iterations

GA		GGA		Reduction	
Sample size	Strata	Sample size	Strata	Sample size	strata
288,638	191,000	84,181	16,555	70.84%	91.33%

377 **3.6 2000 US census data**

378 The Integrated Public Use Microdata Series extract is a 5% sample of the 2000 US census
379 data (Ruggles et al., 2017). The file contains 6,184,483 records. The US Census Data
380 will be very similar to the ACS data as the latter is an annual version of the former. But
381 for this experiment we selected different target and auxiliary variable combinations. The
382 single target variable in this test is usually a key focus of household surveys:

- 383 1. Total household income

384 We used the following information as auxiliary variables (note these are variables which
385 are likely available in administrative data):

- 386 1. Annual property insurance cost
- 387 2. Annual home heating fuel cost
- 388 3. Annual electricity cost
- 389 4. House value

390 The house value variable (VALUEH) reports the midpoint of house value intervals (e.g.
391 5,000 is the midpoint of the interval of less than 10,000), so we have treated it as a cate-
392 gorical variable. As with the 2015 ACS PUMS dataset we have taken a subset for which
393 all values are present. This has resulted in a subset with 627,611 records. The domain

394 for this experiment was Census region and division.

395

Table 3.6: Sample size and strata for the 2000 US census data by Census region and division from the GA and the GGA after 100 iterations

Division	Sampling frame		GA solution		GGA solution	
	Sampling Units	Atomic Strata	Sample sizes	Strata	Sample sizes	Strata
New England	116,045	87,084	81,012	52,628	376	58
Middle Atlantic	183,543	138,470	130,862	86,002	416	75
East North Central	65,480	58,055	53,075	35,794	327	42
West North Central	31,408	29,413	26,525	18,248	324	38
South Atlantic	97,189	83,357	76,716	51,457	440	49
East South Central	21,631	20,429	18,256	12,500	451	62
West South Central	22,582	20,919	18,750	12,730	407	39
Mountain	26,765	25,041	22,161	14,791	351	30
Pacific	62,968	54,864	50,136	33,653	358	49
Total	627,611	517,632	477,493	317,803	3,446	442

396 The results show a sample size of 3,446 for the GGA and a sample size of 477,493 for
397 the GA after 100 iterations.

398 4 An improved Bethel implementation

399 Our GGA was proposed and developed so that it would work with the rest of the functions
400 in *SamplingStrata*. Therefore the rest of the functions in the package remained unchanged.
401 This includes the *bethel.r* function which evaluates the fitness of chromosomes in every
402 iteration and is computationally expensive. For instance, for the PUMS dataset the
403 experiment took approximately 30 days for either GA or GGA with 100 iterations.

404 We searched for performance bottlenecks in *bethel.r* using the R *lineprof* package. Our
405 analysis of results suggested that the function within *bethel.r* called *chromy* appears to
406 take the bulk of computational time. A further examination reveals that *chromy* contains
407 a while loop with a default setting of 200 iterations. Furthermore *bethel.r* itself can be run

408 on each chromosome in any chromosome population on a dataset of any functional size
 409 (which we have the computation power to process) for any number of iterations. Bigger
 410 datasets will take longer to process. We expected that performance would be improved
 411 by converting the *bethel.r* algorithm into C++ then integrating that into R using the
 412 *Rcpp* package (Eddelbuettel, 2013).

Table 4.7: Performance comparison for the above datasets using the R and Rcpp versions of the Bethel-Chromy algorithm

Dataset	Records	Domains	Atomic Strata	Bethel μs	BethelRcpp μs	Speed-up Factor
iris	150	1	8	2684.77	143.13	18.76
swissmunicipalities	2,896	7	641	99,916	10,749.51	9.29
American Community Survey 2015	619,747	51	123,007	565,278,500	47,858,200	11.81
Kiva Loans Data	614,361	73	84,897	826,297,710	82,894,480	9.97
UN Commodity Trade Data 2011	351,057	171	350,895	139,749,810	87,555,870	1.6
US Census Data 2000	627,611	9	517,632	2,686,771	1,303,667	2.06

413 Table 4.7 shows the median time taken to run the Bethel algorithm one hundred times
 414 for the datasets we used to conduct our analysis. Our results confirm that the C++
 415 version of Bethel is faster than the R version. The speed up could make a practical
 416 difference in the number of iterations that can be run in *SamplingStrata* due to the
 417 processing times required for *bethel.r*. However, performance will vary according to the
 418 size and complexity of the problem. The speed up is achieved because C++ enables
 419 communication at a lower level with the computer than R. However, it is also due to the
 420 complexity of the analysis conducted in each for loop as well as the fact that larger data
 421 will restrict the available memory. It should also be noted that the C++ version of Bethel
 422 was compared with the R version as two stand alone functions. The performance of the
 423 C++ version of Bethel within the GGA is not compared with that of the R version in the
 424 GA. This would be part of a larger project to create a C++ version of the *SamplingStrata*
 425 package and integrating it into R.

426 5 Conclusion and further work

427 We created a GGA as an alternative to the existing *SamplingStrata* GA in R. We then
428 compared the two algorithms using a number of datasets. The GGA compares favourably
429 with the GA at finding the correct solution and meeting constraints on smaller datasets,
430 but significantly outperforms the GA on larger datasets where the number of iterations
431 was restricted. This is useful for datasets where the number of iterations has to be
432 constrained owing to computational burden. We have also reported faster processing
433 times by integrating the *bethel.r* function with C++ using the Rcpp package.

434 This work can be developed in several ways. Alternative evaluation techniques to speed
435 up the algorithm could be considered. Further research could also be undertaken into
436 other machine learning techniques for solving this problem.

437 The GGA could be applied to other problems which tackle more general sampling
438 designs with modifications required only for the algorithm evaluating the fitness of chro-
439 mosomes (i.e., the Bethel-Chromy algorithm). For example instead of searching for a
440 stratified simple random sample to meet precision constraints based on population totals
441 or means, the GGA could consider stratified probability proportional to size sampling
442 with an evaluation algorithm that uses more general estimators (e.g., regression or ratio
443 estimators) or more general parameters (e.g., a correlation coefficient).

444 The evaluation algorithm might also be modified to look at scenarios in which the
445 population variances are not known. In these cases, data from previous censuses, ad-
446 ministrative records, or proxy surveys can be used to estimate the population variance.
447 However, estimation of the population variance in a large number of atomic strata requires
448 more careful research.

449 Finally, the groupings of atomic strata by the GGA can be difficult to interpret. For
450 instance, an ordinal auxiliary variable taking values 1 to 4 may be unnaturally separated,
451 where the atomic strata corresponding to values 1 and 3 are grouped in one design
452 stratum and those with values 2 and 4 are grouped in another design stratum. It might

453 be interesting to explore less-than-optimal sample sizes for stratifications that are easier
454 to interpret. For instance, one may impose constraints on the admissible groupings. This
455 would require research into the formulation of appropriate admissibility constraints and
456 their effective implementation in the GGA.

457 **Acknowledgement**

458 We wish to acknowledge Steven Riesz of the Economic Statistical Methods Division of
459 the U.S. Census Bureau and Brian J. McElroy of the Economic Reimbursable Survey
460 Division of the U.S Census Bureau, both of whom answered questions which were of
461 assistance in choosing which US Census Bureau data to use. We would also like to
462 thank Giulio Barcaroli and Marco Ballin, the co-authors of (Ballin and Barcaroli, 2013),
463 for independently testing our GGA. Last but not least we are extremely grateful to the
464 editorial staff and reviewers of Survey Methodology for their constructive suggestions
465 in the review process for this journal submission, especially their suggestions for future
466 work.

467 **References**

- 468 Agustín-Blas, L.E., Salcedo-Sanz, S., Vidales, P., Urueta, G., Portilla-Figueras, J.A.
469 (2011). Near optimal citywide WiFi network deployment using a hybrid grouping ge-
470 netic algorithm. *Expert Systems with Applications* 38(8) 9543–9556.
- 471 Anderson, E. (1935). The irises of the gaspe peninsula. *Bulletin of the American Iris*
472 *society* 59, 2–5.
- 473 Ballin, M. and G. Barcaroli (2013). Joint determination of optimal stratification and
474 sample allocation using genetic algorithm. *Survey Methodology* 39(2), 369–393.

- 475 Barcaroli, G. (2014). SamplingStrata: An R package for the optimization of stratified
476 sampling. *Journal of Statistical Software* 61(4), 1–24.
- 477 Barcaroli, G. (2019 (accessed April 29, 2019)). Optimization of sampling strata with the
478 SamplingStrata package. [https://cran.r-project.org/web/packages/SamplingStrata/
479 vignettes/SamplingStrata.html](https://cran.r-project.org/web/packages/SamplingStrata/vignettes/SamplingStrata.html).
- 480 Bethel, J. W. (1985). *An optimum allocation algorithm for multivariate surveys*. American
481 Statistical Proceedings of the Survey Research Methods, Section, 209-212.
- 482 Bethel, J. (1989). Sample allocation in multivariate surveys. *Survey methodology* 15(1),
483 47-57.
- 484 E.C. Brown and M. Vroblefski. A grouping genetic algorithm for the microcell sectoriza-
485 tion problem. *Engineering Applications of Artificial Intelligence* 17(6), 589–598.
- 486 Bureau, US Census (2013 (accessed February 15, 2017)). *American Community Survey*
487 *Information Guide*. [http://www.census.gov/content/dam/Census/programs-surveys/
488 acs/about/ACS_Information_Guide.pdf](http://www.census.gov/content/dam/Census/programs-surveys/acs/about/ACS_Information_Guide.pdf).
- 489 Bureau, US Census (2016 (accessed February 15, 2017)). *2015 ACS PUMS DATA*
490 *DICTIONARY*. [http://www2.census.gov/programs-surveys/acs/tech_docs/pums/
491 data_dict/PUMSDict15.pdf](http://www2.census.gov/programs-surveys/acs/tech_docs/pums/data_dict/PUMSDict15.pdf).
- 492 Bureau, US Census (2016). *2015 ACS Public Use Microdata Sample (PUMS)*.
493 Washington, D.C. [https://factfinder.census.gov/faces/nav/jsf/pages/
494 searchresults.xhtml?refresh=t#](https://factfinder.census.gov/faces/nav/jsf/pages/searchresults.xhtml?refresh=t#).
- 495 Chromy, J. R. (1987). Design optimization with multiple objectives. *Proceedings of the*
496 *Survey Research Methods Section*.
- 497 De Lit, P., Falkenauer, E., Delchambre A. Grouping genetic algorithms: an efficient
498 method to solve the cell formation problem.

- 499 Eddelbuettel, E. (2013). Seamless R and C++ Integration with Rcpp *ISBN 978-1-4614-*
500 *6867-7* 10.1007/978-1-4614-6868-4
- 501 Falkenauer, E. (1998). *Genetic algorithms and grouping problems*. John Wiley & Sons,
502 Inc.
- 503 Fisher, R. A. (1936). The use of multiple measurements in taxonomic problems. *Annals*
504 *of Eugenics* 7(2), 179–188.
- 505 Galinier, P. , Hao, J. K.. Hybrid Evolutionary Algorithms for Graph Coloring. *Journal*
506 *of Combinatorial Optimization* 3(4):379–397, 1999.
- 507 Hartigan, J. A. and M. A. Wong (1979). Hybrid Evolutionary Algorithms for Graph Col-
508 oring. Algorithm as 136: A k-means clustering algorithm *Journal of the Royal Statistical*
509 *Society. Series C (Applied Statistics)* 28(1), 100–108.
- 510 Hung, C. , Sumichrast, R. T., Brown, E.C. CPGEA: a grouping genetic algorithm for
511 material cutting plan generation. *Computers & Industrial Engineering* 44(4), 651-672.
- 512 James, T., E. Brown, and C. T. Ragsdale (2010). Grouping genetic algorithm for the
513 blockmodel problem. *IEEE Transactions on Evolutionary Computation* 14(1), 103–
514 111.
- 515 Neyman, J. (1934). On the two different aspects of the representative method: the
516 method of stratified sampling and the method of purposive selection. *Journal of the*
517 *Royal Statistical Society* 97(4), 558-625.
- 518 Pelikan, M. and D. E. Goldberg (2000), Genetic Algorithms, Clustering, and the Breaking
519 of Symmetry. *Proceedings of the Sixth International Conference on Parallel Problem*
520 *Solving from Nature*,.
- 521 Prügel-Bennett, A. (2004). Symmetry Breaking in Population-Based Optimization. *IEEE*
522 *Transactions on Evolutionary Computation* 8(1), 63-79.

- 523 Stokes, L. and J. Plummer (2004). Using spreadsheet solvers in sample design. *Compu-*
524 *tational statistics & data analysis* 44 (3), 527–546.
- 525 R Core Team (2015). *R A Language and Environment for Statistical Computing*. Vienna,
526 Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- 527 Ruggles, S., K. Genadek, R. Goeken, J. Grover, and M. Sobek (2017). Integrated public
528 use microdata series: Version 7.0 [dataset]. minneapolis: University of minnesota, 2017.
- 529 Willighagen, E. (2005). Genalg: R based genetic algorithm. *R package version 1*.