1

2    Effect of error and missing data on population structure inference using microsatellite data

3

4

5

6

7    Patrick A. Reeves[1,4], Cheryl L. Bowker[2], Christa E. Fettig[2], Luke R. Tembrock[3], and Christopher

8    M. Richards[1]

9

10

11

12    [1]United States Department of Agriculture, Agricultural Research Service, National Laboratory

13        for Genetic Resources Preservation, 1111 South Mason Street, Fort Collins, CO, 80521,

14        U.S.A.

15    [2]Department of Bioagricultural Sciences and Pest Management, Colorado State University, Fort

16        Collins, CO, 80523

17    [3] Department of Biology, Colorado State University, Fort Collins, CO, 80523

18    [4]Corresponding author (pat.reeves@ars.usda.gov; phone: 970-492-7611; fax: 970-492-7605).

19

20    Keywords:  Bayesian inference, coalescent simulation, Hardy-Weinberg equilibrium, multilocus

21        genotype, neighbor-joining, principal coordinate analysis

22    Running title:  Effect of error and missing data.

23

24                                              ABSTRACT

25          Missing data and genotyping errors are common in microsatellite data sets.  We used

26     simulated data to quantify the effect of these data aberrations on the accuracy of population

27     structure inference.  Data sets with complex, randomly-generated, population histories were

28     simulated under the coalescent.  Models describing the characteristic patterns of missing data and

29     genotyping error in real microsatellite data sets were used to modify the simulated data sets.

30     Accuracy of ordination, tree-based, and model-based methods of inference was evaluated before

31     and after data set modifications.  The ability to recover correct population clusters decreased as

32     missing data increased.  The rate of decrease was similar among analytical procedures, thus no

33     single analytical approach was preferable.  For every 1% of a data matrix that contained missing

34     genotypes, 2–4% fewer correct clusters were found.  For every 1% of a matrix that contained

35     erroneous genotypes, 1–2% fewer correct clusters were found using ordination and tree-based

36     methods.  Model-based procedures that minimize the deviation from Hardy-Weinberg

37     equilibrium in order to assign individuals to clusters performed better as genotyping error

38     increased.  We attribute this surprising result to the inbreeding-like nature of microsatellite

39     genotyping error, wherein heterozygous genotypes are mischaracterized as homozygous.  We

40     show that genotyping error elevates estimates of the level of genetic admixture.  Overall, missing

41     data negatively impact population structure inference more than typical genotyping errors.

42

43                                            INTRODUCTION

44          Short, repetitive regions of the genome, known as microsatellite DNA, simple sequence

45     repeats (SSRs), or short tandem repeats (STRs), are commonly used in molecular population

46     genetic studies (Sunnucks 2000; Guichoux *et al.* 2011).  SSR loci exhibit a unique mutational

2

47    mechanism, slipped-strand mispairing, which causes the duplication or deletion of repeat units,

48    resulting in sequence length variation among alleles (Levinson & Gutman 1987). SSR mutation

49    rates vary widely depending on organism, repeat length, and repeat number, but are generally

50    1E3–1E4 times higher than a typical nucleotide substitution rate of 1E-8 per generation (Dallas

51    1992; Chakraborty *et al.* 1997; Vigouroux *et al.* 2002), thus SSR regions provide highly

52    polymorphic markers, useful for distinguishing individuals, reconstructing population history,

53    and estimating demographic parameters.

54    While single nucleotide polymorphisms (SNPs) may gradually supplant SSRs for certain

55    population genetic applications (Brumfield *et al.* 2003), SSRs remain popular. Over 3500 papers

56    that utilized SSRs were published in 2009 (Guichoux *et al.* 2011), and in 2013–2015, ~33% of

57    articles in the journal *Molecular Ecology* included SSRs as a primary data source. SSR

58    development and typing costs have dropped due to next generation sequencing (Gardner *et al.*

59    2011) and improved multiplexing protocols (Butler 2005; Holleley and Geerts 2009). We expect

60    SSRs to remain in use due to low cost and their ability to outperform SNPs with fewer loci for

61    individual identification (Seddon *et al.* 2005); parentage and sibship analysis (Glaubitz *et al.*

62    2003; Wang & Santure 2009); and population structure inference (Liu *et al.* 2005; Glover *et al.*

63    2010).

64    Most SSR data sets contain missing data and erroneous genotypes. Missing data are

65    entered into an SSR data matrix when a particular sample does not produce an interpretable

66    pattern of DNA fragments after PCR amplification. PCR failure is usually caused by poor

67    quality template DNA or improper PCR conditions, including mispriming due to mutations at

68    primer binding sites (Guichoux *et al.* 2011). The frequency of missing data is elevated across

69    loci for samples with poor quality DNA. Suboptimal PCR conditions increase missing data

3

70   across samples for specific loci.  Consequently, missing genotypes typically occur non-randomly

71   in data matrices, clumped in rows or columns.  Missing data can be minimized by re-extracting

72   DNA and repeating amplifications for problematic individuals, redesigning troublesome primers

73   and optimizing PCR conditions, or excluding individuals and loci with high failure rates from the

74   matrix.

75        SSR genotyping errors arise from three main sources:  the occurrence of null alleles at a

76   locus, preferential amplification of small DNA targets during PCR, and stuttered visualization of

77   amplification products (DeWoody *et al.* 2006; Guichoux *et al.* 2011).  "Null alleles" are

78   genotypic variants that fail to amplify under the conditions specified for the locus.  They often

79   occur due to mutations at primer binding sites that inhibit amplification but may also arise from

80   poor template quality.  Because nothing is amplified, null alleles, by definition, cannot be

81   observed and are not scored.  Consequently, in diploid organisms, the genotype entered into the

82   data matrix becomes—erroneously—homozygous for the other, visible allele, or, if both alleles

83   are nulls, missing data.  Several analytical approaches have been devised to detect null alleles

84   and estimate their frequency (Chakraborty *et al.* 1992; Raymond & Rousset 1995; Van

85   Oosterhout *et al.* 2004; Kalinowski *et al.* 2007) although only about 40% of studies use them

86   (Guichoux *et al.* 2011).  In most cases, eliminating null alleles requires primer redesign outside

87   of highly mutable regions (Dakin & Avise 2004; Chapuis & Estoup 2007).

88        Preferential amplification of short DNA sequences causes "large allele dropout" error.

89   All else equal, short sequences are more efficiently amplified than long.  In a heterozygote with

90   differently-sized alleles this bias may prevent the signal for the large allele from rising above the

91   detection threshold, with the consequence that a heterozygous genotype will be erroneously

92   scored as homozygous for the smaller allele (Wattier *et al.* 1998; Björklund 2005).  Large allele

4

93   dropout can be mitigated in some cases by excluding loci where amplicons exceed 200 bp (Sefc

94   *et al.* 2003).

95        Slipped-strand mispairing during PCR results in the production of shadow peaks around

96   the amplified allele (Murray *et al.* 1993), a phenomenon termed "stutter".  Stutter peaks are

97   usually smaller than the target, and deviate in size by multiples of the repeat unit length, with

98   progressively decreasing signal (Shinde *et al.* 2003).  As with other types of error, stutter causes

99   heterozygotes to be scored as homozygous, but always for the larger of the two alleles, and only

100  when the alleles differ in size by a single repeat unit.  Stutter can be reduced by avoiding SSRs

101  with dinucleotide repeats (Chambers & MacAvoy 2000), decreasing denaturation temperature

102  (Olejniczak & Krzyzosiak 2006), and using highly processive polymerases (Davidson *et al.*

103  2003).

104       The final assembly of an SSR data set requires considerable care.  Filling all cells in a

105  data matrix is time consuming; poor-performing loci must be optimized and recalcitrant

106  individuals must be extracted and genotyped repeatedly.  Consequently, the typical course for

107  dealing with missing data is to eliminate problematic individuals or loci.  This can produce

108  biased sampling because the frequency of missing data may be similar in related populations

109  (Amos 2006).  Some authors have recommended that error rates be reported and efforts made to

110  assess the reliability of conclusions given uncertainty in the genotypes (Bonin *et al.* 2004;

111  Broquet & Petit 2004; Hoffman & Amos 2005).  This practice is increasing—error rate estimates

112  can be found in about a quarter of published papers (Guichoux *et al.* 2011)—but it requires

113  sample replication and extensive post-genotyping data analysis, increasing costs.  If missing data

114  and genotyping error were to impact the accuracy of population structure inference in only minor

115  ways these steps might be avoided.

5

116     The effect of missing data and error on linkage mapping (Hackett & Broadfoot 2003),

117     parentage analysis (Dakin & Avise 2004; Hoffman & Amos 2005; Kalinowski *et al.* 2007), and

118     the estimation of population genetic parameters (Chapuis & Estoup 2007; Hall *et al.* 2012; Peel

119     *et al.* 2013) has been studied in detail, but there have been few investigations of the effect on

120     population structure inference (Pompanon *et al.* 2005; Carlsson 2008; Chapuis *et al.* 2008).  In

121     this study we seek to quantify the extent to which missing data or genotyping error impact

122     population structure inference, to assist researchers in developing strategies to produce SSR data

123     sets that maximize accuracy while minimizing costs.

124

125                                    MATERIALS AND METHODS

126     Highly polymorphic, neutral marker data were simulated using a coalescent model and

127     the software MSMS (Ewing & Hermisson 2010).  The simulation approach is detailed in Reeves

128     *et al.* (2012).  Briefly, 1E4 data sets were generated, each containing 500 diploid individuals

129     equally distributed among 50 populations, and 50 unlinked loci.  An asymmetric island migration

130     model was created by randomly assigning migration rates between populations.  The population

131     scaled mutation rate (θ) was varied from 0–0.5.  The underlying mutation model of MSMS is an

132     infinite sites model.  Unique binary strings output by MSMS were converted into uniquely

133     named alleles following Huelsenbeck and Andolfatto (2007), rendering the infinite sites model

134     as an infinite allele model.  We used an infinite allele model instead of a stepwise mutation

135     model as an expedient, and because it is not clear the extent to which the stepwise mutation

136     model fits real SSR data (Gaggiotti *et al.* 1999).  It is not unusual for SSR data sets to better fit

137     an infinite allele model than stepwise mutation (e.g. Estoup *et al.* 1995; O'Connell *et al.* 1997).

138     The primary defining feature of real SSR data is a limited number of alleles per locus (Paetkau *et*

6

139    *al.* 1997).  Therefore, data sets containing levels of polymorphism comparable to typical SSR

140    loci (on average < 30 alleles per locus, following Kalinowski 2002) were subsampled from the

141    original 1E4 data sets.  A total of 1367 were found and used for further analysis.

142         Simulated data sets were altered to include missing or erroneous genotypes using models

143    describing the distribution of these data aberrations in typical SSR data sets. In the missing data

144    model, the percent of the matrix containing missing data was varied between zero and 25.  A

145    "clumping" parameter was used to bias placement of missing genotypes towards certain loci or

146    individuals.  The clumping parameter was varied among data sets from zero, which caused a

147    uniform distribution of missing data, to ten, which elevated the probability ten-fold that the next

148    missing genotype would occur in a row or column already containing missing data.  Missing

149    genotypes were substituted for known genotypes one by one, with the probability of conversion

150    for each row and column adjusted after each substitution using the clumping parameter, until the

151    specified percentage of missing data was reached.

152         To create data sets that varied in their propensity to contain null alleles, a

153    data-set-specific maximum null allele frequency parameter ($v_d$) was defined, and selected at

154    random from 0 to 20 percent, following Dakin and Avise (2004).  Within each data set, the

155    locus-specific null allele frequency parameter ($v_l$) was chosen at random from 0 to $v_d$ for each

156    locus. The number of null alleles per locus was then defined as $v_l$ multiplied by the number of

157    distinct alleles at the locus, with the alleles that were to act as nulls chosen randomly.  Alleles

158    defined as nulls were treated as unknown, with the consequence that heterozygous genotypes

159    became homozygous for the non-null allele, and homozygous null genotypes became missing

160    data, in the modified data set.

161     To simulate large allele dropout, it was necessary to assign a probability of dropout for

162     each allele that was proportional to its size. A locus-specific maximum probability of dropout

163     ($\delta_l$) was chosen at random for each locus from 0 to $\delta_d$, the data-set-specific dropout probability.

164     The ceiling on $\delta_d$ was set to 0.5 based on empirical studies (Taberlet *et al.* 1996; Gagneux *et al.*

165     1997; Buchan *et al.* 2005). We assumed a curvilinear function relating dropout probability to

166     allele size, with the largest allele at each locus having a dropout probability of $\delta_l$. Coalescent

167     simulations only provide information on allelic state, not allele size, so relative allele sizes ($\sigma_a$)

168     were assigned to all alleles for each locus by randomly sampling an exponential distribution with

169     rate parameter $\lambda$, varied by locus from 0 to 10. The probability of retention (i.e. the probability

170     that an allele does not drop out) was then computed for each allele as one minus the cumulative

171     distribution function of the exponential (CDF $= 1 - e^{-\lambda\sigma_a}$) rescaled to have a maximum value of

172     $\delta_l$. In this way, data sets exhibited varying levels of overall "dropout proneness" governed by the

173     parameter $\delta_d$. Within data sets, the largest alleles at a locus always had the highest dropout

174     probability, but some loci were characterized by dropout probabilities that declined uniformly

175     with consecutively smaller allele size (when $\lambda \rightarrow 0$), while others approximated a threshold effect

176     where alleles above a particular size were highly dropout prone (when $\lambda \rightarrow 10$). If a uniform

177     random number from 0 to 1 exceeded the probability of retention, the allele was made to drop

178     out, and the data set was modified accordingly. Thus, the effect of simulated large allele dropout

179     was to convert heterozygous genotypes to small allele homozygotes with probability $\delta_l(1 -$

180     $e^{-\lambda\sigma_a})$.

181     To model stutter error, we identified all heterozygotes having consecutively-sized alleles

182     using the arbitrary sizes from the large allele dropout model. These genotypes are called

183     "adjacent-allele heterozygotes" (Hoffman & Amos 2005). Stutter error could only affect these

8

184 genotypes, but was not assumed to affect all of them. A model was developed to assign a

185 "probability of stutter error occurrence" to each adjacent-allele heterozygous genotype, thus

186 permitting us to vary the "stutter error proneness" between data sets and among loci. A

187 data-set-specific average probability of stutter error at adjacent-allele heterozygotes ($\overline{\tau_d}$) was

188 randomly chosen from zero to one. $\widetilde{\tau_d}$, the standard deviation of $\overline{\tau_d}$, was randomly set from zero

189 to one. The locus-specific probability of stutter error, $\tau_l$, was sampled from the normal

190 distribution defined by $\overline{\tau_d}$ and $\widetilde{\tau_d}$. The conversion of adjacent-allele heterozygotes into

191 large-allele homozygotes then occurred with probability $\tau_l$, the locus-specific conversion

192 probability.

193 Genotypes in the unmodified matrices were altered in the order that errors would arise

194 during the genotyping process. Null allele errors, which are caused by mispriming, were added

195 first, followed by large allele dropout errors, caused by poor amplification, then by stutter errors,

196 caused by slippage during amplification but attributable primarily to poor scoring procedures.

197 Realized error rates were then recalculated for each error type, for each data set.

198 Three categorically-distinct analytical approaches for inferring population structure were

199 used. First, we applied a class of Bayesian Markov chain Monte Carlo (MCMC) methods

200 introduced by Pritchard *et al.* (2000) in the software STRUCTURE. To avoid *ad hoc* model

201 selection procedures for determining the number of populations (K) (Evanno *et al.* 2005), we

202 used INSTRUCT v1.0 (Gao *et al.* 2007) and STRUCTURAMA v1.0 (Huelsenbeck and

203 Andolfatto 2007) instead of STRUCTURE. For INSTRUCT analyses, we used no-admixture

204 (mode 0) and admixture (mode 1) models, as well as a model that estimates individual

205 inbreeding coefficients simultaneously with individual assignment (mode 5). Modes 0 and 1 of

206 INSTRUCT are comparable to no-admixture and admixture models of STRUCTURE. A single

9

207  Markov chain was run for 1E5 generations and sampled every 25, with the initial 12500

208  generations discarded. The deviance information criterion (DIC) was used to determine the best

209  value of K between 1 and 50 (Gao *et al.* 2011). A discrete assignment was created by assigning

210  individuals to clusters based on the highest assignment probability in the Q-matrix.

211  STRUCTURAMA estimates K alongside individual assignment. Chains were run as for

212  INSTRUCT, the prior on number of populations was set to two, and no admixture was allowed.

213      Second, we used a tree-based approach. Neighbor-joining (NJ) trees were constructed

214  using NTSYS (Rohlf 2008). Inter-individual distances were computed using Lynch's (1990)

215  band sharing coefficient. Populations were counted as correctly inferred when they existed as a

216  monophyletic group in the resulting tree. Third, an ordination method was applied. We used

217  PCOMC, a procedure that couples ordination with cluster analysis to simultaneously determine

218  population number and membership (Reeves & Richards 2009). Principal coordinate analysis

219  was performed using NTSYS. Distance matrices, computed as for NJ, were double-centered

220  prior to the calculation. Principal coordinate values were weighted according to their

221  contribution to the total variance, then subject to the density clustering algorithm PROC

222  MODECLUS (SAS Institute, Cary, NC).

223      Correct and incorrect clusters resulting from application of each analytical method to

224  each data set were counted. Correct clusters were those that contained all 10 individuals that

225  belonged to a single population as specified in the coalescent simulation model, and no others.

226  Incorrect clusters contained some, but not all, members of a population in the model, or

227  individuals from more than one population. The "performance ratio" was defined as the number

228  of correct clusters resulting from analysis of a modified data set divided by the number of correct

229  clusters in the unmodified data set from which it was derived. A performance ratio < 1 indicates

230    that data modification reduces accuracy, while a performance ratio > 1 indicates improved

231    accuracy.  We used the partition distance of Gusfield (2002) as a second, less strict measure of

232    accuracy for methods that produce partitions (INSTRUCT and STRUCTURAMA).  The

233    partition distance ($PD$) has been used previously for this purpose, and is useful because it can

234    quantify partially correct matches between clusters, unlike the performance ratio (Huelsenbeck

235    and Andolfatto 2007; Choi and Hey 2011).  It is defined as the number of elements that must be

236    moved between clusters to make one partition identical to another.  We normalized the $PD$ to the

237    range 0–1 by dividing by the maximum possible $PD$ (Charon $et$ $al.$ 2006), calling the result $PD_n$.

238    We define the "partition distance ratio" as $\frac{1-PD_n^{modified}}{1-PD_n^{unmodified}}$, where numerator and denominator are

239    formulated as similarities to simplify comparison with the performance ratio.  $PD_n^{modified}$ is the

240    normalized partition distance between the partition resulting from analysis of the modified data

241    and the 50 cluster partition defined in the coalescent simulation model (likewise for

242    $PD_n^{unmodified}$).  The performance ratio and the partition distance ratio are non-identical, but

243    correlated, measures of the accuracy of population structure inference (Supplementary Table 1).

244    We preferentially report the performance ratio, because its interpretation is intuitive, and it is

245    applicable to all methods.

246         The false discovery rate ($FDR$) was calculated as the number of incorrect clusters divided

247    by the total number of clusters returned. A ratio, $\frac{FDR^{modified}}{FDR^{unmodified}}$, was calculated to express the

248    difference in probability of recovering incorrect clusters between modified and unmodified data

249    sets.  Multiple regression and likelihood analysis were used to examine the effect of model

250    factors on the accuracy of inference.

251

11

252                                                    RESULTS

253     *Data set validation.* The 1367 simulated data sets contained an average of $14.48 \pm 8.46$ (1 sd)

254     alleles per locus. Population structures ranged from virtually panmictic to highly subdivided

255     (Figure 1). Levels of population subdivision, as quantified by Hedrick's (2005) G'$_{st}$, varied from

256     0.006 to 0.999, with an excess of low G'$_{st}$ values. The degree to which data sets were modified

257     by the missing data model was not significantly related to the level of population differentiation

258     ($r = 0.03$, $p = 0.32$); application of the error models resulted in a slight, but significant,

259     negative correlation between G'$_{st}$ and the proportion of genotypes modified ($r = -0.14$,

260     $p < 0.0001$ (Supplementary Figure 1).

261     *Performance.* The analytical methods differed in accuracy when unmodified data sets were

262     analyzed (Figure 2). For a given level of genetic subdivision, correct clusters were found in NJ

263     trees at a much higher frequency than in INSTRUCT, STRUCTURAMA or PCOMC analyses.

264     This is not surprising because the criterion for identifying correct clusters for NJ was less strict

265     than for other methods—a correct inference occurred when a node existed that defined a

266     population correctly in the tree, but there was no mechanism to determine which nodes defined

267     populations. PCOMC had the lowest rate of correct cluster recovery. STRUCTURAMA

268     recovered more correct clusters and had a lower false discovery rate than INSTRUCT for high

269     levels of genetic subdivision, while INSTRUCT was more accurate at lower levels, regardless of

270     mode.

271          To avoid undefined values when calculating the performance ratio, data sets were

272     excluded when zero correct clusters were inferred with unmodified data. Likewise, for the

273     partition distance ratio, data sets were excluded when $PD_n^{unmodified} = 1$, i.e. when the observed

274     partition was maximally distant from the simulated partition. Because the coalescent model was

275   complex and population subdivision was often low, a substantial number of data sets were

276   excluded using this restriction. The total number of useful data sets ranged from 63 to 473,

277   depending on analytical method, when measured using the performance ratio (Table 1), and 115

278   to 464 for the partition distance ratio (Supplementary Table 1). Population subdivision in the

279   excluded data sets was low (G'$_{st}$ ~5-fold lower than retained data sets), approaching or exceeding

280   the limits of resolution of the methods applied.

281        Performance decreased in a roughly linear manner as missing data increased (Figure 3a–

282   f). The slope of the regression line was significantly different from zero (at $\alpha = 0.01$,

283   Holm-Bonferroni corrected, here and throughout) for all methods except 'INSTRUCT

284   inbreeding' (mode 5) and PCOMC (Table 1). $R^2_{adj}$ values for significant regressions ranged

285   from 0.08 to 0.27 and the slope of performance loss ranged from $m = -1.8 - -3.5$ between

286   analytical approaches (Table 1). The results using the partition distance ratio were similar. A

287   significant negative correlation was found for 'INSTRUCT no admixture' (mode 0) and

288   'INSTRUCT admixture' (mode 1); the correlation was not significant for 'INSTRUCT

289   inbreeding' or STRUCTURAMA (Supplementary Table 1).

290        Taking into account 95% confidence intervals on the slopes, a data matrix with 5%

291   missing data is predicted to result in recovery of at least 72–81% of the correct clusters that

292   could be recovered with a complete data matrix. Based on our data, 95% of recoverable clusters

293   should be found when data matrices contain, for INSTRUCT, 2.5–2.7% missing data, or, for the

294   other methods, 1.5–1.6% missing data. The FDR ratio increased significantly ($\alpha = 0.05$) with

295   missing data for all methods except 'INSTRUCT inbreeding' and PCOMC, indicating that, in

296   addition to fewer correct clusters, users should expect more erroneous clusters as missing data

297   increase.

13

298     The clumping parameter had a statistically significant effect on performance for NJ only

299     (Table 2).  Performance of NJ was higher when missing values were more clumped ($\beta = 0.02$).

300     Using ratios of the Akaike weights, the linear model for percent missing data alone had a much

301     higher probability ($10^{25}$-fold) than the model for clumping parameter for NJ.  Regardless of

302     method, less than 4% of the variance in performance was attributable to clumping—a slight

303     effect—thus the clumping parameter, despite adding realism to the simulation, was largely

304     dispensable.

305     The effect of genotyping error differed between categories of analytical method.  For

306     distance methods NJ and PCOMC, performance declined as erroneous data increased (Figure

307     3k,l).  The slopes of the regressions were significantly different from zero and about half the

308     magnitude found for missing data (Table 1).  A matrix with 5% erroneous data should result in

309     recovery of 85% (NJ) or 82% (PCOMC) of the clusters that would be found if the data set

310     contained no error.  Ninety five percent of recoverable clusters were found when 2.9% (NJ) or

311     3.6% (PCOMC) of the data matrix was erroneous.  For NJ, large allele dropout and stutter had

312     the greatest effect on deteriorating performance (Table 2).

313     In contrast, accuracy of model based methods improved as erroneous data increased.

314     When using the performance ratio, the slope of the regression was positive for all methods

315     (0.41–1.36) (Figures 3g-j, Table 1).  The relationship was statistically significant for

316     STRUCTURAMA.  When using the partition distance ratio, a significant positive correlation

317     was found for all methods (Figure 4, Supplementary Table 1).  Large allele dropout was the most

318     important model effect for INSTRUCT and STRUCTURAMA, explaining 2–10% of the

319     variation in performance (combined model, 5–16%), and holding 18–78% of the linear models'

320    Akaike weight (model probability).  The FDR ratio did not change significantly with increasing

321    error for any method.

322

323                                    DISCUSSION

324        Genomewide genotyping approaches are gaining popularity for studies of population

325    structure but SSRs continue to be used due to low cost and high power of inference.  Although

326    SSR data sets are prone to containing missing data and erroneous genotypes, few studies have

327    examined their impact on analyses of population structure (Pompanon *et al.* 2005).  We explored

328    the effect of missing data and genotyping errors on population structure inference using

329    model-based Bayesian MCMC procedures (INSTRUCT, STRUCTURAMA), a tree-based

330    method (NJ), and an ordination approach (PCOMC).  Our goal is to provide users of SSRs with

331    insight for how much missing data and error might be tolerated in order to achieve a desired

332    level of accuracy.

333        In order to make general recommendations it was necessary to explore a diverse set of

334    population structures.  This was accomplished using coalescent modeling with key parameters—

335    mutation, migration rates, migration directionality—set stochastically, but within plausible

336    ranges (Reeves *et al.* 2012).  The resultant data sets exhibited a large range of complex

337    population structures with widely varying levels of subdivision (Figure 1, Supplementary Figure

338    1).  Nevertheless, the extent to which simulated data can ever accurately represent nature is

339    debatable.  We attempted to produce a realistic subsample from the universe of plausible

340    population structures.  Other models describing population divergence exist, most notably the

341    isolation with migration model (Nielsen & Wakeley 2001).  Our conclusions should be

15

342    interpreted as limited to modified island models similar to those simulated, and should not be

343    expected to precisely predict outcomes for any single data set.

344          Key properties peculiar to SSR data sets were modeled, including the clumped nature of

345    missing data, and the three best understood sources of error:  null alleles, large allele dropouts,

346    and stutter artifacts.  We did not include an assessment of human error, which contributes

347    substantially (Bonin *et al.* 2004), but is not easily modeled.  While most real data sets will be

348    affected by both missing data and error, we elected to study these factors independently because

349    the course of action for mitigation seems distinct.

350    *Impact of missing data.*  For data sets modified with missing data three basic observations were

351    made.  First, the degradation of performance as missing data increased was roughly linear

352    (Figure 3a–f).  Second, the rate at which performance declined was similar among methods

353    (Table 1).  Third, the degree to which missing data were clumped by locus or DNA sample was

354    unimportant (Table 2)—contiguous blocks of missing genotypes did not affect performance

355    much more than an equal number of randomly-distributed missing genotypes.  These

356    observations lead to a simple prediction:  for every 1% of a data matrix that contains missing

357    data, researchers should expect to recover 2-4% fewer correct clusters than would be inferred

358    with a complete data set.

359    *Impact of error on distance methods.*  Error caused a decline in performance for the tree-based

360    NJ method and for the ordination approach PCOMC.  Researchers should expect a ~2% decrease

361    in the number of correct clusters recovered for every 1% of the data matrix impacted by error

362    when using these methods.

363          Although NJ and PCOMC are quite distinct mathematically, the similarity in response

364    between the methods may arise from their use of the same distance metric.  Distance was

365    calculated as $1 - S_{xy}$, where $S_{xy}$ is the ratio of alleles common to two individuals, over the total,

366    or $\frac{|A \cap B|}{|A \cup B|}$, i.e. the Jaccard index (here, A and B represent the sets of unique alleles from each

367    individual). The mechanism by which performance decreases is loss of precision in the genetic

368    distance estimate. Conversion of heterozygotes to homozygotes decreases the number of alleles,

369    lowering $|A \cup B|$. The maximum number of distinct values the index can assume, $|A \cup B| +$

370    1, also decreases, imposing a limit on the assembly of individuals into different clusters.

371    Missing data had a stronger effect than error (Figure 3e,f,k,l). Conversion of a heterozygous

372    locus to missing data may reduce $|A \cup B|$ by two, whereas with error $|A \cup B|$ can only

373    decrease by one, thus the erosion of resolving power is generally more rapid with missing data.

374    *Impact of error on model-based methods.* We observed a positive relationship between the

375    amount of error introduced into a data set and the accuracy of inference when using model-based

376    methods (Figures 3-4, Table 1, Supplementary Table 1). Our data suggest that a matrix

377    containing 25% erroneous data could cause the recovery of 30–70% more correct clusters than

378    an error-free matrix. Hereafter, we attempt to explain this unexpected result.

379        When a population consists of several subpopulations, a "heterozygote deficit" occurs,

380    where the observed heterozygosity of the population analyzed as a whole is lower than predicted

381    under Hardy-Weinberg equilibrium. This "Wahlund effect" (Wahlund 1928) forms the

382    theoretical justification, and source of signal, for a class of model-based clustering methods that

383    includes STRUCTURE, STRUCTURAMA, INSTRUCT, and others (e.g. Dawson and Belkhir

384    2001; Corander *et al.* 2003), which minimize the deviation from Hardy-Weinberg expectations

385    by partitioning individuals into distinct subpopulations. Unfortunately, the signal upon which

386    these methods rely is not uniquely generated by population subdivision. Heterozygote deficit

387    also occurs as a consequence of selfing or inbreeding (Gao *et al.* 2007). Likewise, because SSR

17

388    genotyping error typically involves the conversion of heterozygotes to homozygotes, it too

389    induces a heterozygote deficit.  In our simulation, a decline in observed heterozygosity with

390    increasing error inflated the inbreeding coefficient $F = 1 - \frac{H_o}{H_e}$.  As error increased, the

391    heterozygote deficit $D = \frac{H_o - H_e}{H_e}$ likewise increased (since D is the additive inverse of F).  Thus

392    error produces an inbreeding-like signature in data sets, which magnifies the Wahlund effect and,

393    as a byproduct, inflates the level of genetic subdivision, causing increased signal (Supplementary

394    Figure 2).

395           If this sort of signal amplification is important we would expect model factors ΔF (the

396    change in inbreeding between unmodified and modified data due to error), and/or ΔG'st (the

397    change in genetic subdivision) to drive the positive relationship between error and accuracy.

398    Some evidence supports this.  When accuracy was measured using the partition distance ratio,

399    ΔG'st was the second-most important model effect for all INSTRUCT modes, explaining 3-10%

400    of the variation (Supplementary Table 3).  For 'INSTRUCT no admixture', ΔF was the most

401    important.  ΔG'st and ΔF were also important model effects for STRUCTURAMA, explaining 5-

402    15% of the variation in the partition distance ratio.  However, for STRUCTURAMA, the

403    overwhelming weight of evidence was that ΔK, the difference in the number of populations

404    inferred, was the primary driver of the response.  Thus the improvement in accuracy observed for

405    STRUCTURAMA seems to have an additional underlying cause.

406           STRUCTURAMA uses a Dirichlet process prior to simultaneously infer K and individual

407    assignment.  A "Dirichlet process" describes a distribution of probability distributions and is a

408    mathematical construct commonly used as a prior in MCMC procedures to categorize items into

409    groups when the number of groups is unknown.  During progression of the Markov chain in

410    STRUCTURAMA, a critical decision is made for each individual at each step: whether the

411    individual belongs to an existing cluster or should be assigned to a new cluster. This decision

412    affects both assignment and the inference of K. The algorithm proceeds roughly as follows.

413    First, an individual $i$ is selected at random from the existing partition and removed. Individual $i$

414    is then re-assigned to whichever of the K clusters it fits best, or to a new cluster, all by itself.

415    The probability that the individual is re-assigned to an existing cluster, $k$, is dependent on the

416    number of individuals, $\eta$, in that cluster (large clusters are more attractive) and the marginal

417    posterior probability of drawing $i$'s genotype from cluster $k$. In contrast, the probability of

418    assigning $i$ to a new cluster depends on the concentration parameter, $\alpha$, of the Dirichlet process

419    (the higher $\alpha$ is, the lower the probability that two randomly drawn individuals belong to the

420    same cluster) and the probability of drawing $i$'s genotype from the prior distribution of allele

421    frequencies, where all alleles are equiprobable.

422         In our unmodified data sets, the deviation from Hardy-Weinberg equilibrium is due to

423    subpopulation structure plus some minor, random effects of the coalescent. However, when

424    error is introduced, part of the deviation is then derived from the inbreeding-like effect of SSR

425    error. With error, the relative probability of assigning an individual to an existing cluster versus

426    a new cluster decreases—data with heterozygote deficit, by definition, do not fit nicely into a

427    Hardy-Weinberg scenario—and, accordingly, it becomes more likely that an individual will be

428    assigned to a new cluster, thereby increasing K. Consistent with this expectation, K increased

429    significantly with error (Figure 5a). When more clusters are available in which to distribute

430    individuals, more correct inferences can be made (equivalently, the number of correct inferences

431    possible is suppressed for low levels of error). Consequently, performance improves with

432    increasing error primarily because K increases. While convenient, the STRUCTURAMA

19

433    approach for defining K may be incautious. Miller and Harrison (2014) argued that Dirichlet

434    process models should not be used to infer the number of components.

435        INSTRUCT uses the Deviance Information Criterion (DIC) to select the best K value

436    from a user-defined range. In the present case, a model with K=50 will usually have a higher

437    likelihood than one with, say, K=4, but that does not necessarily mean it is better, it may merely

438    be overfit. The DIC, like other model selection statistics, deals with this by imposing a penalty

439    on the likelihood for adding parameters. The DIC value is, essentially, the mean likelihood

440    across MCMC draws at stationarity, penalized for increasing K. The K value that minimizes the

441    DIC across the range is deemed optimal. This method for choosing K appears largely insensitive

442    to the level of error, and in this sense seems preferable to the Dirichlet process used by

443    STRUCTURAMA. There is currently no consensus on the best procedure to estimate the

444    number of groups in a finite mixture. Until the basic mathematical issues are resolved, the

445    biological problem of estimating K for population structure inference will remain more craft than

446    science.

447    *Impact of error on the estimation of admixture proportions.* Using an admixture model and

448    STRUCTURE, Gao *et al.* (2007) showed that inbreeding causes two spurious results:

449    exaggerated levels of admixture and elevated likelihoods for higher K values. Falush *et al.*

450    (2003, p. 1572) predicted this: "…situations that might cause additional populations to be

451    inferred by STRUCTURE…include a significant frequency of inbreeding, cryptic relatedness

452    within the sample, or the presence of null alleles." The implication is that

453    heterozygote-deficit-inducing factors other than population subdivision might promote the

454    inference of more clusters, or more admixture among clusters, than is actually the case. For

455    'INSTRUCT admixture', the relative level of inferred admixture increased with error, and

20

456  likelihoods were elevated for higher K values (Figure 5b,d). This did not, however, translate into

457  the inference of higher K values under the DIC ($p > 0.19$).

458       To reduce artifacts from conflation of signal types, Gao *et al.* (2007) relaxed the

459  assumption of Hardy-Weinberg equilibrium to estimate an inbreeding parameter alongside

460  subpopulation membership. Relative to 'INSTRUCT admixture,' the relationship between error

461  and admixture for 'INSTRUCT inbreeding' remained positive, but that between K and the model

462  likelihood was reversed: *Ln*L decreased as K increased (Figure 5c,d). Thus the inbreeding

463  model of Gao *et al.* (2007) compensates in some ways for genotyping error, but does not resolve

464  the issue of admixture overestimation. The effects of error might be mitigated by explicitly

465  modeling it as an additional parameter that alters Hardy-Weinberg expectations at each locus

466  independently, distinct from the inbreeding coefficient estimated by INSTRUCT, which alters

467  expectations uniformly across loci.

468       Whether over-estimation of admixture proportions due to genotyping error is a problem

469  in real data sets is unknown, but the potential exists. False confidence in our understanding of

470  the partitioning of genetic variation in nature notwithstanding, artifactually-elevated admixture

471  estimates could have important economic implications, in endangered species management or

472  genomewide association analysis, for example. Newer forms of data, like SNPs, are not immune

473  to the problem. High throughput genotyping and next generation sequencing approaches vary in

474  their capacity to accurately identify heterozygotes (Harismendy *et al.* 2009; Skotte *et al.* 2013),

475  sometimes defaulting to the homozygous condition, akin to SSR error.

476  *Recommendations*. Missing data exert a greater negative effect on the inference of population

477  structure than typical SSR genotyping errors. As missing data increase, the number of correct

478  clusters recovered decreases and the number of incorrect clusters increases. In order to

479 efficiently produce SSR data sets for inferring population structure we recommend limiting the

480 percent of a matrix that contains missing data to ~2%, unless a greater amount can be justified

481 based on the particular hypotheses under examination.  This will allow researchers to retain most

482 of the resolving power of their data while not incurring the extra costs associated with

483 completing a data matrix.  For analyses reliant on simple distance metrics, we recommend

484 limiting the error rate to ~4% of scored genotypes. Model-based population structure inference

485 methods handle genotyping error well.  We recommend the use of admixture models to infer

486 cluster membership, but caution that admixture estimates may be artificially elevated.

487

498                                        REFERENCES
499 Amos W (2006) The hidden value of missing genotypes. Molecular Biology and Evolution, 23,
500        1995–1996.
501 Björklund M (2005) A method for adjusting allele frequencies in the case of microsatellite allele
502        drop-out. Molecular Ecology Notes, 5, 676–679.
503 Bonin A, Bellemain E, Bronken Eidesen P et al. (2004) How to track and assess genotyping
504        errors in population genetics studies. Molecular Ecology, 13, 3261–3273.
505 Broquet T, Petit E (2004) Quantifying genotyping errors in noninvasive population genetics.
506        Molecular Ecology, 13, 3601–3608.
507 Brumfield RT, Beerli P, Nickerson DA, Edwards SV (2003) The utility of single nucleotide
508        polymorphisms in inferences of population history. Trends in Ecology & Evolution, 18,
509        249–256.
510 Buchan JC, Archie EA, Van Horn RC, Moss CJ, Alberts SC (2005) Locus effects and sources of
511        error in noninvasive genotyping. Molecular Ecology Notes, 5, 680–683.
512 Butler JM (2005) Constructing STR multiplex assays. In: Forensic DNA Typing Protocols , pp.
513        53–65. Springer.
514 Carlsson J (2008) Effects of microsatellite null alleles on assignment testing. Journal of Heredity,
515        99, 616–623.

Chakraborty R, Kimmel M, Stivers DN, Davison LJ, Deka R (1997) Relative mutation rates at di-, tri-, and tetranucleotide microsatellite loci. Proceedings of the National Academy of Sciences, 94, 1041–1046.

Chambers GK, MacAvoy ES (2000) Microsatellites: consensus and controversy. Comparative Biochemistry and Physiology Part B: Biochemistry and Molecular Biology, 126, 455–476.

Chapuis M-P, Estoup A (2006) Microsatellite null alleles and estimation of population differentiation. Molecular Biology and Evolution, 24, 621–631.

Chapuis M-P, Lecoq M, Michalakis Y, Loiseau A, Sword GA, Piry S, Estoup A (2008) Do outbreaks affect genetic population structure? A worldwide survey in Locusta migratoria, a pest plagued by microsatellite null alleles. Molecular Ecology, 17, 3640–3653.

Charon I, Denœud L, Guénoche A, Hudry O (2006) Maximum transfer distance between partitions. Journal of Classification, 23, 103–121.

Choi SC, Hey J (2011) Joint inference of population assignment and demographic history. Genetics, 189, 561–577.

Corander J, Waldmann P, Sillanpää MJ (2003) Bayesian analysis of genetic differentiation between populations. Genetics, 163, 367–374.

Dakin EE, Avise JC (2004) Microsatellite null alleles in parentage analysis. Heredity, 93, 504–509.

Dallas JF (1992) Estimation of microsatellite mutation rates in recombinant inbred strains of mouse. Mammalian Genome, 3, 452–456.

Davidson JF, Fox R, Harris DD, Lyons-Abbott S, Loeb LA (2003) Insertion of the T3 DNA polymerase thioredoxin binding domain enhances the processivity and fidelity of Taq DNA polymerase. Nucleic Acids Research, 31, 4702–4709.

Dawson KJ, Belkhir K (2001) A Bayesian approach to the identification of panmictic populations and the assignment of individuals. Genetical Research, 78, 59–77.

DeWoody J, Nason JD, Hipkins VD (2006) Mitigating scoring errors in microsatellite data from wild populations. Molecular Ecology Notes, 6, 951–957.

Estoup A, Garnery L, Solignac M, Cornuet J-M (1995) Microsatellite variation in honey bee (Apis mellifera L.) populations: hierarchical genetic structure and test of the infinite allele and stepwise mutation models. Genetics, 140, 679–695).

Evanno G, Regnaut S, Goudet J (2005) Detecting the number of clusters of individuals using the software STRUCTURE: a simulation study. Molecular Ecology, 14, 2611–2620.

Ewing G, Hermisson J (2010) MSMS: a coalescent simulation program including recombination, demographic structure and selection at a single locus. Bioinformatics, 26, 2064–2065.

Falush D, Stephens M, Pritchard JK (2003) Inference of population structure using multilocus genotype data: linked loci and correlated allele frequencies. Genetics, 164, 1567–1587.

Gaggiotti OE, Lange O, Rassmann K, Gliddon C (1999) A comparison of two indirect methods for estimating average levels of gene flow using microsatellite data. Molecular Ecoloty, 8, 1513–1520.

Gagneux P, Boesch C, Woodruff DS (1997) Microsatellite scoring errors associated with noninvasive genotyping based on nuclear DNA amplified from shed hair. Molecular Ecology, 6, 861–868.

Gao H, Bryc K, Bustamante CD (2011) On Identifying the Optimal Number of Population Clusters via the Deviance Information Criterion. PLoS ONE, 6, e21014.

Gao H, Williamson S, Bustamante CD (2007) A Markov Chain Monte Carlo approach for joint inference of population structure and inbreeding rates from multilocus genotype data. Genetics, 176, 1635–1651.

Gardner MG, Fitch AJ, Bertozzi T, Lowe AJ (2011) Rise of the machines - recommendations for ecologists when using next generation sequencing for microsatellite development. Molecular Ecology Resources, 11, 1093–1101.

Glaubitz JC, Rhodes OE, DeWoody JA (2003) Prospects for inferring pairwise relationships with single nucleotide polymorphisms. Molecular Ecology, 12, 1039–1047.

Glover K, Hansen M, Lien S et al. (2010) A comparison of SNP and STR loci for delineating population structure and performing individual genetic assignment. Bmc Genetics, 11, 2.

Guichoux E, Lagache L, Wagner S et al. (2011) Current trends in microsatellite genotyping. Molecular Ecology Resources, 11, 591–611.

Gusfield D (2002) Partition-distance: a problem and class of perfect graphs arising in clustering. Information Processing Letters, 82, 159–164.

Hackett CA, Broadfoot LB (2003) Effects of genotyping errors, missing values and segregation distortion in molecular marker data on the construction of linkage maps. Heredity, 90, 33–38.

Hall N, Mercer L, Phillips D, Shaw J, Anderson AD (2012) Maximum likelihood estimation of individual inbreeding coefficients and null allele frequencies. Genetics Research, 94, 151–161.

Harismendy O, Ng PC, Strausberg RL et al. (2009) Evaluation of next generation sequencing platforms for population targeted sequencing studies. Genome Biol, 10, R32.

Hedrick PW (2005) A standardized genetic differentiation measure. Evolution, 59, 1633–1638.

Hoffman JI, Amos W (2005) Microsatellite genotyping errors: detection approaches, common sources and consequences for paternal exclusion. Molecular Ecology, 14, 599–612.

Holleley C, Geerts P (2009) Multiplex Manager 1.0: a cross-platform computer program that plans and optimizes multiplex PCR. BioTechniques, 46, 511–517.

Huelsenbeck JP, Andolfatto P (2007) Inference of population structure under a Dirichlet process model. Genetics, 175, 1787–1802.

Kalinowski ST (2002) How many alleles per locus should be used to estimate genetic distances? Heredity, 88, 62–65.

Kalinowski ST, Taper ML, Marshall TC (2007) Revising how the computer program cervus accommodates genotyping error increases success in paternity assignment. Molecular Ecology, 16, 1099–1106.

Levinson G, Gutman GA (1987) Slipped-strand mispairing: a major mechanism for DNA sequence evolution. Molecular biology and evolution, 4, 203–221.

Liu N, Chen L, Wang S, Oh C, Zhao H (2005) Comparison of single-nucleotide polymorphisms and microsatellites in inference of population structure. BMC Genetics, 6, S26.

Lynch M (1990) The similarity index and DNA fingerprinting. Molecular Biology and Evolution, 7, 478–484.

Miller JW, Harrison MT (2014) Inconsistency of Pitman-Yor process mixtures for the number of components. Journal of Machine Learning Research, 15, 3333–3370.

Murray V, Monchawin C, England PR (1993) The determination of the sequences present in the shadow bands of a dinucleotide repeat PCR. Nucleic Acids Research, 21, 2395–2398.

Nielsen R, Wakeley J (2001) Distinguishing migration from isolation: a Markov chain Monte Carlo approach. Genetics, 158, 885–896.

O'Connell M, Danzmann RG, Cornuet J-M, Wright J, Ferguson MM (1997) Differentiation of rainbow trout (Oncorhynchus mykiss) populations in Lake Ontario and the evaluation of the stepwise mutation and infinite allele mutation models using microsatellite variability. Canadian Journal of Fisheries and Aquatic Sciences, 54, 1391–1399.

Olejniczak M, Krzyzosiak WJ (2006) Genotyping of simple sequence repeat factors implicated in shadow band generation revisited. Electrophoresis, 27, 3724–3734.

Paetkau D, Waits LP, Clarkson PL, Craighead L, Strobeck C (1997) An empirical evaluation of genetic distance statistics using microsatellite data from bear (Ursidae) populations. Genetics, 143, 1943–1957.

Peel D, Waples RS, Macbeth GM, Do C, Ovenden JR (2013) Accounting for missing data in the estimation of contemporary genetic effective population size ($N_e$). Molecular Ecology Resources, 13, 243–253.

Pompanon F, Bonin A, Bellemain E, Taberlet P (2005) Genotyping errors: causes, consequences and solutions. Nature Reviews Genetics, 6, 847–846.

Pritchard JK, Stephens M, Donnelly P (2000) Inference of population structure using multilocus genotype data. Genetics, 155, 945–959.

Raymond M, Rousset F (1995) GENEPOP (version 1.2): population genetics software for exact tests and ecumenicism. Journal of Heredity, 86, 248–249.

Reeves PA, Panella LW, Richards CM (2012) Retention of agronomically important variation in germplasm core collections: implications for allele mining. Theoretical and Applied Genetics, 124, 1155–1171.

Reeves PA, Richards CM (2009) Accurate inference of subtle population structure (and other genetic discontinuities) using principal coordinates. PLoS One, 4.

Rohlf FJ (2008) NTSYSpc: Numerical Taxonomy System, ver. 2.11x. Exeter, Setauket, NY.

Seddon JM, Parker HG, Ostrander EA, Ellegren H (2005) SNPs in ecological and conservation studies: a test in the Scandinavian wolf population. Molecular Ecology, 14, 503–511.

Sefc KM, Payne RB, Sorenson MD (2003) Microsatellite amplification from museum feather samples: effects of fragment size and template concentration on genotyping errors. The Auk, 120, 982.

Shinde D, Lai Y, Sun F, Arnheim N (2003) Taq DNA polymerase slippage mutation rates measured by PCR and quasi-likelihood analysis: (CA/GT)n and (A/T)n microsatellites. Nucleic Acids Research, 31, 974–980.

Skotte L, Korneliussen TS, Albrechtsen A (2013) Estimating individual admixture proportions from next generation sequencing data. Genetics, 195, 693–702.

Sunnucks P (2000) Efficient genetic markers for population biology. Trends in Ecology & Evolution, 15, 199–203.

Taberlet P, Griffin S, Goossens B et al. (1996) Reliable genotyping of samples with very low DNA quantities using PCR. Nucleic Acids Research, 24, 3189–3194.

Van Oosterhout C, Hutchinson WF, Wills DPM, Shipley P (2004) MICRO-CHECKER: software for identifying and correcting genotyping errors in microsatellite data. Molecular Ecology Notes, 4, 535–538.

Vigouroux Y, Jaqueth JS, Matsuoka Y et al. (2002) Rate and pattern of mutation at microsatellite loci in maize. Molecular Biology and Evolution, 19, 1251–1260.

Wahlund S (1928) Zusammensetzung von Populationen und Korrelationserscheinungen vom Standpunkt der Vererbungslehre aus betrachtet. Hereditas, 11, 65–106.

652 Wang J, Santure AW (2009) Parentage and sibship inference from multilocus genotype data
653        under polygamy. Genetics, 181, 1579–1594.
654 Wattier R, Engel CR, Saumitou-Laprade P, Valero M (1998) Short allele dominance as a source
655        of heterozygote deficiency at microsatellite loci: experimental evidence at the
656        dinucleotide locus Gv1CT in Gracilaria gracilis (Rhodophyta). Molecular Ecology, 7,
657        1569–1573.
658

659                              DATA ACCESSIBILITY
660 Computer programs used to simulate missing data and genotyping error are available from the
661 authors.
662

663                              AUTHOR CONTRIBUTIONS
664 Designed research, all authors.  Ran analyses, all authors.  Analyzed results, PAR.  Wrote
665 computer programs, PAR.  Wrote the manuscript, all authors.
666

667                              FIGURE LEGENDS
668 Figure 1.  Some complex population structures generated by coalescent simulation, visualized
669 using principal coordinate analysis.  Top row, high population subdivision ($G'_{st} = 0.99$); middle
670 row, intermediate ($G'_{st} = 0.5$); bottom row, low ($G'_{st} = 0.01$).
671

672 Figure 2.  Accuracy of population structure inference for unmodified data sets.  X-axis measures
673 population subdivision, Y-axis shows the proportion of correctly inferred populations out of 50
674 possible (black), and the false discovery rate (grey).  a) INSTRUCT no admixture; b)
675 INSTRUCT admixture; c) INSTRUCT inbreeding; d) STRUCTURAMA; e) Neighbor-joining;
676 f) PCOMC.  False discovery rate not calculable for neighbor-joining.
677

678 Figure 3.  Effect of missing data and error on clustering accuracy using the performance ratio.
679 a,g) INSTRUCT no admixture; b,h) INSTRUCT admixture; c,i) INSTRUCT inbreeding; d,j)
680 STRUCTURAMA; e,k) Neighbor-joining; f,l) PCOMC.  For visual clarity, points are mean
681 values binned with an interval width of 0.01.  Error bars indicate SEM.  Statistics were
682 calculated without binning.  Dotted curves mark the 95% confidence interval on the slope of the
683 regression.  Grey shaded area is the 95% confidence interval on the slope of the regression for
684 the false discovery rate ratio.  False discovery rate not calculable for neighbor-joining.  The
685 range of values in f) is truncated because principal coordinate analysis can accept limited missing
686 data.
687

688 Figure 4. Effect of error on clustering accuracy using the partition distance ratio. a) INSTRUCT
689 no admixture; b) INSTRUCT admixture; c) INSTRUCT inbreeding; d) STRUCTURAMA.
690 Dotted curves mark the 95% confidence interval on the slope of the regression. For clarity, a
691 total of 8 outlying points have been omitted from the plots, but were included in the regression.
692

693 Figure 5. a) As genotyping error increases the optimal value for K is elevated under a Dirichlet
694 process prior using STRUCTURAMA.  b) The level of inferred admixture, as measured by the
695 Shannon index, increases with error for 'INSTRUCT admixture'.  c) The upward bias in the
696 admixture estimate is not corrected by using 'INSTRUCT inbreeding'.  d) Error causes the
697 likelihood to be elevated for higher values of K, plateauing after K > 10.  This effect is reversed

698    by using a model that compensates for inbreeding.  Black line, 'INSTRUCT admixture'; grey
699    line, 'INSTRUCT inbreeding'.
700
701

702

703                                             TABLES

704    Table 1. Regression analysis of clustering accuracy, measured with the performance ratio.
705

| Data set | Method | $n_{obs}$ | $m^a$ | $b^b$ | $R^2_{adj}$ | $p^c$ |
|---|---|---|---|---|---|---|
| Missing | | | | | | |
| | INSTRUCT no admixture | 188 | -1.96 | 1.02 | 0.11 | 2.0E-06*** |
| | INSTRUCT admixture | 133 | -1.97 | 0.95 | 0.08 | 6.6E-04** |
| | INSTRUCT inbreeding | 129 | -1.84 | 0.87 | 0.04 | 1.8E-02 |
| | STRUCTURAMA | 106 | -3.46 | 0.96 | 0.27 | 6.4E-09*** |
| | Neighbor-joining | 473 | -3.13 | 0.84 | 0.24 | 1.1E-30*** |
| | PCOMC | 63 | -3.19 | 0.84 | 8.7E-03 | 0.22 |
| Error | | | | | | |
| | INSTRUCT no admixture | 190 | 0.41 | 0.90 | 0.01 | 8.7E-02 |
| | INSTRUCT admixture | 134 | 1.19 | 0.84 | 0.02 | 3.8E-02 |
| | INSTRUCT inbreeding | 130 | 1.36 | 0.91 | 9.5E-03 | 0.14 |
| | STRUCTURAMA | 128 | 1.05 | 0.79 | 0.11 | 8.7E-05*** |
| | Neighbor-joining | 473 | -1.86 | 0.96 | 0.13 | 1.5E-16*** |
| | PCOMC | 185 | -1.67 | 0.95 | 0.04 | 6.2E-03* |

706    [a]Slope with Y-intercept fixed at 1.
707    [b]Y intercept for natural regression.
708    [c]*significant at Holm-Bonferroni corrected α=0.05; ** α=0.01; *** α=0.001.

28

709

710 Table 2. Multiple regression analysis and linear model likelihoods.

711

| Data set | Method | Model | $R^2_{adj}$ | $\beta$ | SE | p[a] | K | AIC | $\Delta_i$[b] | $w_i$[c] |
|---|---|---|---|---|---|---|---|---|---|---|
| Missing | | | | | | | | | | |
| | INSTRUCT no admixture | | | | | | | | | |
| | | % missing { $\beta_0+\beta_1(x_1)$ } | 0.11 | -1.99 | 0.41 | 3.2E-6*** | 2 | -301.76 | 1.01 | 0.38 |
| | | Clumping parameter { $\beta_0+\beta_2(x_2)$ } | 0.02 | 0.02 | 0.01 | 0.09 | 2 | -282.69 | 20.08 | 2.7E-05 |
| | | Combined { $\beta_0+\beta_1(x_1)+\beta_2(x_2)$ } | 0.12 | | | 2.9E-6*** | 3 | -302.77 | 0 | 0.62 |
| | INSTRUCT admixture | | | | | | | | | |
| | | % missing { $\beta_0+\beta_1(x_1)$ } | 0.08 | -1.51 | 0.48 | 1.8E-3* | 2 | -216.92 | 2.24 | 0.24 |
| | | Clumping parameter { $\beta_0+\beta_2(x_2)$ } | 0.04 | 0.03 | 0.01 | 0.04 | 2 | -211.19 | 7.97 | 0.01 |
| | | Combined { $\beta_0+\beta_1(x_1)+\beta_2(x_2)$ } | 0.1 | | | 3.9E-4** | 3 | -219.16 | 0 | 0.74 |
| | INSTRUCT inbreeding | | | | | | | | | |
| | | % missing { $\beta_0+\beta_1(x_1)$ } | 0.04 | -0.99 | 0.46 | 0.03 | 2 | -226.42 | 2.24 | 0.20 |
| | | Clumping parameter { $\beta_0+\beta_2(x_2)$ } | 0.03 | 0.02 | 0.01 | 0.04 | 2 | -226.05 | 2.61 | 0.17 |
| | | Combined { $\beta_0+\beta_1(x_1)+\beta_2(x_2)$ } | 0.06 | | | 7.6E-3 | 3 | -228.66 | 0 | 0.63 |
| | STRUCTURAMA | | | | | | | | | |
| | | % missing { $\beta_0+\beta_1(x_1)$ } | 0.27 | -3.30 | 0.51 | 2.6E-9*** | 2 | -184.27 | 2.80 | 0.20 |
| | | Clumping parameter { $\beta_0+\beta_2(x_2)$ } | 0.02 | -0.03 | 0.01 | 0.03 | 2 | -152.4 | 34.67 | 2.4E-08 |
| | | Combined { $\beta_0+\beta_1(x_1)+\beta_2(x_2)$ } | 0.30 | | | 5.0E-9*** | 3 | -187.07 | 0 | 0.80 |
| | Neighbor-joining | | | | | | | | | |
| | | % missing { $\beta_0+\beta_1(x_1)$ } | 0.24 | -2.2 | 0.18 | 1.4E-30*** | 2 | -1066.74 | 14.94 | 0.00 |
| | | Clumping parameter { $\beta_0+\beta_2(x_2)$ } | 0.04 | 0.02 | 4.7E-3 | 4.1E-5*** | 2 | -950.623 | 131.06 | 3.5E-29 |
| | | Combined { $\beta_0+\beta_1(x_1)+\beta_2(x_2)$ } | 0.27 | | | 3.7E-33*** | 3 | -1081.68 | 0 | 1.00 |
| | PCOMC | | | | | | | | | |
| | | % missing { $\beta_0+\beta_1(x_1)$ } | 8.7E-3 | -1.77 | 1.26 | 0.17 | 2 | -98.52 | 0 | 0.50 |
| | | Clumping parameter { $\beta_0+\beta_2(x_2)$ } | -0.01 | -0.01 | 0.02 | 0.43 | 2 | -97.15 | 1.37 | 0.25 |
| | | Combined { $\beta_0+\beta_1(x_1)+\beta_2(x_2)$ } | 2.8E-3 | | | 0.34 | 3 | -97.19 | 1.33 | 0.25 |
| Error | | | | | | | | | | |
| | INSTRUCT no admixture | | | | | | | | | |
| | | Null percent { $\beta_0+\beta_1(x_1)$ } | -4.5E-3 | 0.14 | 0.86 | 0.87 | 2 | -208.04 | 5.57 | 0.05 |
| | | Drop out percent { $\beta_0+\beta_2(x_2)$ } | 0.02 | 2.07 | 0.88 | 0.02 | 2 | -213.61 | 0 | 0.78 |
| | | Stutter percent { $\beta_0+\beta_3(x_3)$ } | -3.4E-3 | -1.26 | 2.51 | 0.61 | 2 | -208.25 | 5.36 | 0.05 |
| | | Combined { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)$ } | 0.16 | | | 0.12 | 4 | -209.92 | 3.69 | 0.12 |

29

## INSTRUCT admixture

| Model | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Null percent { $\beta_0+\beta_1(x_1)$ } | -2.3E-3 | 0.99 | 1.49 | 0.51 | 2 | -51.1 | 6.06 | 0.03 |
| Drop out percent { $\beta_0+\beta_2(x_2)$ } | 0.04 | 3.90 | 1.44 | 7.6E-3 | 2 | -57.16 | 0 | 0.59 |
| Stutter percent { $\beta_0+\beta_3(x_3)$ } | 5.3E-3 | -6.45 | 4.49 | 0.15 | 2 | -52.12 | 5.04 | 0.05 |
| Combined { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)$ } | 0.05 | | | 0.03 | 4 | -55.98 | 1.18 | 0.33 |

## INSTRUCT inbreeding

| Model | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Null percent { $\beta_0+\beta_1(x_1)$ } | -7.6E-3 | -0.71 | 1.82 | 0.70 | 2 | 0.24 | 11.09 | 3.1E-03 |
| Drop out percent { $\beta_0+\beta_2(x_2)$ } | 0.05 | 5.35 | 1.68 | 1.8E-3* | 2 | -7.91 | 2.94 | 0.18 |
| Stutter percent { $\beta_0+\beta_3(x_3)$ } | 0.03 | -15.29 | 5.82 | 9.7E-3 | 2 | -4.54 | 6.31 | 0.03 |
| Combined { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)$ } | 0.09 | | | 2.1E-3* | 4 | -10.85 | 0 | 0.78 |

## STRUCTURAMA

| Model | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Null percent { $\beta_0+\beta_1(x_1)$ } | 1.0E-3 | 1.26 | 0.83 | 0.13 | 2 | -200.45 | 15.10 | 3.5E-04 |
| Drop out percent { $\beta_0+\beta_2(x_2)$ } | 0.10 | 3.10 | 0.80 | 1.8E-4** | 2 | -214.15 | 1.40 | 0.33 |
| Stutter percent { $\beta_0+\beta_3(x_3)$ } | 0.02 | 5.12 | 2.60 | 0.05 | 2 | -202.97 | 12.58 | 1.2E-03 |
| Combined { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)$ } | 0.13 | | | 2.0E-4** | 4 | -215.55 | 0 | 0.67 |

## Neighbor-joining

| Model | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Null percent { $\beta_0+\beta_1(x_1)$ } | 0.03 | -1.31 | 0.29 | 1.2E-5*** | 2 | -951.524 | 76.95 | 2.0E-17 |
| Drop out percent { $\beta_0+\beta_2(x_2)$ } | 0.08 | -1.67 | 0.26 | 5.7E-10*** | 2 | -975.431 | 53.04 | 3.0E-12 |
| Stutter percent { $\beta_0+\beta_3(x_3)$ } | 0.07 | -6.69 | 0.99 | 4.2E-11*** | 2 | -971.313 | 57.16 | 3.9E-13 |
| Combined { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)$ } | 0.18 | | | 1.2E-20*** | 4 | -1028.47 | 0 | 1.00 |

## PCOMC

| Model | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Null percent { $\beta_0+\beta_1(x_1)$ } | 2.0E-3 | -1.19 | 0.74 | 0.11 | 2 | -274.075 | 5.05 | 0.05 |
| Drop out percent { $\beta_0+\beta_2(x_2)$ } | 0.02 | -1.49 | 0.66 | 0.03 | 2 | -277.417 | 1.71 | 0.26 |
| Stutter percent { $\beta_0+\beta_3(x_3)$ } | 8.5E-3 | -5.15 | 2.45 | 0.04 | 2 | -275.274 | 3.85 | 0.09 |
| Combined { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)$ } | 0.04 | | | 0.02 | 4 | -279.128 | 0 | 0.61 |

[a] *Significant at Holm-Bonferroni corrected α=0.05; ** α=0.01; *** α=0.001.
[b] Rescaled Akaike information criterion (AIC).
[c] Model probability.

716                                SUPPLEMENTARY FIGURE LEGENDS

717 Supplementary Figure 1. Properties of simulated data sets. a) Frequency of data sets with

718 differing levels of population subdivision. b) The amount of missing data introduced into data

719 sets was unrelated to level of population subdivision. c) The frequency of erroneous genotypes

720 was negatively correlated with population subdivision.

721

722 Supplementary Figure 2. The inbreeding-like effect of SSR genotyping error. a) The level of

723 apparent inbreeding increases with genotyping error. b, c) The increase in F is driven by a

724 decrease in observed heterozygosity rather than an increase in expected heterozygosity, which

725 decreased significantly, but slightly. d) Error increases the apparent level of population

726 subdivision, producing a tendency for more accurate assignment of genotypes to clusters as the

727 amount of error increases.

728

729 Supplementary Table 1. Spearman's rank order correlation of clustering accuracy, measured
730 with the partition distance ratio.
731

| Data set | Method | $n_{obs}{}^a$ | $\rho^{a,b}$ | $p^{a,c}$ | $n_{obs}{}^d$ | $\rho^{d,b}$ | $p^{d,c}$ |
|---|---|---|---|---|---|---|---|
| Missing | | | | | | | |
| | INSTRUCT no admixture | 437 | -0.22 | 2.3E-06*** | 188 | 0.68 | 1.4E-26*** |
| | INSTRUCT admixture | 356 | -0.23 | 1.1E-05*** | 133 | 0.63 | 2.6E-16*** |
| | INSTRUCT inbreeding | 357 | -0.11 | 0.03 | 125 | 0.45 | 1.4E-07*** |
| | STRUCTURAMA | 115 | 0.12 | 0.19 | 100 | 0.05 | 0.61 |
| Error | | | | | | | |
| | INSTRUCT no admixture | 464 | 0.30 | 2.4E-11*** | 190 | 0.51 | 1.0E-13*** |
| | INSTRUCT admixture | 393 | 0.49 | 4.5E-25*** | 134 | 0.41 | 7.0E-07*** |
| | INSTRUCT inbreeding | 382 | 0.53 | 6.2E-29*** | 130 | 0.46 | 2.7E-08*** |
| | STRUCTURAMA | 141 | 0.18 | 0.03 | 115 | 0.57 | 4.4E-11*** |

732 [a] Spearman's rank order correlation between percent data set modification and the partition
733 distance ratio metric.
734 [b] Spearman's coefficient, rho.
735 [c] *Significant at Holm-Bonferroni corrected α=0.05; ** α=0.01; *** α=0.001.
736 [d] Spearman's rank order correlation between the performance ratio and the partition distance
737 ratio.

738 Supplementary Table 2. Multiple regression and likelihood analysis of six competing model effects explaining accuracy of
739 model-based clustering methods subject to erroneous data, measured with the performance ratio.
740

| Method | Model | $R^2_{adj}$ | β | SE | $p^c$ | K | AIC | $\Delta_i^d$ | $w_i^e$ |
|---|---|---|---|---|---|---|---|---|---|
| STRUCTURAMA | | | | | | | | | |
| | % error { $\beta_0+\beta_1(x_1)$ } | 0.11 | 3.16 | 0.95 | 1.2E-03* | 2 | -215.01 | 32.01 | 4.4E-08 |
| | Δ APL[a] { $\beta_0+\beta_2(x_2)$ } | 0.03 | 0.09 | 0.07 | 0.20 | 2 | -204.35 | 42.67 | 2.1E-10 |
| | Δ K { $\beta_0+\beta_3(x_3)$ } | 0.24 | 0.10 | 0.02 | 2.6E-08*** | 2 | -235.96 | 11.06 | 1.6E-03 |
| | Δ G'st { $\beta_0+\beta_4(x_4)$ } | -7.8E-03 | -0.86 | 4.14 | 0.84 | 2 | -199.33 | 47.69 | 1.7E-11 |
| | Δ p(1)[b] { $\beta_0+\beta_5(x_5)$ } | 1.7E-03 | -3.48 | 1.70 | 0.04 | 2 | -200.54 | 46.48 | 3.2E-11 |
| | Δ F { $\beta_0+\beta_6(x_6)$ } | 2.9E-03 | -1.16 | 0.72 | 0.11 | 2 | -200.7 | 46.32 | 3.5E-11 |
| | Best 2 { $\beta_0+\beta_3(x_3)+\beta_5(x_5)$ } | 0.28 | | | 6.4E-10*** | 3 | -240.68 | 6.34 | 1.7E-02 |
| | Best 3 { $\beta_0+\beta_1(x_1)+\beta_3(x_3)+\beta_5(x_5)$ } | 0.30 | | | 2.0E-10*** | 4 | -244.76 | 2.26 | 0.13 |
| | Best 4 { $\beta_0+\beta_1(x_1)+\beta_3(x_3)+\beta_4(x_4)+\beta_5(x_5)$ } | 0.32 | | | 1.4E-10*** | 5 | -247.02 | 0 | 0.40 |
| | Best 5 { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_4(x_4)+\beta_5(x_5)$ } | 0.32 | | | 2.8E-10*** | 6 | -246.68 | 0.34 | 0.33 |
| | All { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_4(x_4)+\beta_5(x_5)+\beta_6(x_6)$ } | 0.32 | | | 1.1E-09*** | 7 | -244.73 | 2.29 | 0.13 |
| INSTRUCT, no admixture | | | | | | | | | |
| | % error { $\beta_0+\beta_1(x_1)$ } | 0.01 | 2.62 | 1.12 | 0.02 | 2 | -210.86 | 12.65 | 5.5E-04 |
| | Δ APL[a] { $\beta_0+\beta_2(x_2)$ } | 0.04 | -1.2E-03 | 0.09 | 0.99 | 2 | -216.59 | 6.92 | 9.6E-03 |
| | Δ K { $\beta_0+\beta_3(x_3)$ } | -3.6E-03 | 1.0E-03 | 0.00 | 0.64 | 2 | -208.21 | 15.30 | 1.4E-04 |
| | Δ G'st { $\beta_0+\beta_4(x_4)$ } | 0.03 | -2.47 | 2.00 | 0.22 | 2 | -214.90 | 8.61 | 4.1E-03 |
| | Δ p(1)[b] { $\beta_0+\beta_5(x_5)$ } | 0.02 | -5.02 | 3.62 | 0.17 | 2 | -212.23 | 11.28 | 1.1E-03 |
| | Δ F { $\beta_0+\beta_6(x_6)$ } | 0.02 | -1.03 | 0.78 | 0.19 | 2 | -212.64 | 10.87 | 1.3E-03 |
| | Best 2 { $\beta_0+\beta_1(x_1)+\beta_4(x_4)$ } | 0.06 | | | 6.2E-04* | 3 | -220.91 | 2.60 | 0.08 |
| | Best 3 { $\beta_0+\beta_1(x_1)+\beta_4(x_4)+\beta_6(x_6)$ } | 0.08 | | | 2.9E-04* | 4 | -223.11 | 0.40 | 0.25 |
| | Best 4 { $\beta_0+\beta_1(x_1)+\beta_4(x_4)+\beta_5(x_5)+\beta_6(x_6)$ } | 0.09 | | | 2.9E-04* | 5 | -223.51 | 0 | 0.30 |

| Model | | | | | | | |
|---|---|---|---|---|---|---|---|
| Best 5 $\{ \beta_0+\beta_1(x_1)+\beta_3(x_3)+\beta_4(x_4)+\beta_5(x_5)+\beta_6(x_6) \}$ | 0.09 | | | 4.0E-04* | 6 | -223.08 | 0.44 | 0.25 |
| All $\{ \beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_4(x_4)+\beta_5(x_5)+\beta_6(x_6) \}$ | 0.09 | | | 8.7E-04* | 7 | -221.30 | 2.21 | 0.10 |

INSTRUCT, admixture

| Model | | | | | | | |
|---|---|---|---|---|---|---|---|
| % error $\{ \beta_0+\beta_1(x_1) \}$ | 0.02 | 5.18 | 2.71 | 0.06 | 2 | -54.77 | 0.22 | 0.19 |
| $\Delta$ APL[a] $\{ \beta_0+\beta_2(x_2) \}$ | 0.02 | 0.04 | 0.16 | 0.78 | 2 | -53.84 | 1.15 | 0.12 |
| $\Delta$ K $\{ \beta_0+\beta_3(x_3) \}$ | -7.6E-03 | -1.9E-04 | 7.7E-03 | 0.98 | 2 | -50.41 | 4.59 | 0.02 |
| $\Delta$ G'st $\{ \beta_0+\beta_4(x_4) \}$ | -3.0E-03 | 8.87 | 7.02 | 0.21 | 2 | -51.02 | 3.97 | 0.03 |
| $\Delta$ p(1)[b] $\{ \beta_0+\beta_5(x_5) \}$ | -4.2E-03 | -6.48 | 4.27 | 0.13 | 2 | -50.85 | 4.14 | 0.03 |
| $\Delta$ F $\{ \beta_0+\beta_6(x_6) \}$ | -4.6E-03 | -3.32 | 1.82 | 0.07 | 2 | -50.80 | 4.19 | 0.03 |
| Best 2 $\{ \beta_0+\beta_1(x_1)+\beta_6(x_6) \}$ | 0.03 | | | 0.04 | 3 | -54.99 | 0 | 0.21 |
| Best 3 $\{ \beta_0+\beta_1(x_1)+\beta_5(x_5)+\beta_6(x_6) \}$ | 0.04 | | | 0.05 | 4 | -54.43 | 0.56 | 0.16 |
| Best 4 $\{ \beta_0+\beta_1(x_1)+\beta_4(x_4)+\beta_5(x_5)+\beta_6(x_6) \}$ | 0.04 | | | 0.05 | 5 | -54.28 | 0.71 | 0.15 |
| Best 5 $\{ \beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_4(x_4)+\beta_5(x_5)+\beta_6(x_6) \}$ | 0.04 | | | 0.09 | 6 | -52.36 | 2.63 | 0.06 |
| All $\{ \beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_4(x_4)+\beta_5(x_5)+\beta_6(x_6) \}$ | 0.03 | | | 0.14 | 7 | -50.36 | 4.63 | 0.02 |

INSTRUCT, inbreeding

| Model | | | | | | | |
|---|---|---|---|---|---|---|---|
| % error $\{ \beta_0+\beta_1(x_1) \}$ | 9.5E-03 | 3.44 | 3.58 | 0.34 | 2 | -1.99 | 0.44 | 0.16 |
| $\Delta$ APL[a] $\{ \beta_0+\beta_2(x_2) \}$ | 9.2E-03 | -0.05 | 0.21 | 0.80 | 2 | -1.95 | 0.48 | 0.15 |
| $\Delta$ K $\{ \beta_0+\beta_3(x_3) \}$ | -3.0E-03 | 0.01 | 0.01 | 0.33 | 2 | -0.37 | 2.06 | 0.07 |
| $\Delta$ G'st $\{ \beta_0+\beta_4(x_4) \}$ | -7.3E-03 | 2.68 | 8.98 | 0.77 | 2 | 0.19 | 2.62 | 0.05 |
| $\Delta$ p(1)[b] $\{ \beta_0+\beta_5(x_5) \}$ | 2.5E-03 | -8.53 | 5.50 | 0.12 | 2 | -1.07 | 1.36 | 0.10 |
| $\Delta$ F $\{ \beta_0+\beta_6(x_6) \}$ | -5.1E-03 | -1.58 | 2.41 | 0.51 | 2 | -0.10 | 2.33 | 0.06 |
| Best 2 $\{ \beta_0+\beta_1(x_1)+\beta_5(x_5) \}$ | 0.02 | | | 0.10 | 3 | -2.43 | 0 | 0.19 |
| Best 3 $\{ \beta_0+\beta_1(x_1)+\beta_5(x_5)+\beta_6(x_6) \}$ | 0.02 | | | 0.14 | 4 | -1.38 | 1.05 | 0.11 |
| Best 4 $\{ \beta_0+\beta_1(x_1)+\beta_3(x_3)+\beta_5(x_5)+\beta_6(x_6) \}$ | 0.02 | | | 0.17 | 5 | -0.34 | 2.09 | 0.07 |
| Best 5 $\{ \beta_0+\beta_1(x_1)+\beta_3(x_3)+\beta_4(x_4)+\beta_5(x_5)+\beta_6(x_6) \}$ | 0.01 | | | 0.27 | 6 | 1.59 | 4.02 | 0.03 |
| All $\{ \beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_4(x_4)+\beta_5(x_5)+\beta_6(x_6) \}$ | 4.2E-03 | | | 0.37 | 7 | 3.52 | 5.95 | 9.9E-03 |

741 [a]Mean number of alleles per locus

742 [b]Mean private allele frequency

743 [c]*significant at Holm-Bonferroni corrected α=0.05; ** α=0.01; *** α=0.001.

744 [d]Rescaled AIC.

745 [e]Model probability.

746

747

748 Supplementary Table 3. Multiple regression and likelihood analysis of six competing model effects explaining accuracy of

749 model-based clustering methods subject to erroneous data, measured with the partition distance ratio metric.

750

| Method | Model | $R^2_{adj}$ | β | SE | $p^c$ | K | AIC | $\Delta_i^d$ | $w_i^e$ |
|---|---|---|---|---|---|---|---|---|---|
| STRUCTURAMA | | | | | | | | | |
| | % error { $\beta_0+\beta_1(x_1)$ } | 0.04 | -0.78 | 0.72 | 0.28 | 2 | -245.03 | 83.80 | 2.4E-19 |
| | Δ APL[a] { $\beta_0+\beta_2(x_2)$ } | 6.8E-03 | -0.01 | 0.05 | 0.84 | 2 | -238.05 | 90.78 | 7.3E-21 |
| | Δ K { $\beta_0+\beta_3(x_3)$ } | 0.43 | 0.12 | 0.01 | 1.2E-15*** | 2 | -318.23 | 10.60 | 1.9E-03 |
| | Δ G'st { $\beta_0+\beta_4(x_4)$ } | 0.05 | -1.94 | 3.15 | 0.54 | 2 | -246.64 | 82.19 | 5.4E-19 |
| | Δ p(1)[b] { $\beta_0+\beta_5(x_5)$ } | 0.01 | -0.13 | 1.37 | 0.92 | 2 | -240.54 | 88.29 | 2.6E-20 |
| | Δ F { $\beta_0+\beta_6(x_6)$ } | 0.15 | 1.51 | 0.56 | 7.6E-03 | 2 | -261.90 | 66.93 | 1.1E-15 |
| | Best 2 { $\beta_0+\beta_3(x_3)+\beta_6(x_6)$ } | 0.48 | | | 1.9E-20*** | 3 | -328.83 | 0 | 0.38 |
| | Best 3 { $\beta_0+\beta_1(x_1)+\beta_3(x_3)+\beta_6(x_6)$ } | 0.48 | | | 6.8E-20*** | 4 | -328.75 | 0.07 | 0.37 |
| | Best 4 { $\beta_0+\beta_1(x_1)+\beta_3(x_3)+\beta_4(x_4)+\beta_6(x_6)$ } | 0.48 | | | 4.0E-19*** | 5 | -327.20 | 1.63 | 0.17 |
| | Best 5 { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_4(x_4)+\beta_6(x_6)$ } | 0.47 | | | 2.4E-18*** | 6 | -325.24 | 3.59 | 0.06 |
| | All { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_4(x_4)+\beta_5(x_5)+\beta_6(x_6)$ } | 0.47 | | | 1.3E-17*** | 7 | -323.25 | 5.58 | 0.02 |
| INSTRUCT, no admixture | | | | | | | | | |
| | % error { $\beta_0+\beta_1(x_1)$ } | 0.06 | 1.32 | 0.52 | 0.01 | 2 | -597.36 | 35.52 | 8.0E-09 |
| | Δ APL[a] { $\beta_0+\beta_2(x_2)$ } | -1.5E-03 | 0.10 | 0.05 | 0.07 | 2 | -567.72 | 65.17 | 2.9E-15 |
| | Δ K { $\beta_0+\beta_3(x_3)$ } | -9.0E-04 | -1.3E-03 | 1.4E-03 | 0.36 | 2 | -567.98 | 64.91 | 3.3E-15 |
| | Δ G'st { $\beta_0+\beta_4(x_4)$ } | 0.10 | 2.70 | 1.45 | 0.06 | 2 | -616.39 | 16.50 | 1.1E-04 |
| | Δ p(1)[b] { $\beta_0+\beta_5(x_5)$ } | 7.7E-03 | -1.01 | 1.21 | 0.40 | 2 | -571.96 | 60.93 | 2.4E-14 |

35

| Model | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| $\Delta$ F { $\beta_0+\beta_6(x_6)$ } | 0.11 | 0.67 | 0.24 | 0.00 | 2 | -621.13 | 11.76 | 1.2E-03 |
| Best 2 { $\beta_0+\beta_3(x_3)+\beta_4(x_4)$ } | 0.13 | | | 1.8E-14*** | 3 | -629.15 | 3.73 | 6.4E-02 |
| Best 3 { $\beta_0+\beta_3(x_3)+\beta_4(x_4)+\beta_6(x_6)$ } | 0.13 | | | 1.8E-14*** | 4 | -630.95 | 1.93 | 0.16 |
| Best 4 { $\beta_0+\beta_1(x_1)+\beta_3(x_3)+\beta_4(x_4)+\beta_6(x_6)$ } | 0.14 | | | 1.5E-14*** | 5 | -632.89 | 0.00 | 0.41 |
| Best 5 { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_4(x_4)+\beta_6(x_6)$ } | 0.14 | | | 4.5E-14*** | 6 | -631.80 | 1.084 | 0.24 |
| All { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_4(x_4)+\beta_5(x_5)+\beta_6(x_6)$ } | 0.14 | | | 1.4E-13*** | 7 | -630.51 | 2.38 | 0.13 |
| **INSTRUCT, admixture** | | | | | | | | |
| % error { $\beta_0+\beta_1(x_1)$ } | 0.12 | 5.79 | 1.04 | 4.7E-08*** | 2 | -45.33 | 6.06 | 1.5E-02 |
| $\Delta$ APL[a] { $\beta_0+\beta_2(x_2)$ } | 8.0E-03 | 0.24 | 0.11 | 0.03 | 2 | 0.16 | 51.54 | 2.0E-12 |
| $\Delta$ K { $\beta_0+\beta_3(x_3)$ } | -1.5E-03 | -3.3E-04 | 5.0E-03 | 0.95 | 2 | 3.85 | 55.24 | 3.1E-13 |
| $\Delta$ G'st { $\beta_0+\beta_4(x_4)$ } | 0.05 | 6.24 | 2.93 | 0.03 | 2 | -16.00 | 35.39 | 6.4E-09 |
| $\Delta$ p(1)[b] { $\beta_0+\beta_5(x_5)$ } | -6.0E-04 | -3.42 | 2.35 | 0.15 | 2 | 3.50 | 54.88 | 3.7E-13 |
| $\Delta$ F { $\beta_0+\beta_6(x_6)$ } | 4.7E-03 | -1.21 | 0.49 | 0.01 | 2 | 1.43 | 52.81 | 1.1E-12 |
| Best 2 { $\beta_0+\beta_1(x_1)+\beta_3(x_3)$ } | 0.13 | | | 1.8E-12*** | 3 | -48.28 | 3.11 | 0.07 |
| Best 3 { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)$ } | 0.13 | | | 1.3E-12*** | 4 | -50.62 | 0.7663 | 0.21 |
| Best 4 { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_4(x_4)$ } | 0.14 | | | 1.9E-12*** | 5 | -51.22 | 0.16 | 0.29 |
| Best 5 { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_4(x_4)+\beta_5(x_5)$ } | 0.14 | | | 3.0E-12*** | 6 | -51.39 | 0.00 | 0.31 |
| All { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_4(x_4)+\beta_5(x_5)+\beta_6(x_6)$ } | 0.14 | | | 1.2E-11*** | 7 | -49.39 | 2.00 | 0.11 |
| **INSTRUCT, inbreeding** | | | | | | | | |
| % error { $\beta_0+\beta_1(x_1)$ } | 0.12 | 6.89 | 1.06 | 3.2E-10*** | 2 | -19.01 | 15.98 | 1.8E-04 |
| $\Delta$ APL[a] { $\beta_0+\beta_2(x_2)$ } | 4.6E-03 | 0.34 | 0.11 | 2.3E-03* | 2 | 26.42 | 61.40 | 2.5E-14 |
| $\Delta$ K { $\beta_0+\beta_3(x_3)$ } | -2.3E-03 | -1.9E-04 | 5.5E-03 | 0.97 | 2 | 29.02 | 64.01 | 6.7E-15 |
| $\Delta$ G'st { $\beta_0+\beta_4(x_4)$ } | 0.03 | 6.35 | 3.01 | 0.04 | 2 | 15.03 | 50.02 | 7.3E-12 |
| $\Delta$ p(1)[b] { $\beta_0+\beta_5(x_5)$ } | -2.5E-03 | -5.82 | 2.41 | 0.02 | 2 | 29.08 | 64.07 | 6.5E-15 |
| $\Delta$ F { $\beta_0+\beta_6(x_6)$ } | -5.0E-04 | -1.61 | 0.52 | 2.2E-03* | 2 | 28.32 | 63.31 | 9.5E-15 |
| Best 2 { $\beta_0+\beta_3(x_3)+\beta_6(x_6)$ } | 0.13 | | | 1.4E-12*** | 3 | -23.86 | 11.13 | 0.00 |

| | | | | | | |
|---|---|---|---|---|---|---|
| Best 3 { $\beta_0+\beta_1(x_1)+\beta_3(x_3)+\beta_6(x_6)$ } | 0.15 | 6.1E-14*** | 4 | -32.04 | 2.95 | 0.12 |
| Best 4 { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_6(x_6)$ } | 0.15 | 1.0E-13*** | 5 | -32.47 | 2.52 | 0.15 |
| Best 5 { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_5(x_5)+\beta_6(x_6)$ } | 0.16 | 5.6E-14*** | 6 | -34.99 | 0 | 0.53 |
| All { $\beta_0+\beta_1(x_1)+\beta_2(x_2)+\beta_3(x_3)+\beta_4(x_4)+\beta_5(x_5)+\beta_6(x_6)$ } | 0.16 | 2.3E-13*** | 7 | -32.99 | 2.00 | 0.20 |

751 [a] Mean number of alleles per locus
752 [b] Mean private allele frequency
753 [c] *significant at Holm-Bonferroni corrected α=0.05; ** α=0.01; *** α=0.001.
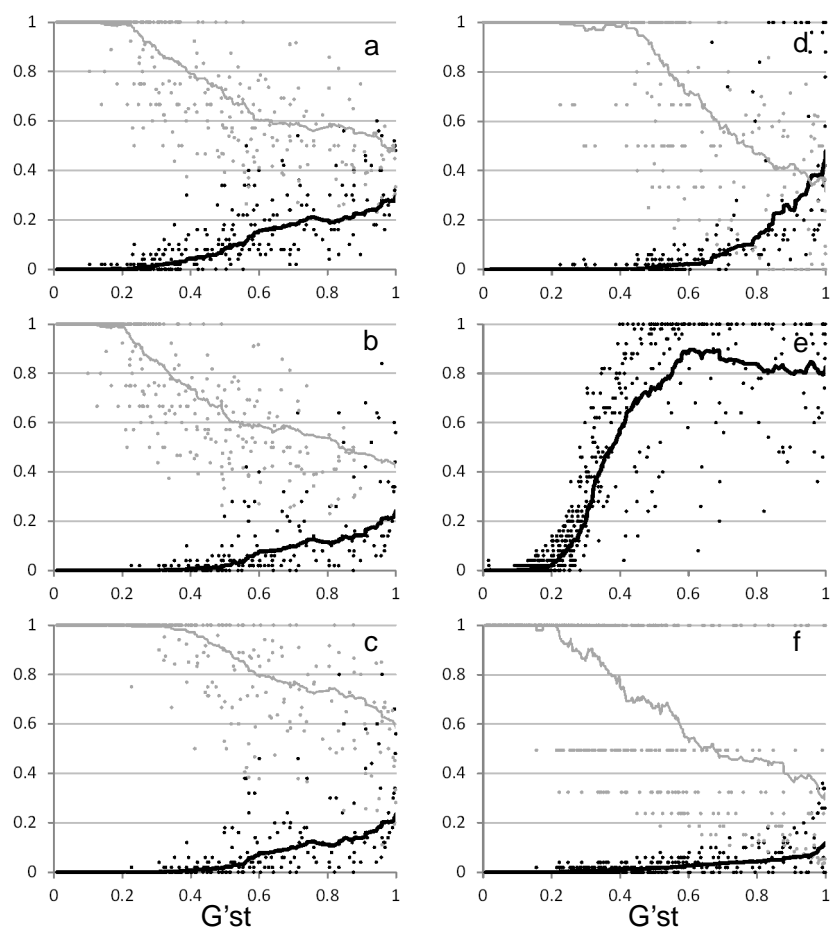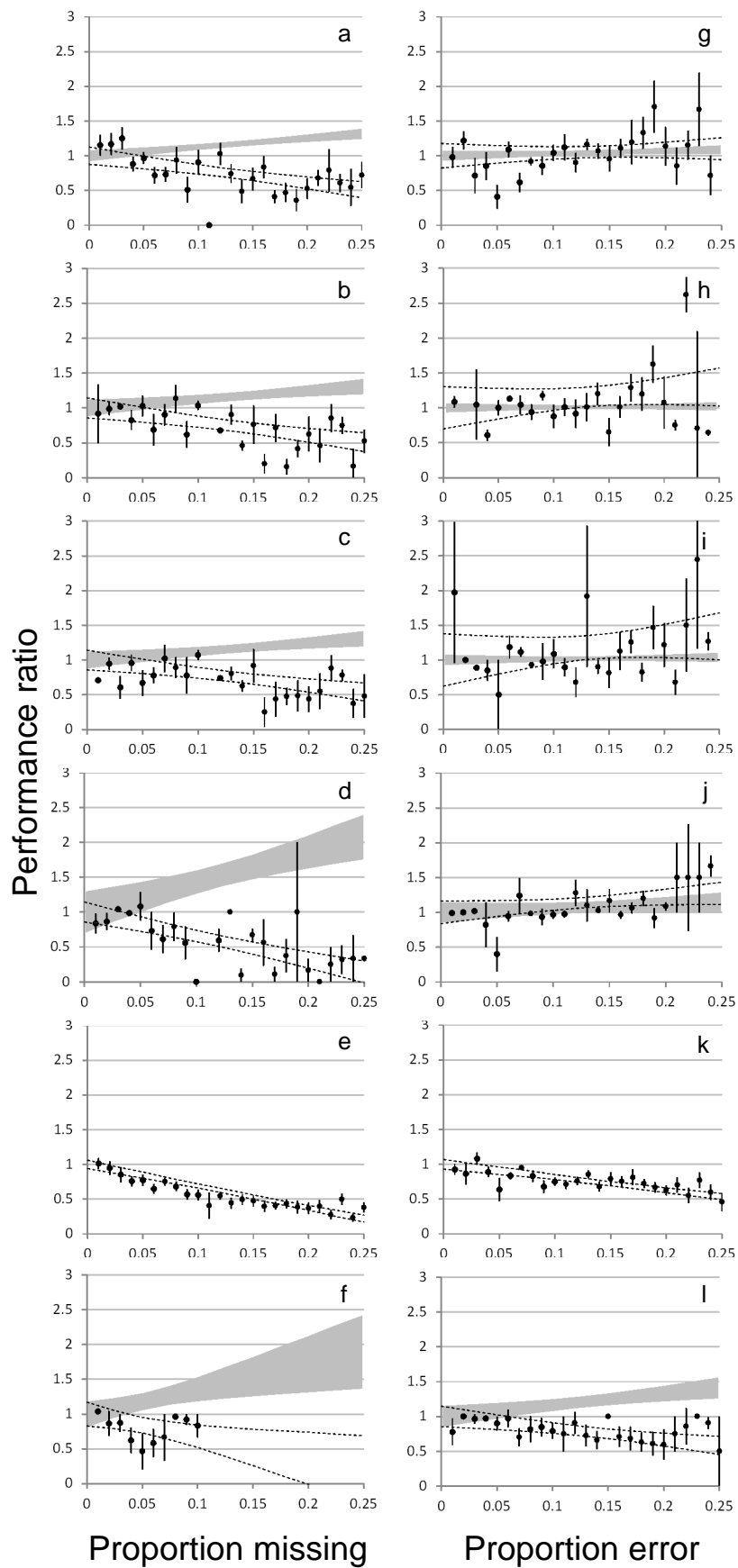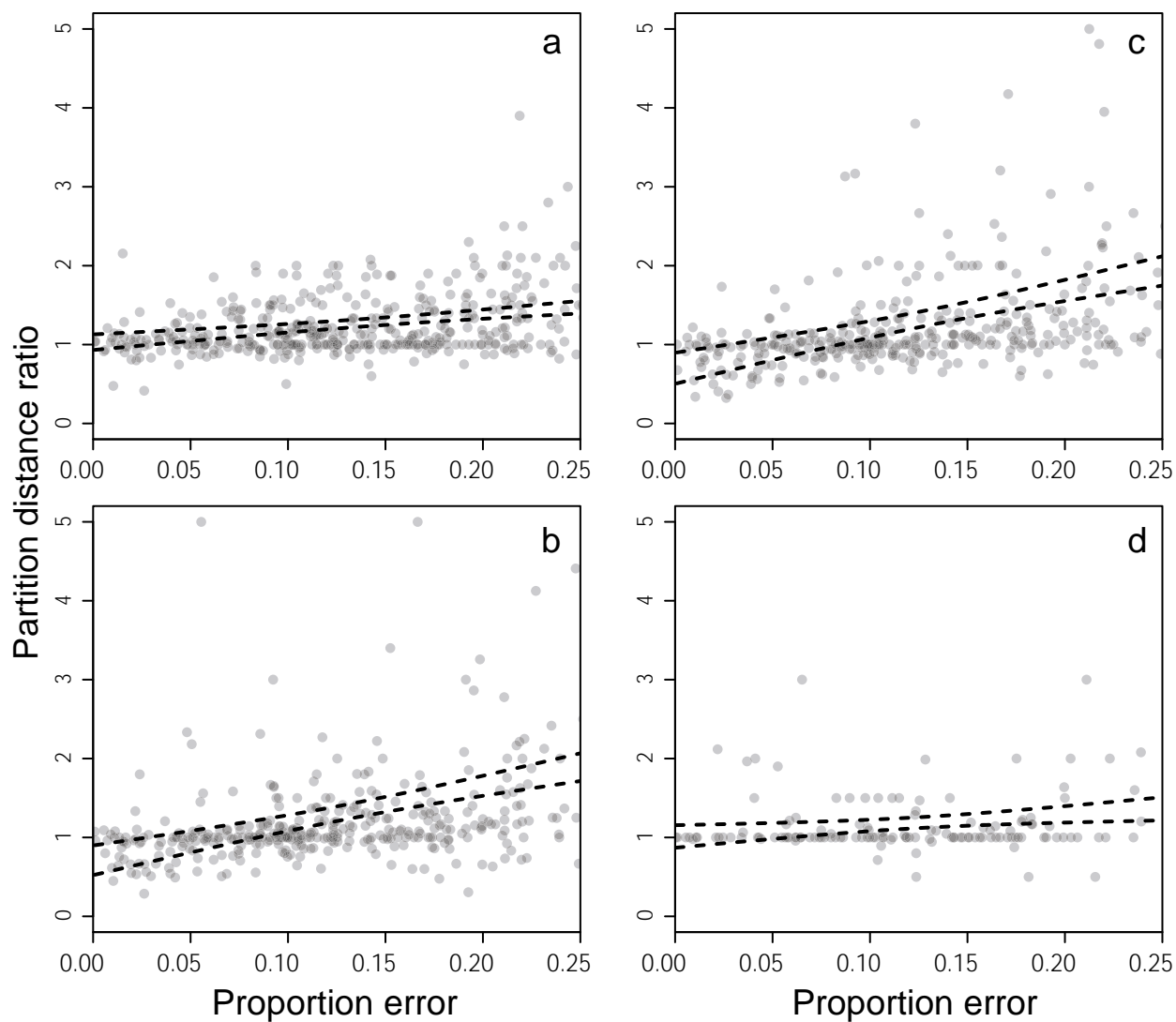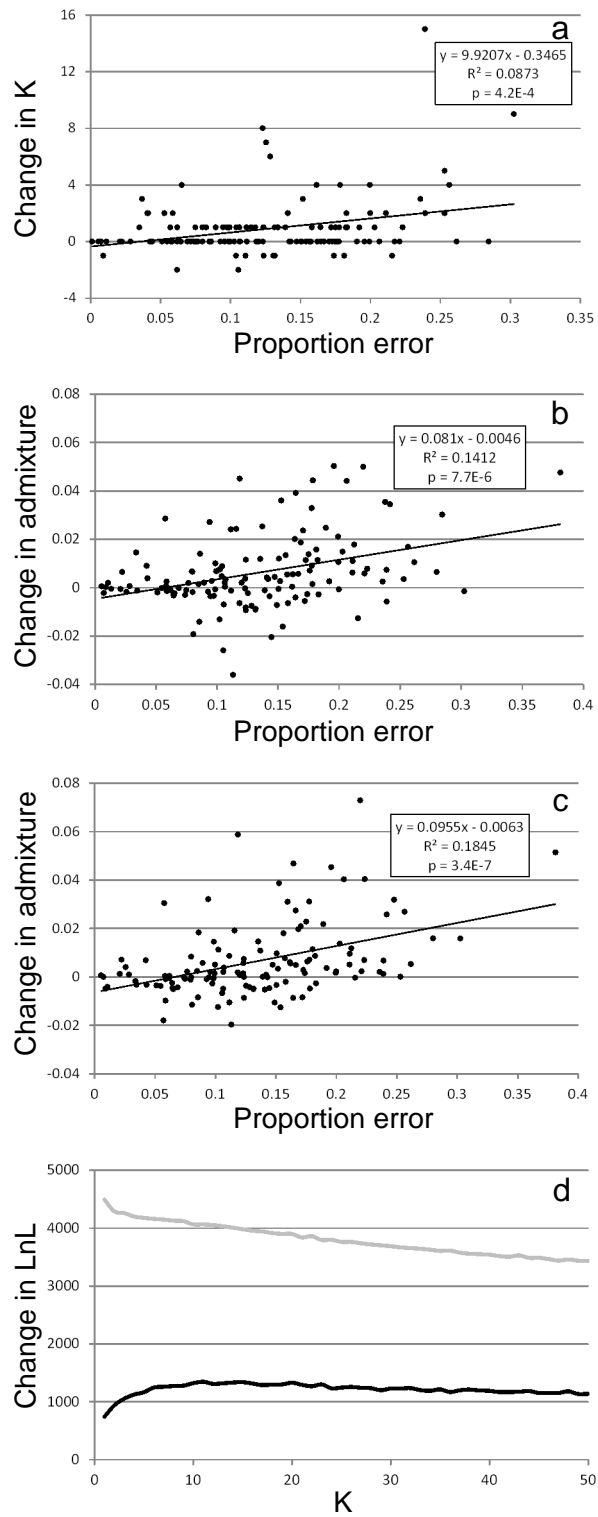754 [d] Rescaled AIC.
755 [e] Model probability.
756

Figure 1.

Figure 2.

Figure 3.

Figure 4.

Figure 5.