

1 A parsimony estimator of the number of populations from a
2 STRUCTURE-like analysis

3

4 Jinliang Wang

5 *Institute of Zoology, Zoological Society of London, London NW1 4RY, United Kingdom*

6

7 *Left running head:* J Wang

8 *Right running head:* Estimating the number of populations

9 *Key words:* STRUCTURE, markers, genetic differentiation, number of populations

10 *Corresponding author:*

11 Jinliang Wang

12 Institute of Zoology

13 Regent's Park

14 London NW1 4RY

15 United Kingdom

16 Tel: 0044 20 74496620

17 Fax: 0044 20 75862870

18 Email: jinliang.wang@ioz.ac.uk

19

20 **Abstract**

21 Population genetics model based Bayesian methods have been proposed and widely applied
22 to making unsupervised inference of population structure from a sample of multilocus
23 genotypes. Usually they provide good estimates of the ancestry (or population membership)
24 of sampled individuals by clustering them probabilistically or proportionally into (anonymous)
25 populations. However, they have difficulties in accurately estimating the number of
26 populations (K) represented by the sampled individuals. This study proposed a new *ad hoc*
27 estimator of K , calculable from the output of a population clustering program such as
28 STRUCTURE or ADMIXTURE. The new criterion, called parsimony index (PI), aims to
29 identify the number of populations (K) which yields consistently the minimal admixture
30 estimates of sampled individuals. Extensive simulated and empirical data were used to
31 compare the accuracy of PI and two popular K estimators based on $\Pr[X|K]$ (i.e. the
32 probability of genotype data X given K) and ΔK (i.e. the rate of change of the probability of
33 data as a function of K) calculated from STRUCTURE outputs, and the accuracy of PI and
34 the cross-validation method calculated from ADMIXTURE outputs. It was shown that PI was
35 more accurate than the other methods consistently in various population structure (e.g.
36 hierarchical island model, different extents of differentiation) and sampling (e.g. unbalanced
37 sample sizes, different marker information contents) scenarios. The ΔK method was more
38 accurate than the $\Pr[X|K]$ method only for hierarchically structured or highly inbred
39 populations, and the opposite was true in the other scenarios. The PI method was
40 implemented in a computer program, KFinder, which can be run on all major computer
41 platforms.

42

43 **Introduction**

44 Traditional population structure analysis uses Wright's (1951) F statistics to describe and
45 understand the patterns of genetic variation in populations. These statistics, F_{IS} , F_{ST} and F_{IT} ,
46 can be estimated from genetic marker data collected from several populations (e.g. Weir and
47 Cockerham 1984), and offer a convenient and elegant means of summarising population
48 structures. However, the analysis of F statistics relies on information of known predefined
49 (e.g. by geographic locations) source populations of sampled individuals. In reality, however,
50 the information might be unavailable, incomplete, or unreliable for sorting individuals into
51 populations. A typical example is mixed stock analysis (Smouse *et al.* 1990), where

52 individuals coming from different source populations to mix in the same feeding/breeding
53 ground are sampled to determine the genetic structure. Another typical example is a batch of
54 seized illegally traded animals or animal parts (e.g. tusks). In both examples, the source
55 populations of the sampled individuals are unknown and are the primary interest of analysis
56 (Hsieh *et al.* 2003; Velo-Anton *et al.* 2007).

57 Pritchard *et al.* (2000) proposed a Bayesian method, based on a population genetics
58 model, to identify populations represented by a sample of individuals and to assign the
59 individuals (or their genomes) probabilistically to the identified source populations using
60 individual genotypic data. This is possible because population genetics theory tells us that
61 individuals from the same source population share the same gene pool, and thus have similar
62 multilocus genotypes that are roughly in Hardy-Weinberg equilibrium and linkage
63 equilibrium. The method, implemented in the computer program STRUCTURE, has largely
64 solved the problem challenging the traditional population structure analysis, and has
65 revolutionized our ability to conduct unsupervised population structure analysis using marker
66 genotypes only. Following Pritchard *et al.* (2000), many similar methods (e.g. Dawson and
67 Belkhir 2001; Dupanloup *et al.* 2002; Corander *et al.* 2003; Guillot *et al.* 2005; Tang *et al.*
68 2005; François *et al.* 2006; Gao *et al.* 2007; Huelsenbeck and Andolfatto 2007; Alexander *et al.*
69 *et al.* 2009; Jombart *et al.* 2010; Raj *et al.* 2014) have been developed to infer population
70 structure with higher computational efficiency and with extended models (e.g. inbreeding
71 models to accommodate inbred individuals and spatial models to use geographic as well as
72 genetic data). The most popular method remains that of Pritchard *et al.* (2000), because of its
73 accurate, robust and versatile models thanks to Pritchard and coworker's original and
74 continued work (e.g. Falush *et al.* 2003; 2007; Hubisz *et al.* 2009). The dominance of
75 STRUCTURE over other programs in marker-based population structure/admixture analyses,
76 even in this genomics era, is readily confirmed by a survey of studies published in the most
77 recent issues of peer-reviewed journals such as Molecular Ecology.

78 STRUCTURE and related methods work well in assigning individuals to their source
79 populations for a given number of populations, K . When population differentiation is
80 substantial or/and marker information is sufficient, they give accurate individual ancestry (or
81 population membership) inferences. However, they have difficulties in identifying source
82 populations and inferring the optimal number of populations, K , represented by the sampled
83 individuals. Pritchard *et al.* (2000) proposed an *ad hoc* procedure to estimate the marginal
84 likelihood $\Pr[X|K]$, the probability of obtaining the genotype data X given K . The K value that

85 maximizes $\Pr[X|K]$ is the best estimate of the number of populations. They demonstrated,
86 using a couple of simple empirical datasets, that the method works well. Evanno *et al.* (2005)
87 found by simulations that the $\Pr[X|K]$ method gives poor estimates of K for hierarchically
88 structured populations, and proposed another *ad hoc* statistic, ΔK (i.e. the rate of change of
89 the probability of data as a function of K), to estimate the number of populations at the
90 uppermost hierarchical level of structure. Alexander and Lange (2011) employed a cross-
91 validation method, implemented in ADMIXTURE software, to identify the best K value as
92 judged by the prediction of systematically withheld data points. Gao *et al.* (2007) and Durand
93 *et al.* (2009) used the deviance information criterion (DIC) for inferring K in InStruct and
94 TESS programs, respectively. Dawson and Belkhir (2001), Pella and Masuda (2006), and
95 Huelsenbeck and Andolfatto (2007) took both K and individual assignments to populations as
96 random variables and used joint priors, such as a Dirichlet process prior, to estimate both in
97 programs PARTITION and STRUCTURAMA. Corander *et al.* (2003, 2004) implemented a
98 split-and-merge algorithm in their program BAPS to estimate K . Patterson *et al.* (2006)
99 proposed an eigenanalysis method, implemented in SmartPCA software, to estimate K as 1
100 plus the number of significant eigenvalues explaining the variation of genotype data. Jombart
101 *et al.* (2010) and Beugin *et al.* (2018) used Akaike information criterion (AIC, Akaike 1998),
102 Bayesian Information Criterion (BIC, Schwarz 1978), Kullback Information Criterion (KIC,
103 Cavanaugh 1999) and their variants to assess the best supported model, and therefore the
104 most likely number of populations. These and other methods were demonstrated to yield
105 good estimates of K in some simple scenarios (e.g. Gao *et al.* 2011), but can be highly
106 inaccurate in difficult situations such as many source populations (say, $K > 10$), unbalanced
107 sample sizes (Wang 2017), hierarchical population structures (Evanno *et al.* 2005), weak
108 differentiation or low marker information (Gao *et al.* 2011), and high admixture.

109 Except for the cross-validation method (Alexander and Lange 2011) and the
110 significant eigenvalue method (Patterson *et al.* 2006), all K estimators described above are
111 based (in one form or another) on the estimated marginal likelihood of the model or the
112 probability of data. This is, in theory, a natural choice for measuring model fit. In practice,
113 however, several problems arise. First, this quantity is difficult to calculate accurately, and as
114 a result, some *ad hoc* approximation is adopted (Pritchard *et al.* 2000). It is unclear how well
115 the approximation works, especially when marker information is insufficient or the inference
116 is difficult (e.g. with many populations, unbalanced sampling, and low differentiation).
117 Second, the likelihood maximization procedures (e.g. expectation maximization algorithm,

118 EM) or the Bayesian Markov Chain Monte Carlo (MCMC) procedures may not converge for
119 this high dimensional optimization problem. The number of parameters to be estimated are
120 roughly $N(K-1)+LK(A-1)$, where N , K , A , and L are the number of sampled individuals,
121 number of populations, average number of alleles at a locus, and the number of loci. The
122 estimated likelihood or probability of data could vary substantially among replicate runs of
123 the same data, especially in difficult situations (above). Third, the criteria of these *ad hoc*
124 methods in selecting the best K value may not be appropriate. For example, AIC, BIC, and
125 KIC are all based on the same principle, assessing model quality by considering its likelihood
126 against its complexity. Apparently, the penalty for model complexity is different among these
127 criteria, and it is unclear which (if any) is the most suitable for this clustering problem. It is
128 possible that none applies in general, and some modified forms of these criteria (e.g. Chen
129 and Chen 2008; Gao and Song 2010) might be more appropriate for this high dimensional
130 clustering problem.

131 In this paper, I propose another *ad hoc* criterion to estimate K , and use extensive
132 simulations and empirical data to show that it is in general more accurate than previous
133 methods. The criterion is based mainly on the quality of individual ancestry estimates from
134 STRUCTURE-like programs, and the best K is the one that consistently yields the minimal
135 mean admixture of sampled individuals. This principle, called minimal admixture or
136 parsimony for simplicity hereafter, is derived from the observation that suboptimal K (i.e.
137 values higher or lower than the true K) usually leads to inconsistent and inflated admixture
138 (mixed ancestry) estimates from Bayesian (e.g. Pritchard *et al.* 2000) or likelihood (e.g. Tang
139 *et al.* 2005; Alexander *et al.* 2009) population clustering analyses. The method is
140 implemented in a computer program, KFinder (<https://www.zsl.org/science/software/KFinder>), to
141 yield the best K given the outputs for a range of K values from a STRUCTURE-like program.
142 Extensive simulations, considering many population scenarios (e.g. differentiation levels,
143 subdivision models) and sampling scenarios (e.g. sample sizes per subpopulation, numbers of
144 markers), were conducted to compare the performances of the parsimony method, the $\Pr[X|K]$
145 method (Pritchard *et al.* 2000) and the ΔK method (Evanno *et al.* 2005) in estimating K from
146 STRUCTURE outputs, and of the parsimony method and the cross-validation method
147 (Alexander and Lange 2011) in estimating K from ADMIXTURE outputs. I showed that the
148 parsimony method improves K estimates consistently and sometimes dramatically over other
149 methods calculated from both programs STRUCTURE and ADMIXTURE.

150 **Methods**

151 ***The parsimony method***

152 I assume a dataset was analysed by STRUCTURE (or related methods such as ADMIXTURE)
 153 under the same set of model parameters (e.g. admixture and correlated allele frequency
 154 models) except for different K values, from a low bound K_L to a high bound K_H . For each K
 155 value ($K = K_L, K_{L+1}, \dots, K_{H-1}, K_H$), a number of $n_r=20$ replicate runs were conducted,
 156 following Evanno *et al.* (2005). Therefore, the total number of runs (and output files) for a
 157 single dataset is $N_R = (K_H - K_L + 1) n_r$. Parameter options for STRUCTURE and
 158 ADMIXTURE runs are detailed below.

159 A parsimony index was calculated for each K in the range K_L to K_H , using the
 160 information in the n_r replicate-run output files. First, the mean estimate of the log probability
 161 of data, $\Pr[X|K]$, was calculated from the n_r output files, and the replicate runs with $\Pr[X|K]$
 162 values smaller than the mean were discarded from further analyses. This is because the
 163 clustering algorithms may not converge, and different runs may end up with highly different
 164 $\Pr[X|K]$ values and individual ancestry estimates. Replicate runs with low $\Pr[X|K]$ values are
 165 expected to give poor estimates of individual ancestries and are thus abandoned. Note in the
 166 best scenario where all n_r replicate runs converge, no runs are abandoned, such that the
 167 number of retained runs, $n'_r = n_r$.

168 Second, an assignment quality score is calculated for each of the n'_r retained replicate
 169 runs. The (main) source population of an individual is determined as the one that has the
 170 largest ancestry coefficient estimate for the individual. The average co-assignment score
 171 within populations is calculated as

$$172 \quad ACS_w = \frac{1}{n_w} \sum_{i=1}^N \sum_{j=1, ij \in \Phi}^i \sum_{l=1}^k q_l^{(i)} q_l^{(j)}, \quad (1)$$

173 where $q_l^{(i)}$ and $q_l^{(j)}$ are the estimated proportions of individual i 's and individual j 's genomes
 174 that originate from population l ($=1, 2, \dots, k$), respectively, $ij \in \Phi$ signifies that individuals i
 175 and j are inferred to come from the same (main) source population, n_w is the number of pairs
 176 of individuals (including an individual with itself) that share the same inferred source
 177 population, and N is the number of sampled individuals. The quantity calculated by (1) is
 178 similar to Dawson and Belkhir's (2001) probability of co-assignment. Both measure the
 179 probability that a set of individuals belong to the same population, independent of the
 180 arbitrary labelling of source populations.

181 The average co-assignment score between populations, ACS_b , is calculated similarly
 182 by (1), except n_w is replaced by n_b , the number of pairs of individuals that have different
 183 source populations, and $ij \in \Phi$ is replaced by $ij \in \emptyset$ which signifies that individuals i and j
 184 are inferred to belong to different source populations.

185 The strength of population structure is characterized by a high value of ACS_w close to
 186 1, and a small value of ACS_b close to 0. When there is no hybridization (admixture) and
 187 individual ancestry inference is perfect, ACS_w reaches its maximal value of 1 and ACS_b
 188 reaches its minimal value of 0. An overall measure of the strength of population structure is

$$189 \quad SPS = ACS_w - ACS_b. \quad (2)$$

190 Third, the harmonic mean of the sizes of well-defined clusters is calculated. For each
 191 individual i ($=1, 2, \dots, N$), its main source population is determined to be l ($=1, 2, \dots, k$) if
 192 $q_l^{(i)}$ is the largest and $q_l^{(i)} \geq Q_{min}$, where Q_{min} is a chosen threshold value (say, 0.8). Its main
 193 source population is undetermined if $q_l^{(i)} < Q_{min}$. The proportion of the N sampled
 194 individuals whose main source population can be determined is then calculated. If this
 195 proportion is not smaller than Q_{min} , then main cluster structure of the sample is obtained.
 196 Otherwise, the value of Q_{min} is halved and used to repeat the above process until the main
 197 cluster structure of the sample is attained. The size of the l th (for $l=1, 2, \dots, k$) cluster, S_l , is
 198 calculated as the number of individuals whose main source populations are determined to be l .
 199 The harmonic mean of the cluster sizes that are larger than 5, $HMCS$, is then calculated,

$$200 \quad HMCS = n_{5+} / \sum_{l=1, S_l > 5}^k S_l^{-1}, \quad (3)$$

201 where n_{5+} is the number of clusters with each containing 5 or more individuals. The
 202 threshold cluster size, 5, is more or less arbitrary. However, it is chosen to reduce the
 203 population splitting errors and the population merging errors.

204 Fourth, the overall strength of the inferred population structure for a given run is
 205 calculated as

$$206 \quad SPS' = \frac{SPS}{HMCS}. \quad (4)$$

207 Among the n'_r runs retained after $\Pr[X|K]$ screening in step 1, the largest value of SPS' , SPS^* ,
 208 is obtained.

209 Fifth, the total number of clusters whose sizes are not larger than 5 (as determined in
 210 step 3) across the n'_r retained runs, n_{5-} , is calculated.

211 Sixth, the consistency of ancestry assignments cross runs is calculated as

$$212 \quad CAA = \frac{1}{N(N-1)/2} \sum_{i=1}^N \sum_{j=1}^{i-1} (\delta_{0,m}^{(ij)} + \delta_{n'_r,m}^{(ij)}), \quad (5)$$

213 where m is the number of runs among a total number of n'_r retained ones in which individuals
 214 i and j are assigned to the same cluster, $\delta_{0,m}^{(ij)} = 1$ and 0 if $m=0$ and $m > 0$ respectively,
 215 $\delta_{n'_r,m}^{(ij)} = 1$ and 0 if $m=n'_r$ and $m < n'_r$ respectively. CAA measures the consistency in main
 216 cluster assignments among replicate runs. Its maximal and minimal values are 1 and 0,
 217 respectively.

218 Seventh, the overall assignment quality for an assumed number of k populations is
 219 measured by the parsimony index

$$220 \quad PI = SPS^* + CAA - \frac{2n_{5-}}{kn_r} - \frac{1}{2k}. \quad (6)$$

221 The last term in (6) is a penalty against small k because both SPS and CAA tend to increase
 222 with a decreasing k value. In the extreme case of $k=1$, $SPS^* \equiv 1$ and $CAA \equiv 1$ because all N
 223 individuals must be inferred to come from the same source population (i.e. $ACS_w \equiv 1$ and
 224 $ACS_b \equiv 0$ in (2) and $\delta_{0,m}^{(ij)} \equiv 0$ and $\delta_{n'_r,m}^{(ij)} \equiv 1$ in (5)).

225 For each k in the range K_L to K_H , a corresponding value of PI is calculated. The k
 226 value that yields the largest PI value is inferred to be the most likely number of populations
 227 represented by the sample, K .

228 **Simulations**

229 Simulated data were generated under different population models and sampling intensities,
 230 and used to evaluate the accuracy of the above described parsimony index (PI) and two
 231 popular methods, Pr[X|K] (Pritchard *et al.* 2000) and ΔK (Evanno *et al.* 2005), in estimating K
 232 from STRUCTURE. Because increasingly large SNP datasets are produced and analysed by
 233 ADMIXTURE and other programs faster than STRUCTURE, I also simulated data with
 234 many SNPs and estimated K using PI and the cross-validation method (Alexander and Lange
 235 2011) calculated from ADMIXTURE outputs.

236 I assumed Wright's (1931) island (IS) model or a two-level hierarchical island (HI)
237 model (Evanno *et al.* 2005) for population structure in simulating genotype data. A number of
238 N_k individuals were drawn at random from population k ($k = 1, 2, \dots, K$), and each sampled
239 individual was genotyped at a number of L loci, each having A codominant alleles.

240 The ancestral allele frequencies at a marker locus l ($l=1, 2, \dots, L$), $\mathbf{p}_{0l} = \{p_{0l1}, p_{0l2}, \dots,$
241 $p_{0lA}\}$, were drawn from a uniform Dirichlet distribution, $\mathcal{D}(1,1, \dots, 1)$. The corresponding
242 allele frequencies of population i under IS model, $\mathbf{p}_{il} = \{p_{il1}, p_{il2}, \dots, p_{ilA}\}$, were drawn from
243 \mathbf{p}_{0l} following the Dirichlet distribution $\mathcal{D}(fp_{0l1}, fp_{0l2}, \dots, fp_{0lA})$, where $f = 1/F_{ST} - 1$
244 (Nicholson *et al.* 2002; Falush *et al.* 2003) and F_{ST} was the genetic differentiation among
245 populations. For HI model, the allele frequencies of an archipelago a , $\mathbf{p}_{al} = \{p_{al1}, p_{al2}, \dots,$
246 $p_{alA}\}$, were sampled from \mathbf{p}_{0l} using $\mathcal{D}(fp_{0l1}, fp_{0l2}, \dots, fp_{0lA})$, where $f = 1/F_{ST1} - 1$ and
247 F_{ST1} was the genetic differentiation among archipelagos. The allele frequencies of an island i
248 within an archipelago a were sampled from \mathbf{p}_{al} using $\mathcal{D}(fp_{al1}, fp_{al2}, \dots, fp_{alA})$, where $f =$
249 $1/F_{ST2} - 1$ and F_{ST2} was the genetic differentiation between islands within the archipelago.

250 Given the allele frequencies of a population (or island) i , the genotype of an
251 individual sampled at random from the population at L loci were generated, assuming Hardy–
252 Weinberg equilibrium and linkage equilibrium. The data, a number of N_k ($k=1, 2, \dots, K$)
253 multilocus genotypes sampled from population k , were then pooled across populations and
254 were subjected to STRUCTURE or/and ADMIXTURE analysis. Data simulated with a few
255 loci and many alleles per locus were analysed by STRUCTURE only, while data simulated
256 with many SNPs ($A=2$) were analysed by both STRUCTURE and ADMIXTURE, or by the
257 latter only.

258 *Simulation 1, differentiation F_{ST}* : The accuracy of a population structure analysis relies on the
259 strength, measured by F_{ST} , of the true structure. Populations of low F_{ST} values (close to 0) are
260 difficult to identify and thus the number of populations represented by a sample of
261 individuals is difficult to estimate from a STRUCTURE-like analysis. This simulation
262 investigated the impact of F_{ST} in an IS model of $K=6$ populations on different K estimators.
263 The 6 populations were assumed to differentiate from the ancestor population to the same
264 extent of $F_{ST}=0.02, 0.04, 0.08$ or 0.16 , or to different extents of $F_{ST}=0.02$ for populations 1
265 and 2, $F_{ST}=0.04$ for populations 3 and 4, and $F_{ST}=0.08$ for populations 5 and 6. Thirty
266 individuals from each population were genotyped at 20 (equal F_{ST}) or 50 (unequal F_{ST}) loci,
267 each having $A=10$ alleles.

268 *Simulation 2, number of loci L:* More markers provide more information and thus should
269 yield more accurate inferences of population structure and K . This simulation considered an
270 IS model of 6 populations differentiated to the same level of $F_{ST}=0.1$. Thirty individuals
271 were sampled from each population and genotyped at a varying number of loci (each having
272 $A=10$ alleles), $L=4, 6, 8, 10, 12, 14, 16, 18, 20$.

273 *Simulation 3, number of populations K:* It becomes increasingly challenging to estimate K
274 accurately with an increasing number of populations represented by a sample of individuals.
275 This simulation generated data from an IS model of $K (=1, 2, \dots, 20)$ populations
276 differentiated to the same level of $F_{ST}=0.1$, and compared the accuracy of different K
277 estimators. Thirty individuals were sampled from each population and genotyped at $L=10$ and
278 $L=20$ loci, each having $A=10$ alleles.

279 *Simulation 4, unbalanced sampling:* Population structure is difficult to infer from a sample
280 containing many individuals from one population but few individuals from another. Heavily
281 represented populations tend to split while lightly represented populations tend to merge in
282 reconstructing the population structure from such an unbalanced sample of individuals. This
283 simulation considered $K=3$ populations in an IS model with $F_{ST}=0.1$. The sample size was
284 fixed at 300 individuals, with a number of X individuals sampled from population 1 and the
285 remaining $300 - X$ individuals sampled equally from populations 2 and 3. X took values of
286 100, 120, 140, ..., 280, resulting in a perfectly balanced sample when $X=100$, and a highly
287 unbalanced sample when $X=280$. Each sampled individual was genotyped at $L=20$ loci, each
288 having $A=10$ alleles.

289 *Simulation 5, hierarchical structure:* The HI model has two true K values, the number of
290 archipelagos (K_a) and the number of islands (K_i). While the ΔK method (Evanno *et al.* 2005)
291 was shown to estimate K_a , it is unclear what the $\text{Pr}[X|K]$ method (Pritchard *et al.* 2000)
292 estimates. Is it K_a , K_i , or neither? The PI method estimates K_i , because the islands have a
293 much smaller harmonic mean cluster size than archipelagos. This simulation considered a HI
294 model of K_a archipelagos, each containing K_a islands (such that $K_i = K_a^2$), where $K_a = 2, 3$ and
295 4. Both F_{ST1} and F_{ST2} were assumed 0.1, and 30 individuals were sampled from each island
296 (total sample size $N= 30K_a^2$). Each sampled individual was genotyped at $L=20$ loci, each
297 having $A=10$ alleles.

298 *Simulation 6, hybridization:* PI index was partially based on minimizing the estimated
299 admixture, and its accuracy might be compromised for a sample containing many hybrid

300 individuals. This simulation considered an IS model of $K=3$ populations with $F_{ST}=0.1$ and
301 different degrees of hybridization. A sample contained 50 individuals from each population,
302 among which a proportion H were either F1 or F2 hybrids (at equal probabilities) between the
303 resident population and any of the other populations (with an equal probability). Each
304 sampled individual was genotyped at $L=20$ loci, each having $A=10$ alleles.

305 *Simulation 7, inbreeding:* Inbreeding causes correlation between the homologous genes at a
306 locus within an individual, and thus a loss of information in inferring population structures.
307 This simulation considered different degrees of inbreeding (due to selfing) in an IS model of
308 $K=5$ populations with $F_{ST}=0.1$. A sample contained 30 individuals from each population,
309 each individual being produced by self-reproduction at a rate s ($s=0, 0.05, 0.1, 0.2, 0.4, 0.8$)
310 or by outbreeding at a rate $1-s$. Each sampled individual was genotyped at $L=10$ or 20 loci,
311 each having $A=10$ alleles.

312 *Simulation 8, many SNPs and low F_{ST} :* With genomic data of many SNPs, it is now possible
313 to infer population structure even when it is rather weak (i.e. F_{ST} small). This simulation
314 considered an IS model of $K=5$ populations with $F_{ST}=0.01$. A sample of 20 individuals were
315 drawn from each population, and each sampled individual was genotyped at $L=100, 200, 400,$
316 $800, 1600, 3200, 6400, 12800,$ and 204800 loci, each having $A=2$ alleles. The data were
317 analysed by both STRUCTURE (except for $L=204800$) and ADMIXTURE, and K was
318 inferred by PI, $\Pr[X|K]$, ΔK and cross-validation methods.

319 *Simulation 9, many SNPs and variable F_{ST} :* Simulation 8 showed that the cross-validation
320 method always inferred $K=1$, while the truth is $K=5$, even when $L=204800$ loci were used in
321 ADMIXTURE analysis which yielded almost perfect population assignments under $K=5$.
322 Simulation 9 was conducted to investigate whether the extremely poor performance of
323 ADMIXTURE's cross-validation method was due to the low F_{ST} (0.01) or not. For this
324 purpose, $K=5$ populations in the island model with variable F_{ST} values (0.01, 0.02, 0.03, 0.04,
325 0.05, 0.06, 0.07, 0.08, 0.09, 0.10) were simulated. A sample of 20 individuals were drawn
326 from each population, and each sampled individual was genotyped at $L=1000$ loci, each
327 having $A=2$ alleles. The data were analysed by ADMIXTURE, and K was inferred by the PI
328 and cross-validation methods.

329 ***Structure analysis***

330 The simulated data were analysed by STRUCTURE program (version 2.3.4, Pritchard *et al.*
331 2000), and the analysis results were further analysed by the $\Pr[X|K]$ method (Pritchard *et al.*
332 2000), the ΔK method (Evanno *et al.* 2005) and the new PI method for the most likely
333 number of populations, K . The model and parameter settings adopted in the analyses were
334 admixture model and correlated allele frequency model, INFERALPHA=1, ALPHA=1.0,
335 POPALPHAS=1, UNIFPRIORALPHA=1, ALPHAMAX=10.0, ALPHAPROPSD=0.025.
336 The alternative prior for individual ancestry was adopted, because it gave much more
337 accurate STRUCTURE analysis results than the default prior when sampling was highly
338 unbalanced (Wang 2017). The burn-in length was 10^4 , 10^5 and 10^6 iterations in analysing data
339 simulated with $K < 6$, $6 \leq K < 10$, and $K \geq 10$ populations, respectively. More populations lead
340 to more parameters for STRUCTURE to estimate and therefore pose more challenges for the
341 MCMC algorithm to converge. These burn-in lengths were obtained by experimenting with
342 many pilot analyses of data with different simulated K values. The run length was 10^4
343 iterations. For all other parameters not mentioned above, their default values were used.

344 For a dataset simulated with a given K , STRUCTURE analyses were conducted for
345 each assumed number of populations in the range $[K_L, K_H]$, where $K_L = \text{Max}[K-3, 1]$ and K_H
346 $= K+3$. This narrow range was adopted to reduce the computational burden, and because the
347 primary interest was whether the true simulated K value was recovered or not. For each
348 assumed K value k , a number of $n_r=20$ replicate runs were conducted.

349 The simulated data with many SNPs were also analysed by ADMIXTURE program
350 (version 1.3.0, Alexander and Lange 2011). The SNP data were first reformatted by PLINK
351 (version 1.9.0, Purcell *et al.* 2007), and then analysed by ADMIXTURE using the program's
352 default settings. For each dataset and each assumed K in the range $[K_L, K_H]$ where K_L
353 $= \text{Max}[K-3, 1]$ and $K_H = K+3$, a number of 10 independent replicate runs of ADMIXTURE
354 were conducted. The most likely K was inferred from ADMIXTURE outputs using its default
355 cross-validation method and the PI method.

356 All simulations and data analyses described above were conducted on a large Linux
357 cluster, using many cores in parallel by MPI.

358 ***Accuracy assessment***

359 For each simulation (i.e. set of parameters), a number of 50 replicate datasets were generated.
360 Each replicate dataset was analysed by STRUCTURE or/and ADMIXTURE assuming k in

361 the range $[K_L, K_H]$, and was replicated with $n_r=20$ for each assumed K value. The total
362 number of STRUCTURE-like analyses for a single dataset was thus $(K_H - K_L + 1) \times 20$, which
363 was 80, 100, 120, 140 when the simulated K value was 1, 2, 3, and ≥ 4 , respectively. The
364 results in STRUCTURE's output files were analysed by the ΔK method as described by
365 Evanno *et al.* (2005), the $\text{Pr}[X|K]$ method (Pritchard *et al.* 2000) as described by Wang (2017),
366 and the PI method as described above to obtain the estimates of K , denoted by \hat{K}_{Ev} , \hat{K}_{Pr} , and
367 \hat{K}_{PI} respectively. The ADMIXTURE program outputs were used to yield the K estimates
368 from both PI method and the cross-validation method (Alexander and Lange 2011), denoted
369 by \hat{K}_{AL} .

370 The accuracy of an estimator was measured by the proportion of replicate datasets in
371 which the estimator was equal to the simulated true K . These accuracy measurements were
372 denoted as $\text{Pr}(\hat{K}_{Ev} = K)$, $\text{Pr}(\hat{K}_{Pr} = K)$, $\text{Pr}(\hat{K}_{PI} = K)$ and $\text{Pr}(\hat{K}_{AL} = K)$ for estimators \hat{K}_{Ev} ,
373 \hat{K}_{Pr} , \hat{K}_{PI} and \hat{K}_{AL} , respectively.

374 ***A human dataset***

375 The performance of the three K estimators, \hat{K}_{Ev} , \hat{K}_{Pr} , and \hat{K}_{PI} , was also compared by
376 analysing a human dataset, published in Wang *et al.* (2007). The dataset contains 1484
377 individuals sampled from 78 world-wide populations, each individual being genotyped at 678
378 microsatellite loci. It proves to be difficult to reconstruct the population structure
379 unambiguously from the 1484 sampled individuals, even using all of the 678 highly
380 polymorphic microsatellites (Wang *et al.* 2007). For demonstration purpose and for reducing
381 computational burden of a bootstrapping analysis, I choose to analyse a sub-dataset
382 composing of 24 Basque individuals sampled from France, 19 Melanesian individuals from
383 Bougainville, 21 Surui individuals from Brazil, 22 Mandenka individuals from Senegal, and
384 29 Japanese individuals from Japan. These populations are well differentiated and have
385 balanced sample sizes, and as a result can be distinguished using about 20 markers. A number
386 of 100 replicate datasets were generated by bootstrapping over loci, for $L = 10, 20, 40, 80$ and
387 160. Each dataset was analysed by 20 replicate STRUCTURE runs for each assumed K value
388 from 1 to 10 (including the true value of $K=5$). The three K -estimators were then applied to
389 the STRUCTURE outputs. Accuracy of the estimators was evaluated by calculating the
390 proportions of the replicate datasets in which $K < 5$, $K = 5$, and $K > 5$ were obtained.

391

392 **Results**

393 ***Simulation 1, differentiation F_{ST}***

394 All three estimators become more accurate with an increasing F_{ST} (Figure 1). The accuracy of
395 estimators \hat{K}_{Ev} and \hat{K}_{Pr} is similar, and is consistently lower than that of \hat{K}_{Pl} .

396 Estimator \hat{K}_{Pl} has an accuracy increasing rapidly with an increasing F_{ST} value (Figure
397 1), and it recovers the simulated K value completely (i.e. accuracy = 100%) when $F_{ST} \geq 0.08$.
398 The low accuracy at a small F_{ST} value of 0.02 is due to the insufficient marker information.
399 Keeping $F_{ST} = 0.02$ and all the other parameters but increasing the number of loci to $L=100$
400 increases the accuracy of \hat{K}_{Pl} to 100%, and that of \hat{K}_{Ev} and \hat{K}_{Pr} to 72% and 82% respectively.

401 Unequal F_{ST} among populations does not affect the accuracy order of the three K
402 estimators. When $F_{ST} = \{0.02, 0.02, 0.04, 0.04, 0.08, 0.08\}$, $L=50$ and the other parameters
403 have the same values as in Figure 1, the accuracy is $\Pr(\hat{K}_{Ev} = K) = 66\%$, $\Pr(\hat{K}_{Pr} = K) = 100\%$
404 and $\Pr(\hat{K}_{Pl} = K) = 100\%$.

405 ***Simulation 2, number of loci L***

406 \hat{K}_{Ev} is less accurate than \hat{K}_{Pr} and \hat{K}_{Pl} , especially when the number of loci L is small (Figure
407 2).

408 ***Simulation 3, number of populations K***

409 \hat{K}_{Ev} cannot distinguish a panmictic population ($K=1$) from a structured or subdivided
410 population ($K > 1$), because its statistic is undefined at $K=1$. When both K and L are small
411 (Figure 3, left panel), \hat{K}_{Ev} is less accurate than the other two estimators. \hat{K}_{Pr} and \hat{K}_{Pl} have a
412 similar accuracy when $K > 1$. When $K=1$ (i.e. a single panmictic population), however, \hat{K}_{Pr}
413 and \hat{K}_{Pl} yield the correct estimate at a frequency of about 45% and 100%, respectively. When
414 $K > 10$, all three estimators become highly inaccurate (Figure 3, left panel), and require more
415 markers to yield accurate population structure inferences (Figure 3, right panel).

416 ***Simulation 4, unbalanced sampling***

417 \hat{K}_{Pr} and \hat{K}_{Pl} are far more accurate than \hat{K}_{Ev} in the case of unbalanced sampling (Figure 4).
418 \hat{K}_{Ev} is very vulnerable to the unevenness of sample sizes. It gives poor estimates of K even
419 when samples from different populations are only slightly different in size.

420 ***Simulation 5, hierarchical structure***

421 Figure 5 shows that \hat{K}_{Ev} estimates the number of archipelagos while \hat{K}_{Pr} and \hat{K}_{Pl} estimate
422 the number of islands. It also shows that \hat{K}_{Pr} is less accurate than \hat{K}_{Ev} and \hat{K}_{Pl} for different K
423 values. Although \hat{K}_{Ev} and \hat{K}_{Pl} estimate K_a and K_i respectively, they have a similar accuracy.

424 ***Simulation 6, hybridization***

425 \hat{K}_{Pl} is robust to the presence of hybrids (Figure 6). Perfect K estimates were obtained even
426 when 32% of sampled individuals are either F1 or F2 hybrids. The other two methods also
427 perform well, but relatively \hat{K}_{Ev} is the least accurate method even when hybrid frequencies
428 are low. Both \hat{K}_{Ev} and \hat{K}_{Pr} show a dip in accuracy when hybrid rate is high.

429 ***Simulation 7, inbreeding***

430 \hat{K}_{Pl} is little affected by inbreeding (Figure 7). Almost perfect K estimates were obtained even
431 when selfing rate is 80%. In contrast, \hat{K}_{Pr} works well only when inbreeding is absent or very
432 low. Its accuracy decreases rapidly with an increasing selfing rate. It outperforms \hat{K}_{Ev} when
433 selfing rate is negligibly small, but quickly becomes less accurate than \hat{K}_{Ev} with an increasing
434 selfing rate. More markers do not help. Actually when selfing rate is substantial, $L=20$ loci
435 leads to less accurate \hat{K}_{Pr} than $L=10$ (Figure 7).

436 ***Simulation 8, many SNPs and low F_{ST}***

437 At a low differentiation of $F_{ST}=0.01$, all K estimators calculated from both STRUCTURE
438 and ADMIXTURE outputs perform poorly when the number of SNPs (L) is small ($L<3200$)
439 (Figure 8). With more SNPs, \hat{K}_{Pl} quickly reaches the maximal accuracy of 100%, no matter it
440 is calculated from STRUCTURE or ADMIXTURE outputs. The other three estimators, \hat{K}_{Ev} ,
441 \hat{K}_{Pr} and \hat{K}_{CV} , still perform poorly even when a large number of SNPs (12800 and 204800 for
442 STRUCTURE and ADMIXTURE, respectively) are used. The cross-validation estimator,
443 \hat{K}_{CV} , consistently yields an estimate of $K=1$ (i.e. no population structure), using 100-204800
444 SNPs.

445 ***Simulation 9, many SNPs and variable F_{ST}***

446 The cross-validation estimator, \hat{K}_{CV} , is inaccurate when F_{ST} is small (Figure 9), as observed
447 before (Figure 8). With $F_{ST} \leq 0.06$, the cross-validation method always yields $\hat{K}_{CV} = 1$, much

448 smaller than the truth of $K=5$. The accuracy of \hat{K}_{CV} increases rapidly with increasing F_{ST}
449 when it is larger than 0.06, and reaches the maximum 100% when $F_{ST}=0.08$. In contrast, \hat{K}_{PI}
450 becomes perfect when $F_{ST} \geq 0.02$.

451 **Human SSR data**

452 \hat{K}_{Ev} consistently underestimates K , irrespective of the number of loci (Figure 10). The vast
453 majority of estimates are $\hat{K}_{Ev}=2$, with population Surui forming a cluster and the remaining
454 four populations forming the other cluster. \hat{K}_{Pr} underestimates and overestimates K when L is
455 low and high respectively, yielding a maximal accuracy (i.e. the frequency of estimates of
456 $K=5$) of 60% at $L=20$. It is a bit bizarre that \hat{K}_{Pr} gives fewer estimates of $K=5$ with an
457 increasing L when $L > 20$. The accuracy of \hat{K}_{PI} always increases with L . The estimator is
458 (when $L > 10$) or is close to (when $L = 10$) the most accurate, and yields perfect K estimates
459 when $L > 20$.

460

461 **Discussion**

462 In this study I proposed an *ad hoc* estimator of the number of populations represented in a
463 sample of individuals (K), which can be calculated from the results of a STRUCTURE-like
464 analysis. While previous estimators (e.g. Pritchard *et al.* 2000; Evanno *et al.* 2005) rely on the
465 estimated likelihood or probability of data, the new method, in contrast, evaluates and
466 employs the individual ancestry assignment quality as the criterion in choosing the most
467 likely K . It is based on the observation that, in a STRUCTURE-like analysis, assuming a
468 higher and lower than true K value leads to the splitting and merging of source populations,
469 respectively. In both cases, the individual ancestry assignment quality is usually undermined,
470 as characterized by inflated admixture within each replicate run and increased inconstancy
471 across replicate runs. Loosely speaking, the parsimony index method estimates the most
472 likely K by identifying the number of populations which yields the most consistent and the
473 minimal average admixture.

474 My extensive simulations under a variety of population structure and sampling
475 scenarios show that the new estimator (\hat{K}_{PI}) outperforms the current popular estimators
476 overall. In some difficult situations such as unbalanced sampling (Figure 4), low population
477 differentiation (Figures 8, 9, 1), low marker information (Figure 2), hierarchical structure

478 (Figure 5), and inbreeding (Figure 7), the new estimator improves K estimation substantially.
479 I also show that the new estimator is accurate when hybridization is present (Figure 6), and is
480 more accurate than other estimators when hybridization is high. This seems to be a bit of
481 surprising, given that \widehat{K}_{PI} is partially based on minimizing average admixture. However, true
482 admixture is different from false admixture. While the former is consistently inferred across
483 different K values and across different replicate runs for a given K value, the latter is
484 estimated only when the assumed K deviates from the truth and is estimated inconsistently
485 across replicate runs. Therefore, minimizing average admixture still leads to the recovery of
486 the actual population structure in the presence of true admixture.

487 My simulation confirms the conclusion that estimator \widehat{K}_{Ev} is more accurate than \widehat{K}_{Pr}
488 for populations in the hierarchical island model (Evanno *et al.* 2005). While \widehat{K}_{Ev} estimates the
489 number of archipelagos (K_a), \widehat{K}_{Pr} tends to estimate the number of islands (K_i). \widehat{K}_{Ev} estimates
490 K_a with an accuracy consistently and considerably higher than \widehat{K}_{Pr} estimates K_i (Figure 5).
491 The simulations also show \widehat{K}_{Ev} is more accurate than \widehat{K}_{Pr} for highly inbred populations
492 (Figure 7). Unfortunately, however, the accuracy advantage of \widehat{K}_{Ev} is lost in other scenarios
493 of population structure and sampling (Figures 1~4, 6). In several realistic scenarios (e.g.
494 unbalanced sampling, Figure 4), \widehat{K}_{Pr} is much more accurate than \widehat{K}_{Ev} . It is unfortunate that
495 \widehat{K}_{Ev} is now widely favoured over \widehat{K}_{Pr} in estimating K , even when there is no clear evidence
496 that the populations are in a hierarchical structure or highly inbred.

497 The confusion as to which estimator gives a better K estimate arises because all
498 estimators are *ad hoc* and their accuracies must be evaluated using simulated or empirical
499 datasets. Due to the heavy computational burden, however, few studies (e.g. Evanno *et al.*
500 2003; Gao *et al.* 2011) were conducted to compare the accuracy of different estimators under
501 various population structure and sampling scenarios. Typically a simulation study (like the
502 present one) considers many different sets of parameter combinations, and simulates and
503 analyses a large number of replicate datasets for a given parameter combination. Each
504 simulated dataset must be analysed with different assumed K values, and for each assumed K
505 value, a number of replicate runs (say, 20) must be conducted. The total number of
506 STRUCTURE runs for a single dataset is $(K_H - K_L + 1) \times 20$, which is 140 when the simulated
507 $K > 3$, $K_H = K + 3$, and $K_L = K - 3$. If a typical run takes about 0.5 hours, this means analysing a
508 single dataset takes about 70 hours. Analysing 50 replicate datasets simulated under a given
509 set of parameters would take about 3500 hours. A typical figure with 8 plotting points (values

510 on the x axis) would take about 28000 hours, and the 8 figures from STRUCTURE analyses
511 shown in this study would take about 224000 hours. It is obviously impossible to conduct a
512 simulation study like the present one on a desktop computer. My simulation was carried out
513 on a Linux cluster using 512 cores in parallel.

514 This study focussed on applying different K estimators to STRUCTURE (Pritchard *et al.*
515 *et al.* 2000) analyses. Other methods for population structure inference (e.g. Corander *et al.*
516 2003; Tang *et al.* 2005; Gao *et al.* 2007; Huelsenbeck and Andolfatto 2007; Alexander *et al.*
517 2009) use the same genotype data and give similar outputs such as individual ancestry
518 proportions. For the analysis of genomic SNP data with thousands to millions of loci,
519 STRUCTURE is too slow and much faster methods are increasingly used. Alexander *et al.*
520 (2009) improved Tang *et al.*'s (2005) expectation maximization algorithm of a likelihood
521 model, and implemented the algorithm in a computer program ADMIXTURE. The program
522 runs several orders faster than STRUCTURE, and yet provides similarly good results of both
523 ancestry assignments and K estimates (by the cross-validation method, \hat{K}_{CV}) in some tested
524 situations (Alexander *et al.* 2009; Alexander and Lange 2011). This study simulated genomic
525 data and compared the performances of the cross-validation method and other methods
526 (Figure 8). It is clear that \hat{K}_{CV} calculated from ADMIXTURE outputs behaves similarly to
527 \hat{K}_{Pr} and \hat{K}_{Ev} calculated from STRUCTURE outputs, when populations are little differentiated
528 and many SNPs are used. The accuracy of the three estimators is rather poor, compared with
529 that of the parsimony estimator \hat{K}_{PI} calculated from the outputs of both programs. It seems
530 \hat{K}_{CV} is very sensitive to the F_{ST} , and becomes accurate only for highly differentiated
531 populations (Figure 9). At low differentiation, it is conservative and always falsely infers a
532 single ($\hat{K}_{CV} = 1$) panmictic population even when many markers are used (Figures 8 and 9).

533 The parsimony K estimator described in this study was implemented in a computer
534 program, KFinder, freely downloadable from <https://www.zsl.org/science/software/KFinder>.

535

536

537

538 **References**

539 Akaike H (1998) Information theory and an extension of the maximum likelihood principle.
540 In E. Parzen, K. Tanabe, & G. Kitagawa (Eds.), Selected papers of Hirotugu Akaike (pp.
541 199–213). Springer Series in Statistics. New York, NY: Springer.

542 Alexander DH, Novembre J, Lange K (2009) Fast model-based estimation of ancestry in
543 unrelated individuals. *Genome Research*, **19**, 1655-1664.

544 Alexander DH, Lange K (2011) Enhancements to the ADMIXTURE algorithm for individual
545 ancestry estimation. *BMC Bioinformatics*, **12**, 246.

546 Beugin MP, Gayet T, Pontier D, Devillard S, Jombart T (2018) A fast likelihood solution to
547 the genetic clustering problem. *Methods in Ecology and Evolution*, **9**, 1006-1016.

548 Cavanaugh JE (1999) A large- sample model selection criterion based on Kullback's
549 symmetric divergence. *Statistics & Probability Letters*, **42**, 333–343.

550 Chen J, Chen Z (2008) Extended Bayesian information criteria for model selection with large
551 model spaces. *Biometrika*, **95**, 759–771.

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

554 Corander J, Waldmann P, Marttinen P, Sillanpää MJ (2004) BAPS2: enhanced possibilities
555 for the analysis of population structure. *Bioinformatics* **20**, 2363-2369.

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

558 Dupanloup I, Schneider S, Excoffier L (2002) A simulated annealing approach to define the
559 genetic structure of populations. *Molecular Ecology*, **11**, 2571-2581.

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

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

564 Falush D, Stephens M, Pritchard JK (2007) Inference of population structure using multilocus
565 genotype data: dominant markers and null alleles. *Molecular Ecology Resources*, **7**, 574-
566 578.

567 François O, Ancelet S, Guillot G (2006) Bayesian clustering using hidden Markov random
568 fields in spatial population genetics. *Genetics*, **174**, 805–816.

569 Gao X, Song P (2010) Composite likelihood Bayesian information criteria for model
570 selection in high-dimensional data. *Journal of the American Statistical Association*, **105**,
571 1531–1540.

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

575 Gao H, Bryc K, Bustamante CD (2011) On identifying the optimal number of population
576 clusters via the deviance information criterion. *PloS one*, **6**, e21014.

577 Guillot G, Mortier F, Estoup A (2005) GENELAND: a computer package for landscape
578 genetics. *Molecular Ecology Resources*, **5**, 712-715.

579 Hsieh HM, Huang LH, Tsai LC, Kuo YC, Meng HH, Linacre A, *et al.* (2003) Species
580 identification of rhinoceros horns using the cytochrome b gene. *Forensic Science*
581 *International*, **136**, 1–11.

582 Hubisz MJ, Falush D, Stephens M, Pritchard JK (2009) Inferring weak population structure
583 with the assistance of sample group information. *Molecular Ecology Resources*, **9**, 1322-
584 1332.

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

587 Jombart T, Devillard S, Balloux F (2010) Discriminant analysis of principal components: a
588 new method for the analysis of genetically structured populations. *BMC Genetics*, **11**, 94.

589 Nicholson G, Smith AV, Jonsson F, Gustafsson O, Stefansson K, Donnelly P (2002)
590 Assessing population differentiation and isolation from single nucleotide polymorphism
591 data. *Journal of the Royal Statistical Society Series B*, **64**, 695–715.

592 Patterson NJ, Price AL, Reich D (2006) Population structure and eigenanalysis. *PLoS*
593 *Genetics*, **2**, e190.

594 Pella J, Masuda M (2006) The gibbs and split-merger sampler for population mixture analysis
595 from genetic data with incomplete baselines. *Canadian Journal of Fisheries and Aquatic*
596 *Sciences*, **63**, 576–596.

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

- 599 Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, Maller J, Sklar P, De
600 Bakker, PI, Daly MJ, Sham PC (2007) PLINK: a tool set for whole-genome association
601 and population-based linkage analyses. *The American Journal of Human Genetics*, **81**,
602 559-575.
- 603 Raj A, Stephens M, Pritchard JK (2014) fastSTRUCTURE: variational inference of
604 population structure in large SNP data sets. *Genetics*, **197**, 573-589.
- 605 Schwarz G (1978) Estimating the dimension of a model. *Annals of Statistics*, **6**, 461-464.
- 606 Smouse PE, Waples RS, Tworek JA (1990) A genetic mixture analysis for use with
607 incomplete source population data. *Canadian Journal of Fisheries and Aquatic Sciences*,
608 **47**, 620-634.
- 609 Tang H, Peng J, Wang P, Risch NJ (2005) Estimation of individual admixture: analytical and
610 study design considerations. *Genetic Epidemiology*, **28**, 289-301.
- 611 Velo-Anton G, Godinho R, Ayres C, Ferrand N, Rivera AC (2007) Assignment tests applied
612 to relocate individuals of unknown origin in a threatened species, the European pond
613 turtle (*Emys orbicularis*). *Amphibia-Reptilia*, **28**, 475-84.
- 614 Wang J (2017) The computer program STRUCTURE for assigning individuals to populations:
615 easy to use but easier to misuse. *Molecular Ecology Resources*, **17**, 981-990.
- 616 Wang S, Lewis Jr CM, Jakobsson M, Ramachandran S, Ray N, Bedoya G, Rojas W, Parra
617 MV, Molina JA, Gallo C, Mazzotti G (2007) Genetic variation and population structure
618 in Native Americans. *PLoS Genetics*, **3**, e185.
- 619 Weir BS, Cockerham CC (1984) Estimating F -statistics for the analysis of population
620 structure. *Evolution*, **38**, 1358-1370.
- 621 Wright S (1931) Evolution in Mendelian populations. *Genetics*, **16**, 97-159.
- 622 Wright S (1951) The genetical structure of populations. *Annals of Eugenics*, **15**, 323-354.

623

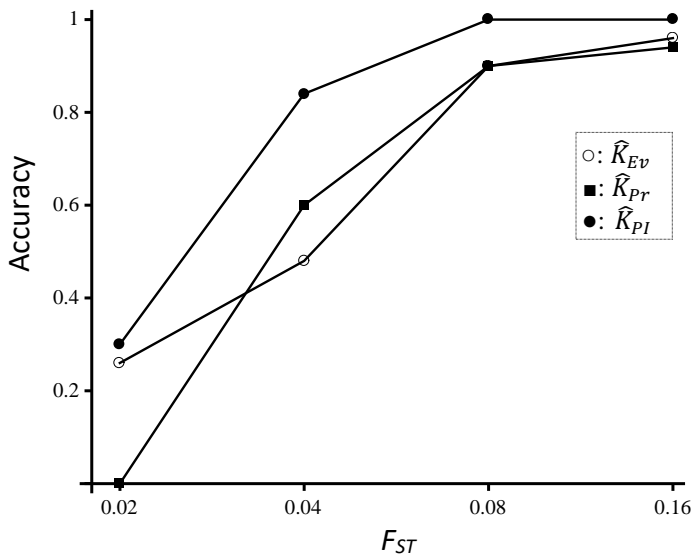
624 **Data Accessibility**

625 Source code (in Fortran 9x) for simulating genotype data, preparing input files for
626 STRUCTURE, running STRUCTURE, and calculating the three K estimators: DRYAD entry
627 DOI: #####

628 **Acknowledgements**

629 I am grateful to the editor, N Barton, and two anonymous reviewers for helpful comments on
630 previous versions of this article.

631



632

633

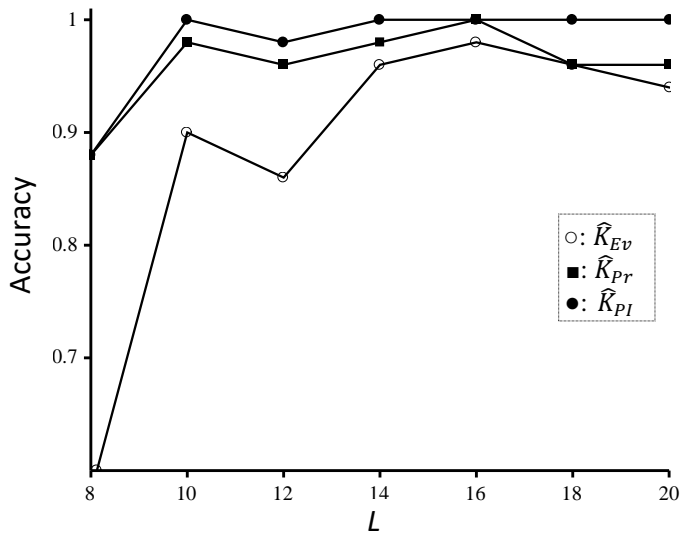
634

635

636

637 Figure 1: Accuracy of three K estimators (calculated from STRUCTURE outputs) as a
638 function of F_{ST} . A number of $K=6$ populations in the island model with $F_{ST}=0.02, 0.04, 0.08$
639 or 0.16 (x axis) were simulated. Thirty individuals from each population were sampled and
640 genotyped at 20 loci, each having $A=10$ alleles.

641



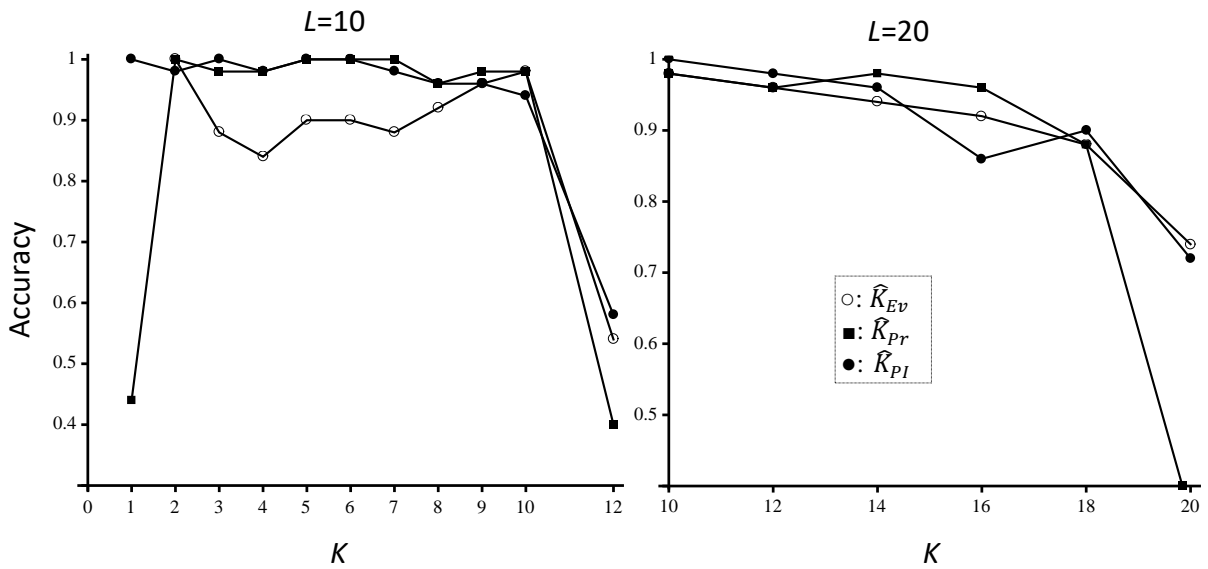
642

643

644

645 Figure 2: Accuracy of three K estimators (calculated from STRUCTURE outputs) as a
 646 function of the number of loci L . A number of $K=6$ populations in the island model with F_{ST}
 647 $=0.1$ were simulated. Thirty individuals from each population were sampled and genotyped at
 648 $L=8, 10, 12, 14, 16, 18$ and 20 loci, each having $A=10$ alleles.

649



650

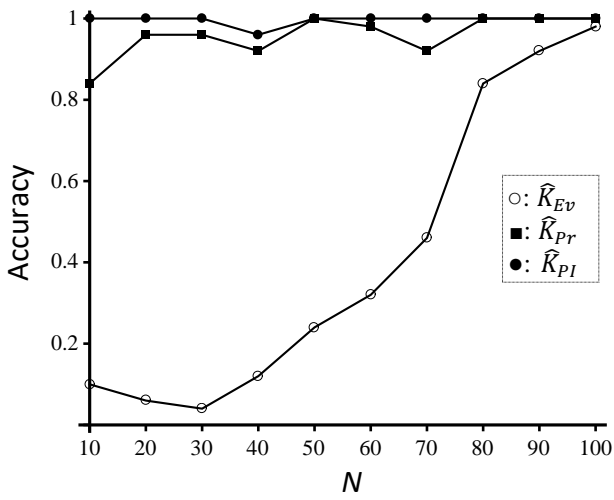
651

652

653 Figure 3: Accuracy of three K estimators (calculated from STRUCTURE outputs) as a
 654 function of the simulated (true) number of populations (K). A number of K (x axis)
 655 populations in the island model was simulated, assuming $F_{ST}=0.1$. A number of 30
 656 individuals were sampled from each population, and each individual was genotyped at either
 657 $L=10$ (left panel) or $L=20$ (right panel) loci, each having 10 alleles.

658

659



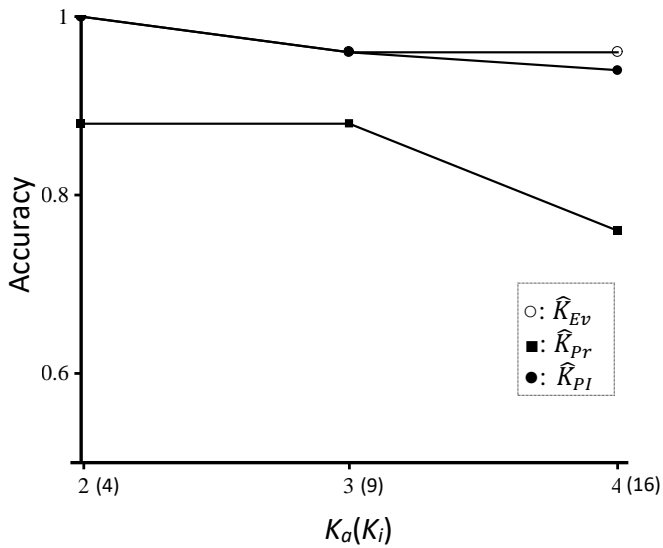
660

661

662

663 Figure 4: Accuracy of three K estimators (calculated from STRUCTURE outputs) as a
664 function of the smaller sample size N . Three populations in the island model was simulated,
665 assuming $F_{ST}=0.1$. A number of $300-2N$, N , and N individuals were sampled from
666 populations 1, 2, and 3, respectively. The individuals were genotyped at $L=20$ loci, each
667 having 10 alleles.

668



669

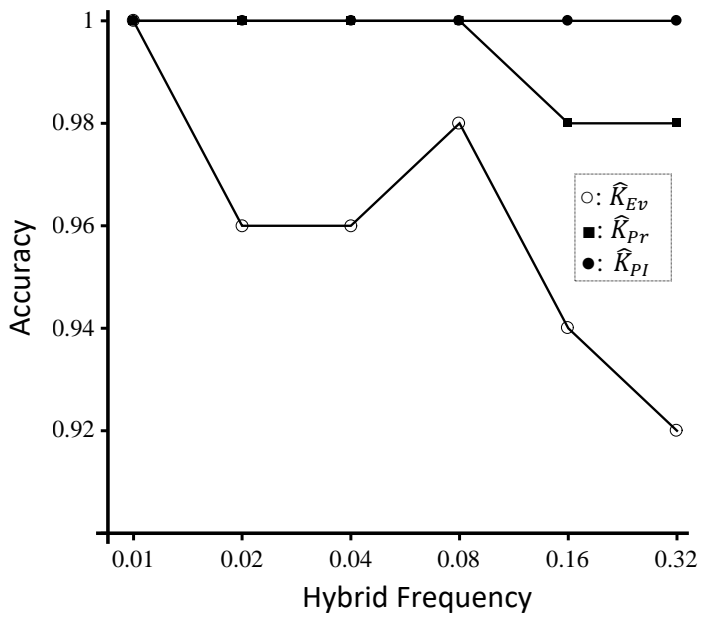
670

671

672

673 Figure 5: Accuracy of three K estimators (calculated from STRUCTURE outputs) as a
 674 function of the number of archipelagos K_a or the number of islands K_i in a HI model. Note,
 675 the HI model has two true K values, the number of archipelagos (K_a) and the number of
 676 islands (K_i), the latter being equal to the square of the former in the simulations. Accuracy
 677 $\Pr(\hat{K}_{Ev} = K)$, $\Pr(\hat{K}_{Pr} = K)$ and $\Pr(\hat{K}_{PI} = K)$ is calculated with K being K_a for \hat{K}_{Ev} , and
 678 being K_i for \hat{K}_{Pr} and \hat{K}_{PI} . The F_{ST} values among archipelagos and among islands within an
 679 archipelago are both 0.1. Thirty individuals were sampled from each island, and were
 680 genotyped at $L=20$ loci, each having 10 alleles.

681



682

683

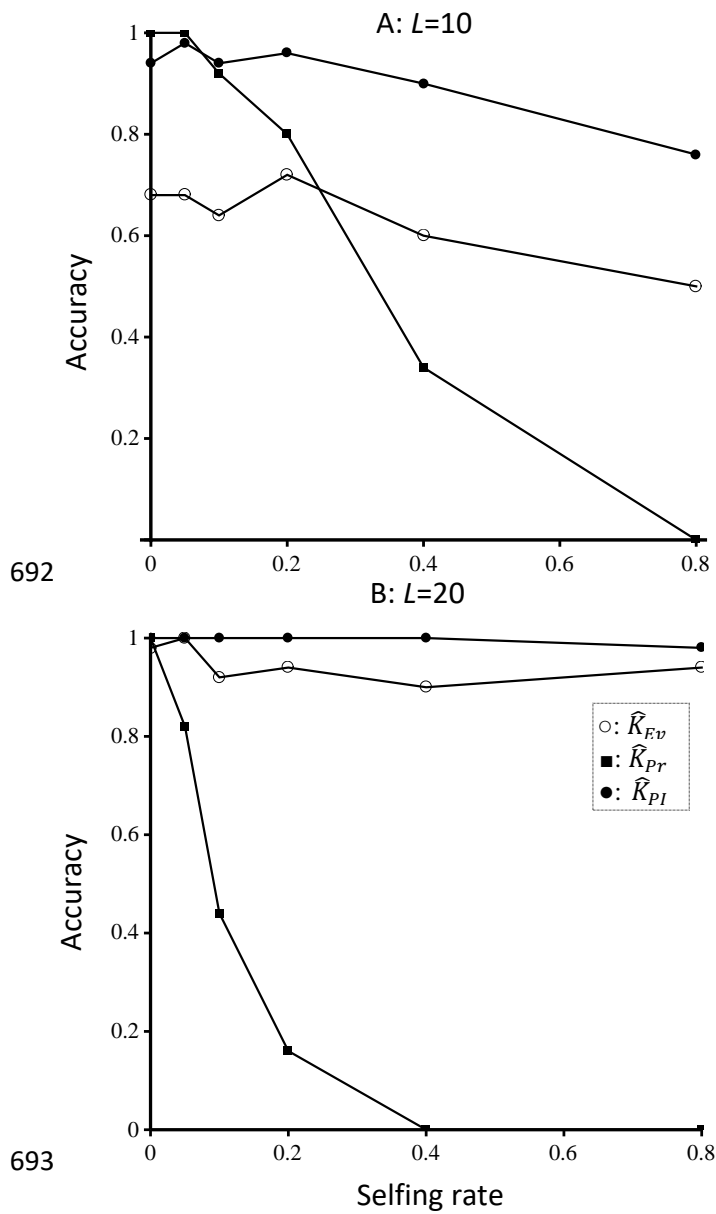
684

685

686

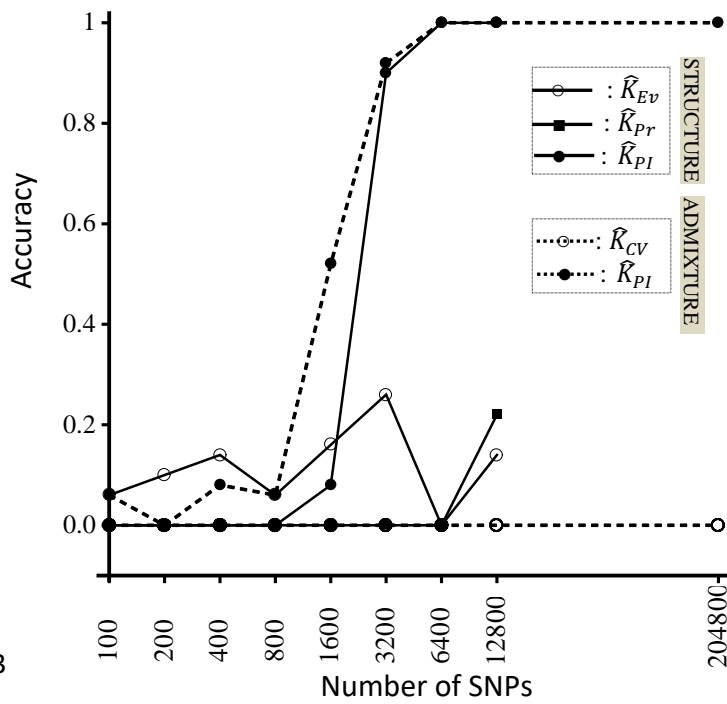
687 Figure 6: Accuracy of three K estimators (calculated from STRUCTURE outputs) as a
 688 function of the frequency of hybrids (F1 and F2) in a sample of individuals. Three
 689 populations in the island model with $F_{ST}=0.1$ were simulated, and 30 individuals were
 690 sampled from each population and were genotyped at $L=20$ loci, each having 10 alleles.

691



697 Figure 7: Accuracy of three K estimators (calculated from STRUCTURE outputs) as a
 698 function of the selfing rate of sampled individuals. Five populations in the island model with
 699 $F_{ST}=0.1$ were simulated, and 30 individuals were sampled from each population and were
 700 genotyped at $L=10$ (upper panel) or 20 (lower panel) loci, each having 10 alleles.

702



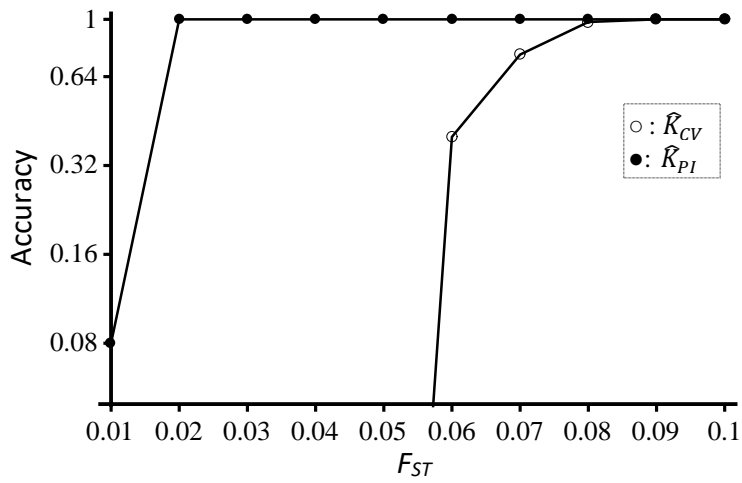
703

704

705

706 Figure 8: Accuracy of estimators \hat{K}_{EV} , \hat{K}_{Pr} , \hat{K}_{PI} calculated from STRUCTURE outputs and \hat{K}_{CV}
 707 and \hat{K}_{PI} calculated from ADMIXTURE outputs as a function of the number of SNPs used in
 708 population structure analyses. Five populations in the island model with $F_{ST}=0.01$ were
 709 simulated, and 20 individuals were sampled from each population and were genotyped at
 710 $L=100 - 204800$ (x axis) loci, each having 2 alleles.

711



712

713

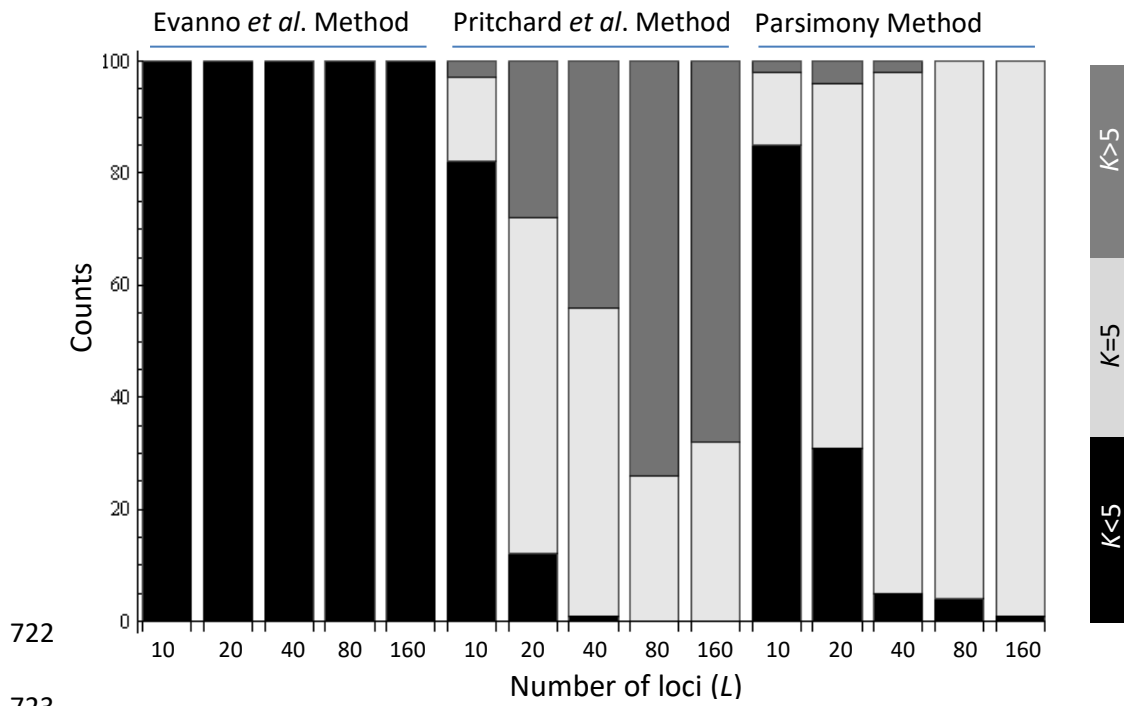
714

715

716

717 Figure 9: Accuracy of estimators \hat{K}_{CV} and \hat{K}_{PI} calculated from ADMIXTURE outputs as a
 718 function of F_{ST} . Five populations in the island model with F_{ST} varying (x axis) in the range
 719 $[0.01, 0.10]$ were simulated, and 20 individuals were sampled from each population and were
 720 genotyped at $L=1000$ loci, each having 2 alleles.

721



722

723

724

725

726 Figure 10: Accuracy of three K estimators as a function of the number of loci, L . A number of
 727 100 replicate datasets were obtained by bootstrapping (over loci) for each number of loci (x
 728 axis) from a human SSR dataset, and were analysed by STRUCTURE for estimating K . The
 729 original dataset has 117 individuals sampled from 5 (true $K=5$) populations.