

1

Supplementary Information:

2

Multi-objective optimized genomic breeding

3

strategies for sustainable food improvement

4

Deniz Akdemir^{a,**}, William Beavis^b, Roberto Fritsche-Neto^c,

5

Asheesh K.Singh^b, and Julio Isidro-Sánchez^d

6

^aCornell Statistical Consulting Unit, Cornell University,

7

Ithaca, NY, USA

8

^bDepartment of Agronomy, Iowa State University, Ames, Iowa,

9

USA.

10

^cDepartment of Genetics, “Luiz de Queiroz” Agriculture

11

College, University of Sao Paulo, Brazil

12

^bDepartment of Agronomy, Iowa State University, Ames, Iowa,

13

USA

14 ^cAgriculture & Food Science, Animal and Crop section,

15 University College Dublin, Dublin, Ireland

16 ^{**}Correspondance: E-mail: akdemir@cornell.edu

17

August 9, 2018

18 **Methods**

19 **Model for estimating breeding values**

20 More specifically, the model that is used is given by

$$Y = XB + ZG + E, \tag{1}$$

21 where Y is the $n \times d$ response variable, X is the $n \times q$ design matrix of $q \times d$ the

22 fixed effects B , Z is a $n \times q$ design matrix of the $q \times d$ random effects G , and

23 E is the $n \times d$ matrix of residual effects. The random effects and the residual

24 are independently distributed, and have matrix variate distributions ($G \sim$

25 $N_{q \times d}(0_{q \times d}, K, \Sigma)$ and $E \sim N_{n \times d}(0_{n \times d}, R, \Sigma_E)$) where K is a $q \times q$ relationship

26 matrix, Σ is a $d \times d$ covariance matrix, R is a $n \times n$ covariance matrix, Σ_E is a

27 $d \times d$ covariance matrix. An early reference to this multi-trait model appears
28 in [4].

29 **Multi-objective optimization concepts**

30 A vector $u = (u_1, \dots, u_k)$ is said to dominate another vector $v = (v_1, \dots, v_k)$
31 (written as $u \succeq v$) if and only if u is partially less than v , i.e., $\forall i \in 1, \dots, k$,
32 $u_i \leq v_i \wedge \exists i \in 1, \dots, k$ $u_i < v_i$. Pareto optimal solutions are those which,
33 when evaluated, produce vectors whose performance f_i cannot be improved
34 without adversely affecting another f_j , $i \neq j$. In a minimization problem, a
35 solution \mathbf{x} is said to be Pareto optimal if and only if there is no \mathbf{x}' for which
36 $F(\mathbf{x}')$ dominates $F(\mathbf{x})$, i.e., there exists no feasible vector \mathbf{x}' which would
37 decrease some criterion without causing a simultaneous increase in at least
38 one other criterion.

39 For a multi-objective problem, $F(\mathbf{x})$, the Pareto Optimal Set, P^* , is de-
40 fined as: $P^* := \{\mathbf{x} : \neg \exists \mathbf{x}' F(\mathbf{x}') \succeq F(\mathbf{x})\}$.

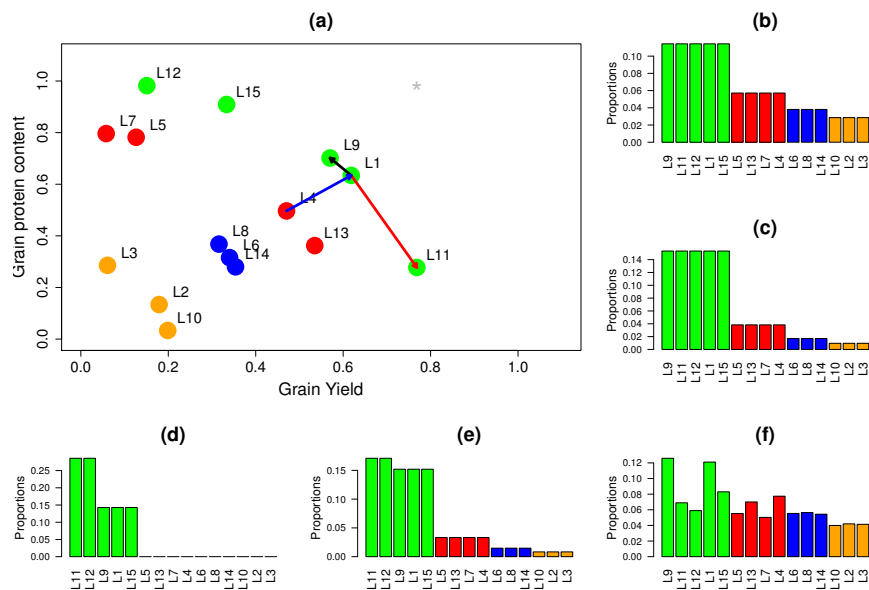
41 The Pareto front PF^* is defined as: $PF^* := \{u = F(\mathbf{x}) | \mathbf{x} \in P^*\}$.

42 Non-dominance ordering and assignment of parental contribution pro-
43 portions based on the above ideas are demonstrated with an hypothetical
44 example in Figure 1. (a-f). In this figure, genotype L1 is nondominated

45 (so as L12, L15, L9, L11) with respect to GY and GPC since there are no
46 other genotypes in this set when replaced with L1 that would not involve
47 decrease in at least one of BVs. For example, replacing L1 with L11 gives an
48 increase in GY value but a decrease in the GPC (red arrow). Replacing L1
49 with L9 gives an increase in GPC but a decrease in the GY (black arrow).
50 Genotype L4 is dominated by L1, because by replacing it with L1 we can
51 increase both BVs. The genotypes in dominance Level 1 are the set of non-
52 dominated individuals. The genotypes in dominance Level 2 (red points) are
53 obtained as the non-dominated genotypes in the smaller subset obtained by
54 removing the genotypes in Level 1. This process is continued until all the
55 genotypes are assigned to their dominance levels. The gray star refers to an
56 example of an ideal genotype with respect to the set of genotypes. Weights
57 can be assigned inversely proportional to the distance from ideal solution.
58 (b) Parental proportions are inversely proportional to dominance levels. (c)
59 Parental proportions are inversely proportional to squared dominance levels.
60 (d) Parental proportions are obtained with respect to the non-dominance
61 counts over combinations of traits. L12 is non dominated with respect to
62 Grain protein and also when we consider the two traits at the same time, so
63 gets a weight proportional to 2. L1 is only non-dominated with respect to

64 two traits at a time, so gets a weight proportional to 1. (e) Is obtained by
 65 taking the mean of the proportions in (c) and (d). (e) Parental proportions
 66 inversely proportional to the distances obtained using formula 2 using the
 67 ideal solution, $p = 2$ and equal weights for traits. Bar-plots are colored with
 respect to dominance levels.

Figure 1: (a) This is a scatterplot for breeding values (BVs) of grain yield (GY) and grain protein content (GPC) for 15 genotypes.



68

69 Multi-objective optimization techniques

70 In general, it is not possible to find an analytic expression of the line or
71 surface that contains the Pareto optimal solutions. Several techniques are
72 used to find representative points on the Pareto frontier:

73 Scalarization (Linear combination)

74 If a solution \mathbf{x} to the general multi-objective problem is non inferior,
75 then there exist $w_l = 0, l = 1, 2, \dots, k$ (w_r is strictly positive for some
76 $r = 1, 2, \dots, k$), and $\lambda_i = 0, i = 1, 2, \dots, m$, such that:

$$\sum_{l=1}^k w_l \delta f_l(\mathbf{x}) - \sum_{i=1}^m \lambda_i \delta g_i(\mathbf{x}) = 0.$$

77 This condition is named Kuhn-Tucker Condition (KTC) and is necessary
78 for a non-inferior solution. When all of the $f_l(\mathbf{x})$ are concave and \mathbf{x} be-
79 longs to a convex set, they are sufficient as well. Since KTC is sufficient
80 for non-inferiority, non-inferior solutions might be found by solving a scalar
81 optimization problem in which the objective function is a weighted sum of
82 the components of the original vector-valued function. That is, the solution
83 to the problem: $\sum_{i=1}^k w_i f(\mathbf{x})$, where $w_i \geq 0$ for all i and strictly positive for
84 at least one objective, is usually non inferior. Then on inferior set and the
85 set of non-inferior solutions can be generated by varying the weights w_i in

86 the objective function.

87 The reduction of the problem to a single-objective function means to
88 make all alternatives comparable with a preference framework that becomes
89 a total order. Hence w_i values choice is very important to achieve the final
90 decision and, for this reason, value choice is made by the decision maker.
91 However, the decision maker, in order to choose the coefficients, must have
92 a clear perception of how this choice effects all the functions with respect to
93 each other.

94 The main advantages of this method are its simplicity (in implementation
95 and use) and its efficiency (computationally speaking). Its main disadvan-
96 tage is the difficulty to determine the appropriate weight coefficients to be
97 used when enough information about the problem is not available (this is an
98 important concern, particularly in real-world applications). Also, a proper
99 scaling of the objectives requires a considerable amount of extra knowledge
100 about the problem. To obtain this information could be a very expensive
101 process. A more serious drawback of this approach, is that it cannot be used
102 to generate certain portions of the Pareto front when the conditions of KTC
103 are not satisfied, regardless of the weights combination used. Nevertheless,
104 aggregating functions could be very useful to get a preliminary sketch of

105 the Pareto front of a certain problem or to provide prior information to be
106 exploited by another approach.

107 Other scalarization methods include L_p -norm, Chebyshev and the single-
108 objective product formulation. For each of these scalarizations, a charac-
109 terization of the Pareto set can be obtained by varying the scalarization
110 parameters and solving many single-objective optimization problems. Ide-
111 ally, the points returned by the scalarized problems should be sufficiently
112 spread out in the efficient frontier.

113 **The ϵ -constraint method**

114 Besides the scalarization by linear combination approach, the ϵ -constraint
115 method is probably the best known technique to solve multi-objective opti-
116 mization problems. There is no aggregation of criteria, instead only one of
117 the original objectives is minimized while the others are transformed to con-
118 straints. The idea was introduced by [2]. Through this approach among p
119 objective function only one is kept as such, the other $p - 1$ are transformed in
120 constraints fixing threshold values ϵ_k (with $k = 1, \dots, p, k \neq j$) over them (if
121 functions must be minimized). Therefore, the problem: $\min F(\mathbf{x})$ is substi-
122 tuted by the ϵ -constraint problem: $\min f_j(\mathbf{x}) \quad f_k(\mathbf{x}) \leq \epsilon_k, k = 1, \dots, p, k \neq j$.
123 The main disadvantage of this approach is its (potentially high) computa-

124 tional cost, also due to the preliminary individuation of ϵ_i values.

125 **Finding non-dominated solutions**

126 Several algorithms exist for finding the non-dominated set in from a larger
127 set. Popular Kung algorithm ([5]) involves first sorting the population in de-
128 scending order in accordance to first objective function. Afterwards, the pop-
129 ulation is recursively partitioned as top (T) and bottom (B) sub-populations.
130 As top half (T) is better in objective in comparison to bottom half (B) in
131 first objective, so we check the bottom half for domination with top half.
132 The solution of B which are not dominated by solutions of T are merged
133 with members of T to form merged population M. Another algorithm is the
134 Jun Du Algorithm ([3]).

135 **Selecting a "good" solution on the frontier surface**

136 **Ideal solution concept and global criterion:**

137 Let \mathbf{x}_i^* be a vector which optimizes the i th objective function $f_i(\mathbf{x})$ for $i =$
138 $1, 2, \dots, k$. Then the vector $[f(\mathbf{x}_1^*), f(\mathbf{x}_2^*), \dots, f(\mathbf{x}_k^*)]'$ is ideal for an multi-
139 objective problem and is consequently called the ideal vector.

140 The global criterion method aims to minimize a function (global criterion)
141 which is a measure of how close the DM can get to the ideal vector. A measure

142 of closeness to the ideal solution is a family of L_p -metrics defined as follows:

$$L_p(F) = \left[\sum_{l=1}^{k+1} \left| \frac{f(\mathbf{x}_1^*) - f(\mathbf{x})}{f(\mathbf{x}_1^*)} \right|^p \right]^{1/p}.$$

143 If the functions values are normalized to the range $[0, 1]$, then the above

144 formula becomes

$$L_p(F) = \left[\sum_{l=1}^{k+1} |1 - f(\mathbf{x})|^p \right]^{1/p}.$$

145 Finally, if weights are attached to the functions $f_1(\mathbf{x}), f_1(\mathbf{x}), \dots, f_{k+1}(\mathbf{x})$ a

146 weighted version can be written as

$$L_p^w(F) = \left[\sum_{l=1}^{k+1} w_l |1 - f(\mathbf{x})|^p \right]^{1/p}. \quad (2)$$

147 In the remaining of this manuscript, we have used $p = 2$ which coincides with

148 the Euclidean distance.

149 **Ideal Solution for GEBVs**

150 A simple estimator, say $\max f_i$, for the ideal solution $f(\mathbf{x}_i^*)$, for trait i in

151 a certain breeding population is the maximum observed value for that trait.

152 It is also possible to estimate this quantity by calculating the maximal esti-

153 mated genomic value using the marker effects estimates. We use the former

154 approach in the remaining of this article. When calculating the distance from

155 the ideal solution, the objective function values were scaled to the range $[0, 1]$

156 using the transformation

$$f^*(\mathbf{x}) = \frac{f_i(\mathbf{x}) - \min f_i}{\max f_i - \min f_i},$$

157 for trait i ; and $\min f_i$ is the estimate for the worst value of the trait i , it is

158 calculated in the same fashion as $\max f_i$.

159 **Other multi-trait breeding approaches**

160 Supplementary Figure 5 displays the parental contribution proportions ob-
161 tained by non-dominance counts for 100 lines with the highest proportions
162 for the improving yield and protein in the four environments.

163 Supplementary Figures 8 and 9 display the individuals that would be
164 identified by the classical multi-trait breeding schemes culling, tandem and
165 index selection for the wheat and barley datasets. These can be contrasted
166 with the Figures 2 a and b and also with the Figure 4 and Supplementary
167 Figure 4. Note that among the classical methods index selection will give
168 the closest results to the multi-objective optimized breeding methods.

169 **Multi-objective training population design**

170 Two examples with different sets of selection of training populations related
171 optimality criteria where we display Pareto fronts for training populations of

Figure 2: This figure represents the non-dominance ordering of the individuals in the barley dataset for three traits. Barley data: Dominance ordering based on three traits, 13 levels of dominance. The GEBVs for height, grain yield and protein from barley data are plotted with the dominance ordering of these individuals indicated by the lines of different color.

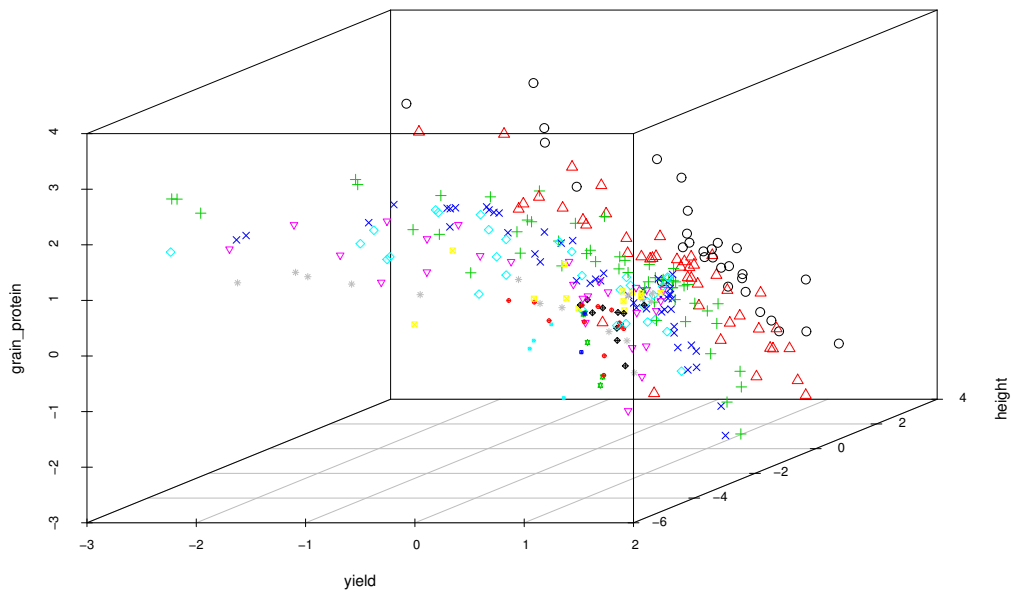
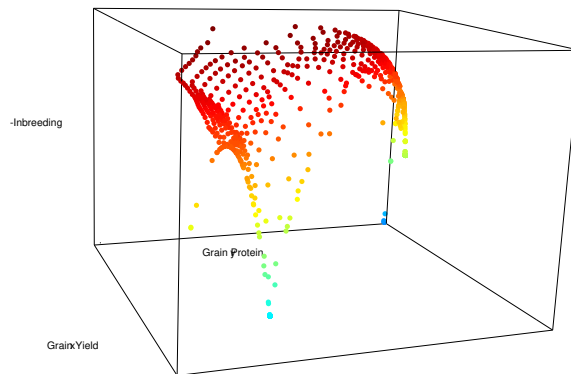


Figure 3: Pareto optimal solutions for parental contributions (Barley data) obtained by solving the optimization problem giving in Equation (1) for improving grain yield (GY) and grain protein content (GPC) while controlling coancestry, i.e, we assume we want to maximize GY, GPC and the negative of inbreeding. The redness of the points indicates closeness to ideal solutions as calculated by the formula in Supplementary Equation (2).



172 size 100 (selected from the remaining genotypes in the wheat dataset after
173 reserving a random sample of genotypes as a target population). Some exam-
174 ples of optimality criteria include determinant optimality criterion (minimizes
175 the determinant of the covariance matrix of the model coefficient estimates

Figure 4: Three 'good' solutions on the barley frontier curve obtained from Supplementary Figure 3. Red points indicated the individuals that have non-zero parental proportions. The sizes of the points are proportional to the magnitude of the parental contributions. The figures on the right side, represent the same information but on the first two principal components of the genotyping marker space.)

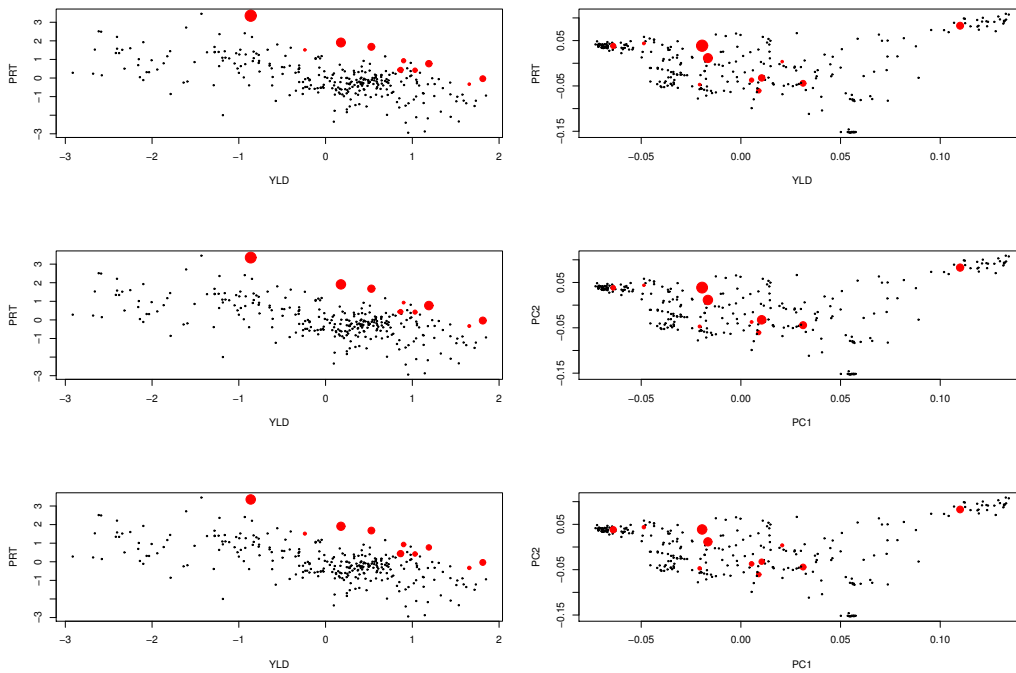


Figure 5: Barley data: parental contribution proportions obtained by non-dominance counts for 100 lines with the highest proportions for the improving yield in four environments.

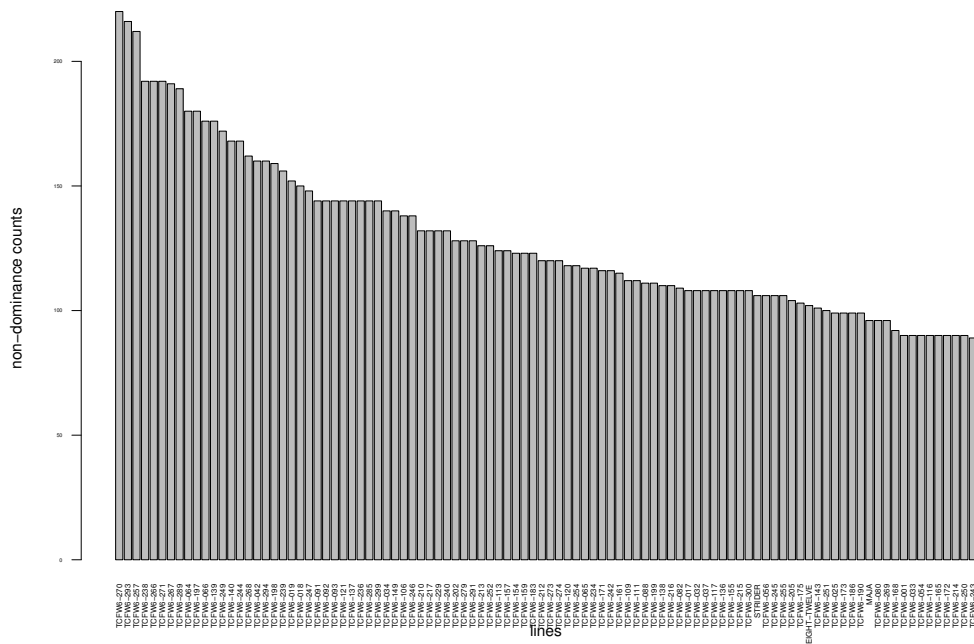
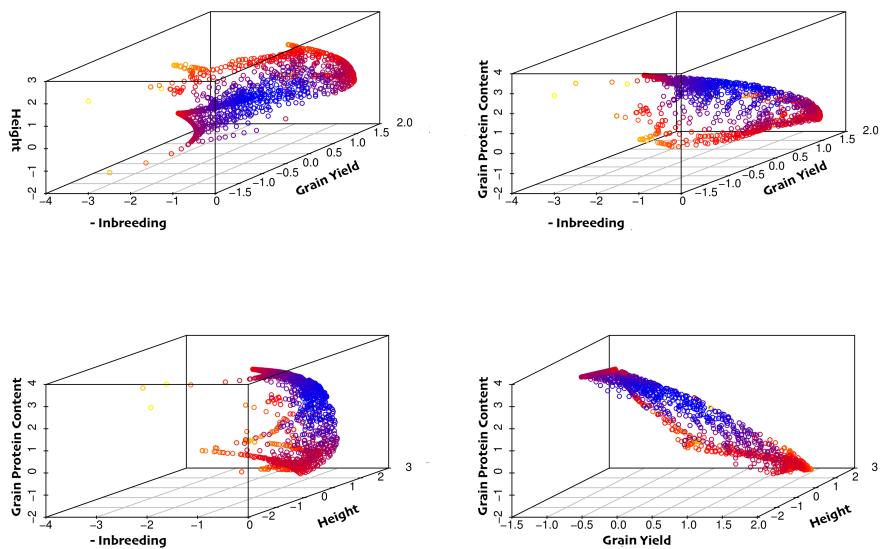
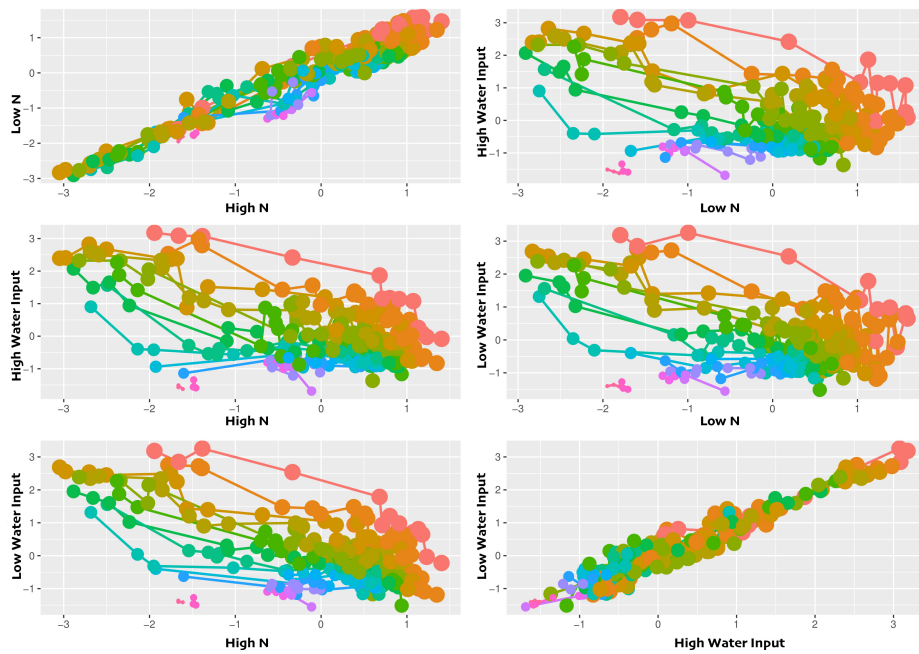


Figure 6: Barley data: The dimensions correspond to negative of inbreeding, and average gains based on GEBVs for grain yield (GY), height and grain protein content (GPC). Each point on the 3 dimensional scatterplots correspond to a Pareto optimal solution for parental contributions. Blueness of the points measure the closeness to the ideal solution as calculated by Equation (2).



176 for a principal components regression model based on the markers; genomic
 177 distance based criteria such as mean or the minimum distances among geno-
 178 types in the training set, mean or maximum distance to the target set of
 179 genotypes. In the examples below, determinant optimality criterion was

Figure 7: Barley data: Dominance ordering based on yield in 4 environments (dry-irrigated \times high-low nitrogen). The environment specific GEBVs for grain yield for barley data are plotted with the dominance ordering of these individuals indicated by the lines of different color.



180 calculated using the first 50 principal components of the marker matrix as
 181 suggested in the R package STPGA [1] and the distance based criteria were
 182 calculated using the Euclidean distance matrix calculated using the marker
 183 matrix. The multi-objective optimization problem is setup such that we are
 184 seeking solutions (lists of training sets of size 100) to minimize these different
 185 criteria.

Figure 8: Wheat data: Other multi-trait breeding approaches

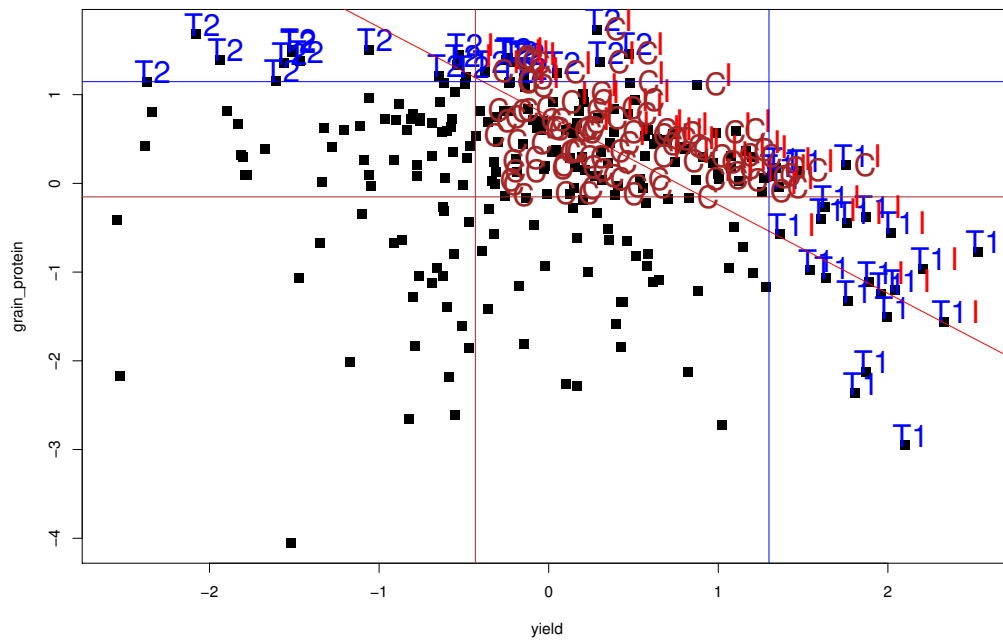


Figure 9: Barley data: Other multi-trait breeding approaches

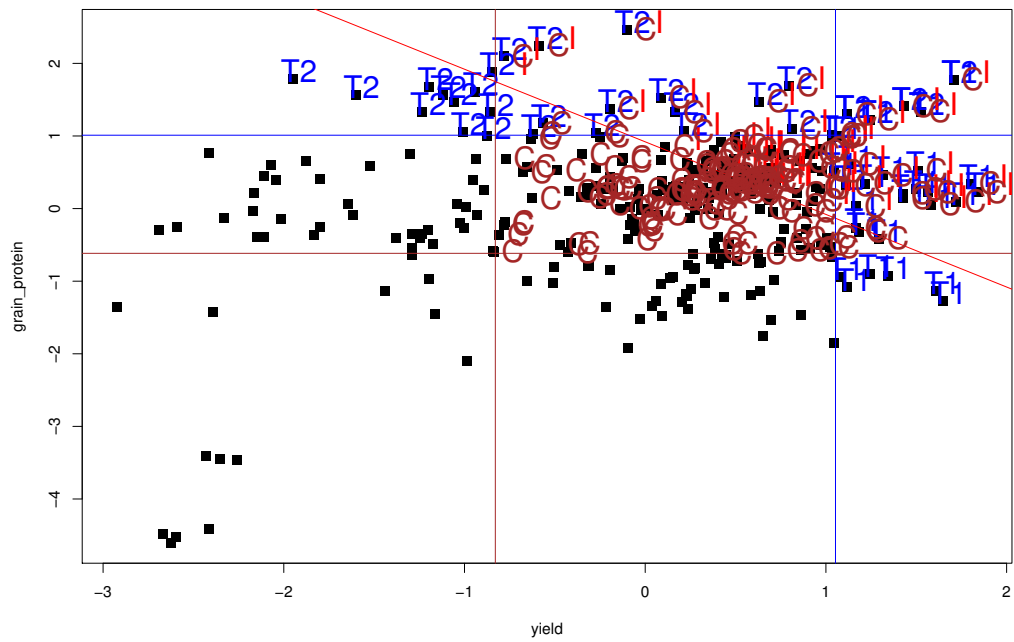


Figure 10: Barley data: Four dimensional Pareto front represented in two dimensional plots. The dimensions correspond to negative of inbreeding, and average gains based on genomically estimated breeding values for yield, height and protein content. Each point on the scatterplots correspond to a Pareto optimal solution for parental contributions. Blueness of the points measure the closeness to the ideal solution.

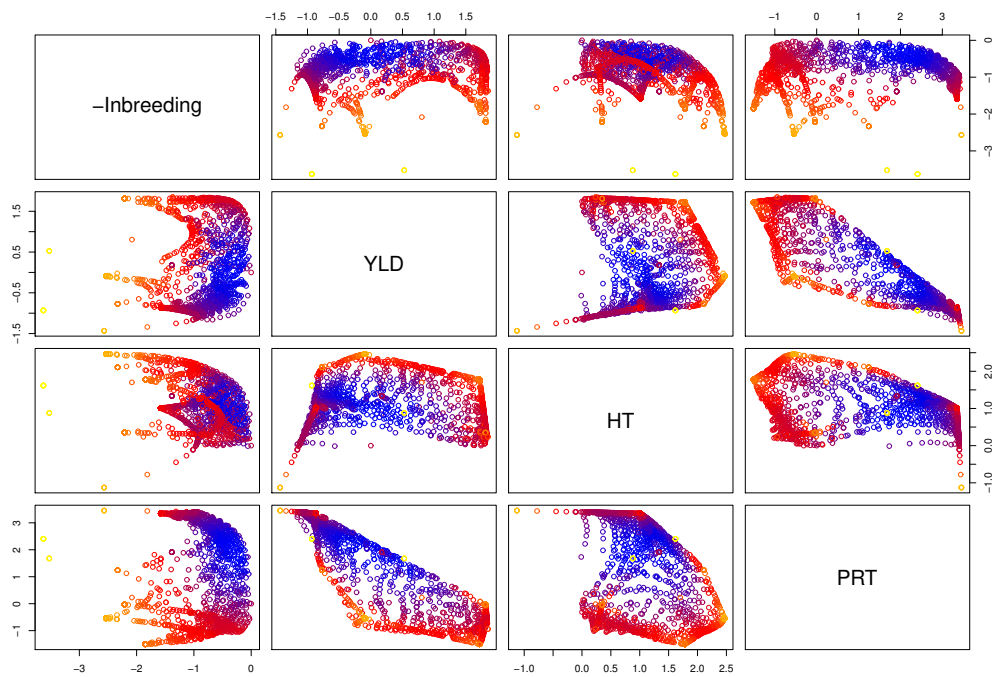


Figure 11: Barley data: SOM plot for barley data for GEBVs over three traits, yield, protein content and height. In this example, the figure with the title "Mapping" displays the mapping of the genotypes into clusters that are obtained using SOM analysis. These clusters are displayed in a two dimensional grid and the topology preserving mapping property means that closely located clusters contain genotypes with similar properties in terms of the three traits. The three dimensional dendrogram shows the closeness of these clusters in the SOM space. The change in the average values of these three traits can be observed from the "SOM Plot". In addition the surface plots display the change in the individual traits along the mapping directions of the "Mapping" plot.

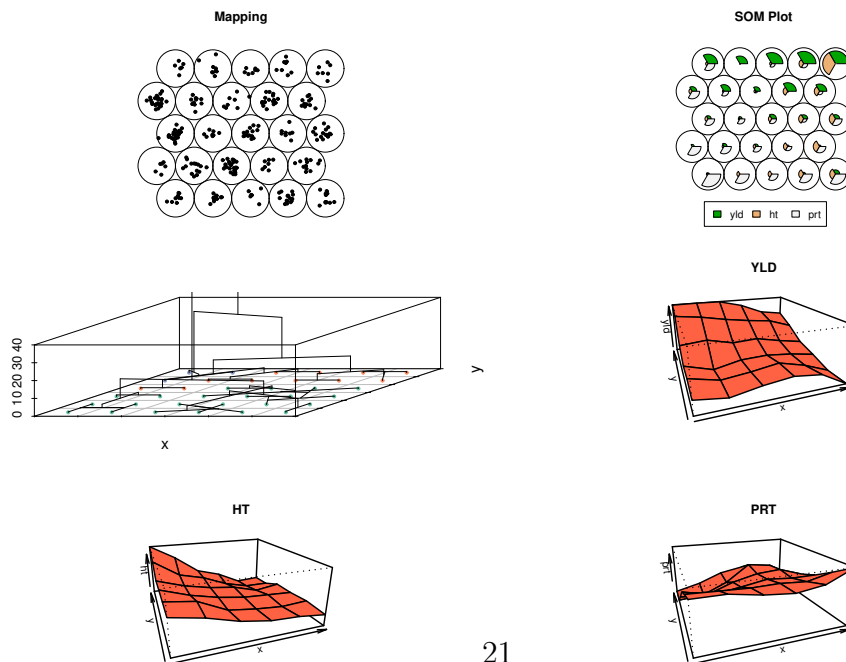


Figure 12: **Simulated data: Gain, Usefulness, Coancestry:** Genomic mating approach is also a multiobjective optimized breeding approach. The frontier surface in the figure represents the tradeoff between the measures of gain, usefulness, coancestry for pareto optimal mating plans.

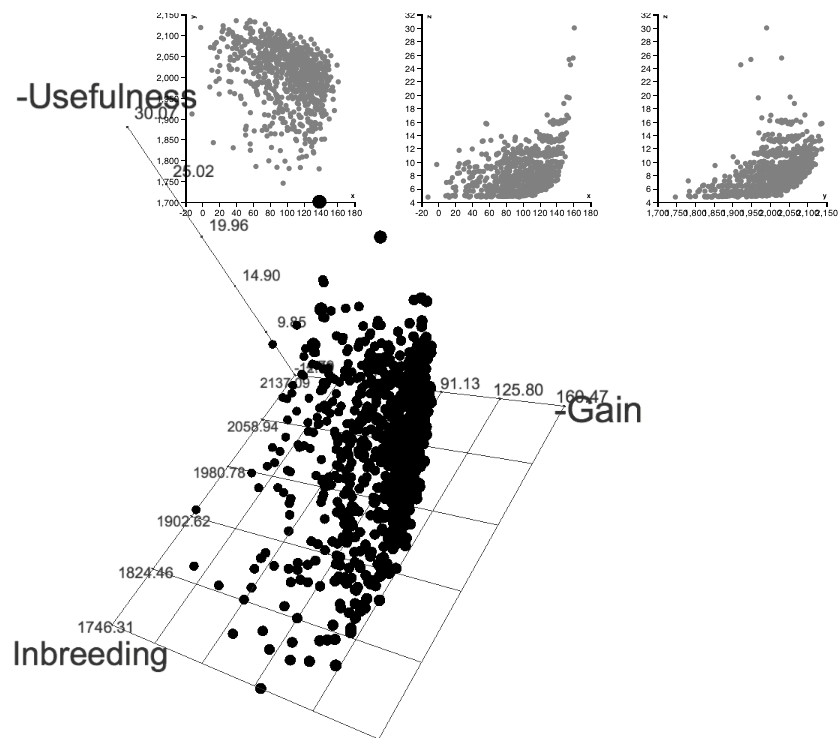
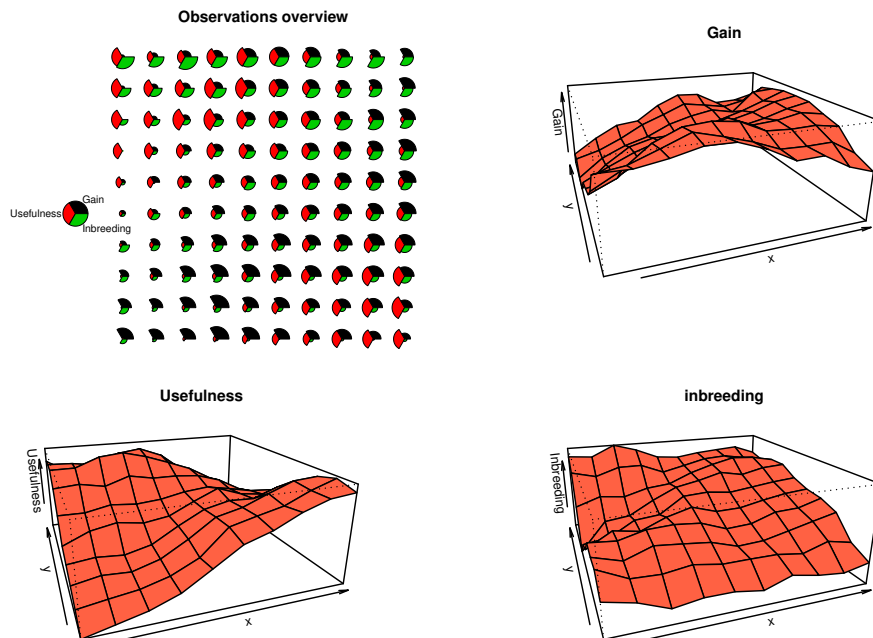


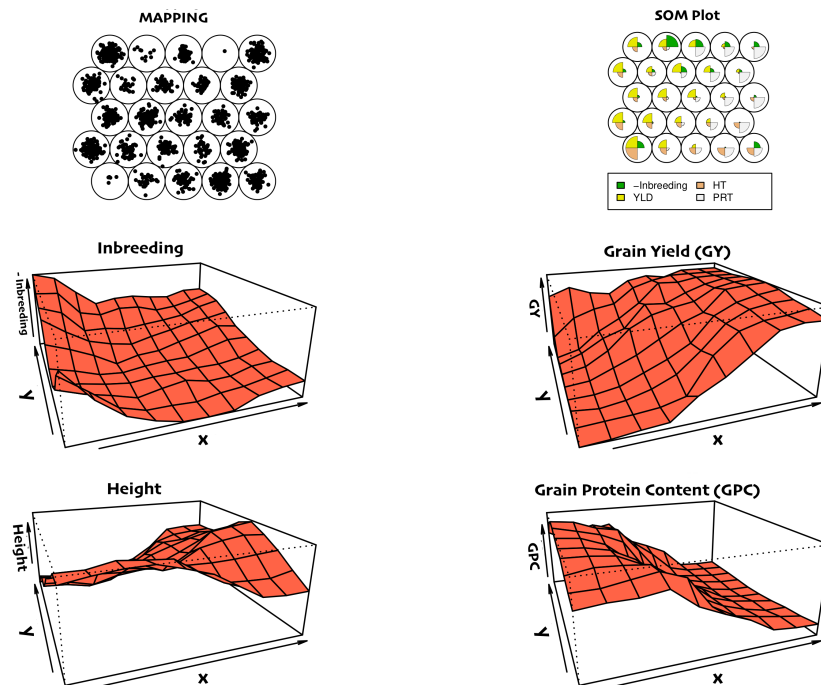
Figure 13: **Simulated data: Gain, Usefulness, Coancestry:** The top left figure "Observations overview" is a representation of all the solutions on the Pareto frontier in Supplementary Figure 12 using the mapping of these solutions into clusters that are obtained using SOM analysis. These clusters are displayed in a two dimensional grid and the topology preserving mapping property means that closely located clusters contain solutions with similar properties in terms of gain, usefulness, and coancestry. The remaining graphs are showing how each measure behaves individually on the two SOM features that correspond to the dimensions of the "Observations overview".



186 **References**

- 187 [1] D. Akdemir. *STPGA: Selection of Training Populations by Genetic Al-*
188 *gorithm*, 2018. R package version 5.0.
- 189 [2] V. Chankong and Y. Y. Haimes. *Multiobjective decision making: theory*
190 *and methodology*. Courier Dover Publications, 2008.
- 191 [3] J. Du and Z. Cai. A sorting based algorithm for finding a non-dominated
192 set in multi-objective optimization. In *Natural Computation, 2007. ICNC*
193 *2007. Third International Conference on*, volume 4, pages 436–440. IEEE,
194 2007.
- 195 [4] C. Henderson and R. Quaas. Multiple trait evaluation using relatives’
196 records. *Journal of Animal Science*, 43(6):1188–1197, 1976.
- 197 [5] H.-T. Kung, F. Luccio, and F. P. Preparata. On finding the maxima of
198 a set of vectors. *Journal of the ACM (JACM)*, 22(4):469–476, 1975.

Figure 14: Self-Organizing Maps plot for barley data for parental proportions over three traits, grain yield, grain protein content and height and negative of inbreeding. This is another representation of the fourth dimensional Pareto surface in Supplementary Figure 6. The dimension x and y are found by the SOM algorithm to preserve the distances measured by the four dimensions of the frontier curve. The "Mapping" shows solutions on the frontier curve clustered with respect to x and y units (SOM dimensions), the "SOM Plot" gives the average values of the four dimensions in each cluster, the three dimensional surfaces show individual response surfaces for the four dimensions with respect to SOM dimensions.



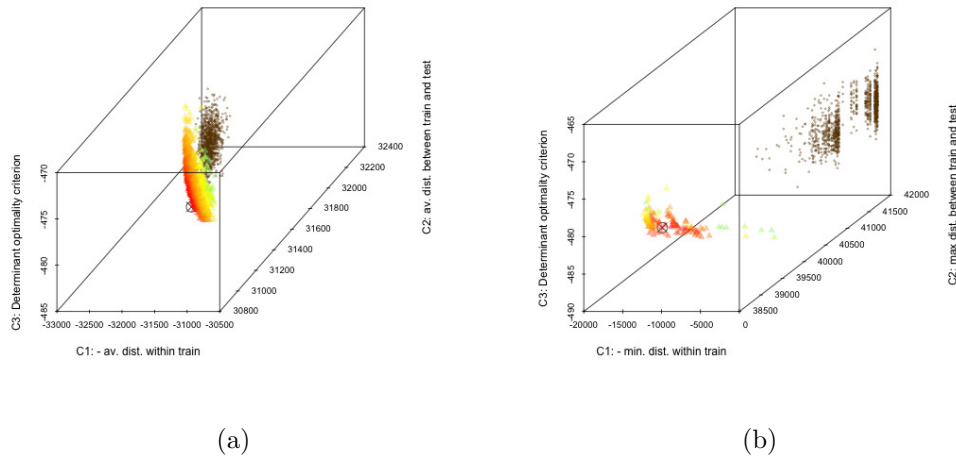


Figure 15: MOO selection of training populations for Wheat data: A random subset of 50 genotypes were selected as target population. For training sets of size 100 selected from the remainder genotypes after removing the target set. The brown circles in the graphs represent criteria values for the 1000 initial solutions where the colored triangles represent criteria values for solutions on the Pareto front. The solution which has the shortest distance to the ideal solution (the solution formed by combining the best solutions for each of the criteria taken individually) is marked by a black \otimes symbol. (a) Optimization Problem 1: Look for solutions that minimize DOPT criterion, negative mean genetic distance in the training population and mean genetic distance of training to target. (b) Optimization Problem 2: Look for solutions that minimize DOPT criterion, negative of minimum genetic distance between pairs of individuals in the training population and maximum genetic distance of individuals in training set to individuals in the target set.