

Title: Optimal Rate of Correct Assignment with backward elimination
locus selection

Version:
1.0

Authors: J. Jasper and W. Templin

Date: December 14, 2010

Introduction

As part of the locus selection process proposed for chum salmon in WASSIP, we propose using f_{ORCA} (Rosenberg et al. 2003; Rosenberg 2005) with backward elimination as one of the marker selection methods for choosing SNPs for the chum salmon baseline (Tech Doc 8). Results from this analysis are proposed to provide 30% of the locus-selection weight, the most of any analysis. The information measure, f_{ORCA} , returns the Optimal Rate of Correct Assignment (ORCA) for a particular locus set with respect to a specific baseline. At each iteration of the routine, a randomly drawn individual is assigned to a population for which its genotypic probability is a maximum. We propose adapting f_{ORCA} to allow us to determine the best set of loci to provide separation among reporting groups taking advantage of potential synergy among loci. To do this we propose implementing a backward elimination algorithm similar to that described in BELS (Bromaghin 2008). However, we opted not to use the program BELS because it is too time-consuming. Even though the Gene Conservation Laboratory does proportional allocation (as does BELS) rather than individual assignment (as does f_{ORCA}), we feel that f_{ORCA} with backward elimination has merit under a Bayesian mixed stock analysis routine because it attempts to select a suite of markers that optimizes the genotypic probabilities of potential mixture individuals, and BAYES (Pella and Masuda 2001) uses these probabilities to stochastically assign the mixture individuals each iteration.

Current f_{ORCA} Algorithm

While a closed form solution of f_{ORCA} is available (Rosenberg et al. 2003), it becomes impractical for large locus sets. Therefore, Rosenberg (2005) provided an iterative algorithm for estimating f_{ORCA} . This algorithm can be explained as follows.

1. Uniformly draw a population at random from the baseline.
2. Randomly generate a multi-locus genotype based on the allele frequencies of the population chosen in the first step.

¹ This document serves as a record of communication between the Alaska Department of Fish and Game Commercial Fisheries Division and the Western Alaska Salmon Stock Identification Program Technical Committee. As such, these documents serve diverse ad hoc information purposes and may contain basic, uninterpreted data. The contents of this document have not been subjected to review and should not be cited or distributed without the permission of the authors or the Commercial Fisheries Division.

- 32 3. Assign that genotype to the population for which its genotypic probability is a
33 maximum.
- 34 4. Repeat Steps 1-3 10,000 times.
- 35 5. After repeating this process multiple times, f_{ORCA} is calculated as the proportion of
36 times that the assignment in Step 3 is the same population drawn in Step 1.

37 While f_{ORCA} is typically used to evaluate how well a marker set can assign individuals
38 back to the correct population, it could also be adapted for evaluating how well a marker
39 set can be used to assign individuals back to the correct region. With this application the
40 algorithm would be as follows.

- 41 1. Uniformly draw a population at random from the baseline.
- 42 2. Determine the region to which the population belongs.
- 43 3. Randomly generate a multi-locus genotype based on the allele frequencies of the
44 population chosen in the first step.
- 45 4. Assign that genotype to the population for which its genotypic probability is a
46 maximum.
- 47 5. Determine the region to which the assignment population belongs.
- 48 6. Repeat Steps 1-5 10,000 times.
- 49 7. After repeating this process multiple times, f_{ORCA} is calculated as the proportion of
50 times that the assignment in Step 5 is the same region drawn in Step 2.

51

52 **Backward Elimination Locus Selection Algorithm**

53 Rosenberg's f_{ORCA} algorithm provides a means of evaluating the performance of a locus
54 set, but it does not provide us with an algorithm for selecting sets of markers to evaluate.
55 Rosenberg (2005) does provide four such algorithms and discusses the advantages and
56 limitations of each: 1) Exhaustive evaluation, 2) Univariate accumulation, 3) Greedy
57 accumulation, and 4) Maxmin accumulation.

58 One locus selection algorithm that Rosenberg failed to discuss is the method used in the
59 Backward Elimination Locus Selection (BELS) algorithm laid-out by Bromaghin (2008).
60 This algorithm has the advantages of being both simple to implement and it exploits
61 synergies among loci. However, Bromaghin (2008) does not use f_{ORCA} to evaluate marker
62 sets; rather he uses actual maximum likelihood mixed stock analysis and bootstrap
63 simulations to evaluate performance in the software BELS. While we agree that this is a
64 relevant measure, unlike f_{ORCA} , it suffers from being prohibitively slow and may be biased
65 in some circumstances (Anderson 2008).

66 We suggest that marker selection applications with large numbers of populations and loci
67 should employ the BELS algorithm for selecting marker panels to evaluate, but use the
68 f_{ORCA} function to do the evaluation. For the purposes of WASSIP, we will use the correct
69 assignment to region algorithm described above.

70 This would be accomplished by the following:

- 71 1. Start with entire set of L potential markers.
72 2. Create L sub-sets of L-1 markers by removing each marker, in turn, from full the
73 set.
74 3. Evaluate f_{ORCA} on all L sub-sets using correct assignment to region.
75 4. Identify sub-set with maximum f_{ORCA} .
76 5. Record which locus was removed.
77 6. Return to Step 1 using the sub-set identified in Step 4 as the new full set of L-1
78 loci.

79 This process is continued until no markers remain. The loci can be ranked according to
80 the order in which they were removed or scored according to their f_{ORCA} value.

81 This algorithm has been implemented in R for use with the chum salmon SNP selection
82 process described in Technical Document 8, “Chum salmon SNP selection process
83 outline.”

84 The limitations of f_{ORCA} are: 1) it (likely) suffers from providing an optimistic rate of
85 correct assignment, and; 2) spurious differences in allele frequencies can lead to falsely
86 identifying some loci as influential. An extension of f_{ORCA} that may alleviate its
87 limitations would be to implement a “leave-one-out” approach by which we randomly
88 draw an individual from the ascertainment baseline, recalculate the allele frequencies
89 without that individual, then assign the individual based on the recalculated allele
90 frequencies. While more difficult to implement, this version may be a more viable
91 solution. We are currently working on programming this extension.

92

93 Citations

- 94 Anderson E.C., R.S. Waples, S.T. Kalinowski. 2008. An improved method for estimating
95 the accuracy of genetic stock identification. *Canadian Journal of Fisheries and*
96 *Aquatic Sciences* 65:1475-1486.
- 97 Bromaghin, JF. 2008. BELS: backward elimination locus selection for studies of mixture
98 composition or individual assignment. *Molecular Ecology Resources* 8: 568-571
- 99 Rosenberg, NA, LM Li, R Ward, & JK Pritchard. 2003. Informativeness of Genetic
100 Markers for Inference of Ancestry. *American Journal of Human Genetics* 73
101 (1421):1402-1422
- 102 Rosenberg, NA. 2005. Algorithms for Selecting Informative Marker Panels for
103 Population Assignment. *Journal of Computational Biology* 12 (9):1183–1201

104

105

Technical Committee review and comments

106
107

108 General comments: In general the approach seems reasonable, but we have some specific
109 comments as detailed below.

110
111

Minor comments:

112 Line 13: "At each iteration of the routine, a randomly drawn individual is
113 assigned to a population for which its genotypic probability is a maximum." How is this
114 individual chosen? What is the pool of candidate individuals?

115 Line 29: "Uniformly draw a population at random from the baseline." What
116 exactly does this mean? Each population has equal weight, and then the draw is random?

117 Line 63: "While we agree that this is a relevant measure, unlike *fORCA*, it suffers
118 from being prohibitively slow and may be biased in some circumstances (Anderson
119 2008)." After "unlike *fORCA*", two attributes are listed but only one (being slow) is
120 unlike *fORCA*. The bias described by Anderson et al. (2008) is equally applicable to
121 *fORCA*. See below for more on this point.

122
123

Responses to specific questions:

124
125

1. *Is our approach to linkage disequilibrium and HWE reasonable?*

126

For the most part, but we have several comments to consider.

127
128

1) For both types of analyses, it is important to ensure that the baseline
populations represent single panmictic populations. If not, a Wahlund effect
could cause both HW and LD departures that appear to be data quality issues
but actually reflect population mixture.

129
130
131

2) For both types of analyses, be careful about only using results of tests of
statistical significance. You are really interested in the magnitude of the
effect size here, but P values also depend heavily on sample sizes. Also, the
direction of departure (e.g., heterozygotes excess or deficiency) can be
informative about potential causes.

132
133
134

3) The LD analyses will consider pairs of loci, of which there are $n(n-1)/2$
possible comparisons for n loci. Since n could be 200 or more, this represents
a huge number of pairwise comparisons, each of which could be conducted
for many different populations. Using the Bonferroni correction here would
require consideration of tiny P values, which could lead to unpredictable
results. It is probably more useful to screen for pairs of loci that are
consistently out of equilibrium (using the nominal alpha level) in multiple
populations. Some consideration of effect size (the magnitude of LD) would
also be useful in evaluating how serious a problem any deviations are likely to
cause.

135
136
137
138
139
140
141
142
143
144
145
146

147
148

2. *Is our method to determine the relative value of different treatments of linked markers
advisable? Is the use of *fORCA* as a measure appropriate?*

149
150

The general procedure described at lines 56-68 of Document 8 seems reasonable,
as does the logic for using a procedure that assigns entire individuals rather than making

151 fractional assignments. With the caveats noted below, *fORCA* should be ok as a means to
152 assess *relative* power for correct assignment.

153

154 3. *Are the tests appropriately structured to provide a set of SNPs that will perform well*
155 *for WASSIP?*

156 The proposed methods should produce a set of SNPs with high power to resolve
157 stock identification problems in Western Alaska.

158

159 4. *Does the weighting applied to each set of tests seem reasonable?*

160 The weights chosen are obviously somewhat arbitrary but do not appear to be
161 unreasonable. Because of the applied focus of this project, it is appropriate to assign
162 greater weight to markers that have high power for the local areas of interest. However,
163 we were pleased to see that the criteria include non-trivial weight to markers with wider
164 geographic relevance (10% weight for Pacific Rim individual populations, plus 6% for
165 major non-Alaska groups). This will help ensure that the considerable efforts here to
166 develop markers will have much broader application to the scientific and fishery
167 management communities.

168

169 Minor comments:

170 In the proposed PCA analysis for Pacific-wide assessments, part (iii) is partially
171 redundant as it will include information already used for (i) and (ii)

172 Outside Alaska: we don't necessarily disagree with the particular comparisons
173 proposed, but the rationale for choosing them is not given.

174

175 5. *Are there other measures that would be more appropriate?*

176 Can't think of any offhand.

177

178 General comments about bias and *fORCA*

179 It is important to distinguish between two different types of biases that can
180 potentially arise in evaluations such as those proposed here.

181 The first type of bias, described by Anderson et al. (2008), occurs when one is
182 interested in assessing the power of a particular set of markers to resolve the composition
183 of a mixture comprised of individuals from a specified group of source populations. The
184 ideal way to do this is to create simulated mixtures of individuals, with the genotype of
185 each individual being chosen based on actual allele frequencies in one of the (randomly
186 chosen) source populations. The bias arises because we never know the actual allele
187 frequencies—we only have samples. Because of random sampling error, allele
188 frequencies in samples from the baseline populations will on average be more divergent
189 than are the true population allele frequencies. On average, this factor inflates F_{st} among
190 baseline samples by the magnitude $1/(2S)$, where S is the baseline sample size. When
191 simulated mixtures are constructed using these baseline allele frequencies (which appear
192 more different than the populations actually are), the population assignments will tend to
193 be overly optimistic. Furthermore, the relative importance of sampling error (and hence
194 the bias) will be larger when true genetic differences among populations are very small—
195 as occurs with Western Alaska chum salmon. Anderson et al. (2008) described a simple

196 leave-one-out procedure that eliminates the bias, but the routine described at lines 41-50
197 of Document 10 would be subject to this type of bias.

198 The second type of bias, described by Anderson (2010), applies to locus-selection
199 programs. The bias is not in the locus selection *per se*, but rather in the evaluation of
200 power of the resulting set of loci for population assignment. Anderson (2010) showed
201 that the bias arises because none of the commonly-used software programs for locus
202 selection (including BELS) use proper cross validation. Instead, some of the information
203 used to select the panel of loci is also used to evaluate its performance, and this leads to
204 an overly optimistic assessment of assignment power. We did not see any indication that
205 the combined *FORCA*-BELS approach proposed in Document 10 would *not* be subject to
206 this type of bias. Also, although the authors list 4 methods Rosenberg (2005) evaluated
207 for selecting subsets of loci, they don't explain why they did not consider any of them for
208 the current project.

209 One reason that proper cross-validation is often not done is that it is costly in
210 terms of information content. The "gold standard" of cross validation is to split the data
211 in half: the first half is used to develop the algorithm, the second half to evaluate its
212 performance. However, doing this means that the algorithm is likely to be less precise
213 because it is based on less data. Researchers are thus typically faced with a trade-off
214 between precision in developing the best algorithm (use all the data in the first step) and
215 the downstream consequences (subsequent assessments of performance using the same
216 data will tend to be overly optimistic). Anderson (2010) suggested a simple modification
217 to the cross-validation procedure that retains most of the information without leading to
218 appreciable bias in assessing performance.

219 In summary, both types of biases can lead to overly optimistic assessments of
220 power, which should be a concern given the stated goals of the project. For applications
221 that only consider relative power, these biases might not be important. Also, it might be
222 the case that the proposed locus-selection approach is perfectly fine for selecting an
223 optimal panel of loci, but that the estimates of power to be expected when that panel is
224 applied to real data are biased upwards.

225 Text at lines 84-91 of Document 10 seems to acknowledge at least the bias
226 problem identified by Anderson et al. (2008), but it is not clear that both of the potential
227 sources of bias described above have been fully considered in the documents we
228 reviewed. This topic merits closer scrutiny to determine the optimal way to proceed
229 given project goals.

230
231

232 Anderson, E.C., R.S. Waples, S.T. Kalinowski. 2008. An improved method for
233 estimating the accuracy of genetic stock identification. *Canadian Journal of Fisheries
234 and Aquatic Sciences* 65:1475-1486.

235

236 Anderson, E.C. 2010. Assessing the power of informative subsets of loci for population
237 assignment: standard methods are upwardly biased. *Molecular Ecology Resources*
238 10:701-710.