

1 *Molecular Ecology Ressources – subject area: Methodological Advances*

2

3 **Estimation of demo-genetic model probabilities with Approximate Bayesian**

4 **Computation using linear discriminant analysis on summary statistics**

5

6 ARNAUD ESTOUP¹, ERIC LOMBAERT², JEAN-MICHEL MARIN³, THOMAS

7 GUILLEMAUD², PIERRE PUDLO^{1,3}, CHRISTIAN P. ROBERT^{4,5,6} and JEAN-

8 MARIE CORNUET¹

9

10 ¹ INRA UMR Centre de Biologie et de Gestion des Populations (INRA / IRD / Cirad /
11 Montpellier SupAgro), Montferrier-sur-Lez, France.

12 ² Equipe "Biologie des Populations en Interaction", UMR 1301 IBSV INRA-CNRS-
13 Université de Nice-Sophia Antipolis, Sophia-Antipolis, France.

14 ³ I3M, UMR CNRS 5149, Université Montpellier 2, France.

15 ⁴ Université Paris Dauphine, CEREMADE, Paris, France.

16 ⁵ Institut Universitaire de France, Paris, France.

17 ⁶ CREST, Paris, France.

18

19

20 **Keywords:** Approximate Bayesian Computation, coalescence, discriminant analysis,
21 model probability, evolutionary scenario, molecular markers, population genetics

22

23

24 **Correspondence:** Arnaud Estoup

25 Centre de Biologie et de Gestion des Populations

26 Institut National de la Recherche Agronomique

27 Campus International de Baillarguet CS 30 016

28 34988 Montferrier / Lez, FRANCE

29 Fax: +33 (0)4.99.62.33.45

30 E-mail: estoup@supagro.inra.fr

31

32

33 **Running title:** ABC model probabilities using LDA-transformation

34 **Abstract**

35

36 Comparison of demo-genetic models using Approximate Bayesian Computation (ABC) is
37 an active research field. Although large numbers of populations and models (i.e. scenarios)
38 can be analysed with ABC using molecular data obtained from various marker types,
39 methodological and computational issues arise when such numbers become too large.

40 Moreover, Robert *et al.* (2011) have shown that the conclusions drawn on ABC model
41 comparison cannot be trusted per se and required further simulation analyses. Monte Carlo
42 inferential techniques to empirically evaluate confidence in scenario choice are very time
43 consuming, however, when the numbers of summary statistics (Ss) and scenarios are large.

44 We here describe a methodological innovation to process efficient ABC scenario
45 probability computation using linear discriminant analysis (LDA) on Ss before computing
46 logistic regression. We used simulated pseudo-observed datasets (*pods*) to assess the main
47 features of the method (precision and computation time) in comparison to more traditional
48 probability estimation using raw (i.e. not LDA-transformed) Ss. We also illustrate the
49 method on real microsatellite datasets produced to make inferences about the invasion
50 routes of the coccinellid *Harmonia axyridis*. We found that scenario probabilities computed
51 from LDA-transformed and raw Ss were strongly correlated. Type I and II errors were
52 similar for both methods. The faster probability computation that we observed (speed gain
53 around a factor 100 for LDA-transformed Ss) substantially increases the ability of ABC
54 practitioners to analyze large numbers of *pods* and hence provides a manageable way to
55 empirically evaluate the power to discriminate among a large set of complex scenarios.

56

57 **Introduction**

58

59 One prospect of current biology is that molecular data will help us to reveal the complex
60 demographic processes that have acted on natural populations. The extensive availability
61 of various molecular markers and increased computer power have promoted the
62 development of inferential methods and associated softwares (e.g. Beaumont & Rannala
63 2004; Excoffier & Heckel 2006). Among these novel methods, Approximate Bayesian
64 Computation method (ABC; Beaumont *et al.* 2002) is increasingly used to make inferences
65 from large datasets for complex models in population and evolutionary biology (e.g.
66 Estoup *et al.* 2004; Fagundes *et al.* 2007; Jakobsson *et al.* 2006; Rosenblum *et al.* 2007;
67 Neuenschwander *et al.* 2008; Toni *et al.* 2009; Verdu *et al.* 2009; Bazin *et al.* 2010; Estoup
68 & Guillemaud 2010; Ascundes *et al.* 2011). The use of ABC techniques has also been
69 envisaged and successfully processed in other research fields, such as infectious disease
70 epidemiology (e.g. Luciania *et al.* 2009) and systems biology (e.g. Ratmann *et al.* 2009).

71 General statistical features, practical aspects, and applications of ABC in
72 evolutionary biology have been reviewed in at least three recent papers (Bertorelle *et al.*
73 2010; Csilléry *et al.* 2010; Beaumont 2010). Briefly, ABC constitutes a recent approach to
74 carrying out model-based inference in a Bayesian setting in which model likelihoods are
75 difficult to calculate (due to the complexity of the models considered) and must be
76 estimated by massive simulations. In ABC, the posterior probabilities of different models
77 and/or the posterior distributions of the demographic parameters under a given model are
78 determined by measuring the similarity between the observed dataset (i.e. the target) and a
79 large number of simulated datasets; all raw datasets (i.e. multilocus genotypes or individual
80 sequences) are summarized by so called summary statistics (Ss). Examples of such Ss in
81 population genetics are the mean number of alleles or heterozygosity per population and
82 F_{ST} or genetic distances between pairs of populations. In practice, ABC users can base
83 their analysis on simulation programs and then use statistical software to post-process their

84 simulation outputs. Several ABC programs have recently been developed to provide non-
85 specialist users with integrated solutions. They vary in the extent to which they are user-
86 friendly and they can be used for both data simulation and some post-processing steps (see
87 Table 1 in Bertorelle *et al.* 2010).

88 Although the methodology presented here is of more general interest, the present
89 work focuses on population genetics applications and applies to the model choice question.
90 In this context, models are evolutionary scenarios for which relative supports are compared
91 through their posterior probabilities. Choosing among a finite set of scenarios is crucial
92 when doing inferences about evolutionary history and processes for at least two reasons: (i)
93 it allows making general conclusions about major evolutionary events (e.g. admixture
94 between populations, occurrence of bottleneck events or identification of source
95 populations) and (ii) it makes it possible to estimate posterior probabilities of parameters
96 assuming a single scenario if the later is strongly supported (see the reviews of Bertorelle
97 *et al.* 2009, Csilléry *et al.* 2010 and Estoup & Guillemaud 2010 for various illustrations
98 regarding model choice). When processing ABC analyses, all the models are generally
99 simulated the same number of times. This is equivalent to giving the same prior probability
100 to each model under comparison and zero probability to any other model. In the final set of
101 retained simulations (those that have S_s close to the target's), the datasets produced by the
102 more supported models will be overrepresented and the datasets produced by other models
103 will be under-represented or even absent. Intuitively, the probability of a model is related
104 to the relative frequency of the datasets it produces that are among the retained simulations
105 (Weiss & von Haeseler 1998; Pritchard *et al.* 1999). This frequency may be taken as an
106 estimate of the posterior probability of a model, but this estimate is rarely accurate in
107 complex models when, inevitably, the retained simulations are either too few or also
108 contain datasets not closely matching the observed data (e.g. Guillemaud *et al.* 2010).

109 Recently, Leuenberger & Wegmann (2010) proposed the use of a parametric General
110 Linear Model to adjust the model frequencies in the retained simulations. However, the
111 most used and tested method, also available in integrated ABC packages such as DIYABC
112 (Cornuet *et al.* 2008, 2010), is the adjustment based on the polychotomous logistic
113 regression introduced by Beaumont (2008) (see also Fagundes *et al.* 2007; Cornuet *et al.*
114 2008). The coefficients for the regression between a model indicator (response) variable
115 and the simulated Ss (the explanatory variables) can be estimated, allowing the estimation
116 of the posterior probability for each model at the intercept condition where observed and
117 simulated Ss coincide. Confidence intervals (i.e. 95 % CI) of the probabilities can be
118 computed as suggested by Cornuet *et al.* (2008).

119 Large numbers of populations and loci can be analysed with ABC, and there is no
120 limit to the number and complexity of the models (hereafter named scenarios) considered.
121 However, several issues arise when the number of populations becomes too large. The
122 number of Ss to be manipulated increases considerably with the number of populations.
123 This is especially true when different types of markers requiring different types of Ss are
124 considered in the same analysis. A too large number of Ss may be of concern because ABC
125 algorithms attempt to sample from a small multidimensional sphere around the observed
126 statistics. The more Ss, the more difficult it becomes to match the observations closely and
127 increasing the number of simulations may not be sufficient to deal with this issue
128 (Beaumont *et al.* 2002). This phenomenon, which may potentially degrade the estimations
129 of posterior distributions of demo-genetic parameter as well as those of model posterior
130 probabilities, is often referred to as the “curse of dimensionality” (e.g. Beaumont *et al.*
131 2002; Blum & François 2009). There may be also a problem of co-linearity among
132 explanatory variables (Ss) resulting in instability of the regression when (too) many Ss are
133 introduced (Besley *et al.* 2004; Bazin *et al.* 2010). Recent improvements of ABC get round

134 these problems by using dimension reduction techniques, including a non-linear feed-
135 forward neural network (Blum & François, 2009) and partial least squares (PLS)
136 regression (Wegmann *et al.* 2009; see also Bazin *et al.* 2010). At least some algorithms of
137 this type have been implemented in the package ABCtoolbox (Wegmann *et al.* 2010). The
138 added value of such algorithms in the context of complex models and large datasets
139 remains, however, to be thoroughly tested (Bertorelle *et al.* 2010). Most importantly,
140 although the model itself can be considered as an additional parameter to infer, the PLS
141 dimension reduction technique applies to a continuous response variable. Therefore, this
142 technique can be applied to the estimation of posterior distributions of demographic and
143 genetic parameters under a given model and not to the computation of posterior
144 probabilities of models, the latter corresponding to a discrete response variable. Initially
145 developed for the estimation of posterior distributions of demographic and genetic
146 parameters, neural networks might theoretically be applied to model choice (Ripley 1996),
147 but, to our knowledge, this has not been tested and achieved in practice, at least in the
148 context of complex models and large datasets.

149 Robert *et al.* (2011) have shown that, because ABC algorithms involve an unknown
150 loss of information induced by the use of insufficient summary statistics, the conclusions
151 drawn on model comparison cannot be trusted per se and required further simulation
152 analyses. As pointed by Bertorelle *et al.* (2010) and Robert *et al.* (2011) among others,
153 confidence in model choice may be nevertheless empirically evaluated by processing
154 Monte Carlo evaluation of false allocation rates (type I and II errors) based on ABC
155 posterior probabilities computed from simulated pseudo-observed datasets. A version of
156 this exploratory analysis is already provided in the DIYABC software (Cornuet *et al.* 2008,
157 2010). This evaluation, based on the simulation and analysis of pseudo-observed datasets
158 (hereafter named *pods*), represents a useful and manageable quality assessment for

159 practitioners but is very time consuming. The polychotomous logistic regression used to
160 estimate scenario probabilities requires the computation of a matrix involving a very large
161 number of loops (i.e. [number of compared scenarios]² x [number of Ss]² x [number of
162 selected simulated datasets close to the target dataset]) at each iteration of the Newton-
163 Raphson method (Cornuet *et al.* 2008). This makes computation particularly time
164 consuming when the number of scenarios and Ss become large. Moreover, computations
165 involve several large matrices and probabilities that are sometimes simply not computable
166 when the computer memory space is not large enough. This is of particular concern when
167 type I and type II errors have to be computed from a large number of *pods*. As previously
168 stressed, such computations are nevertheless more and more requested by ABC experts for
169 assessing the power to discriminate among scenarios (e.g. Fagundes *et al.* 2007; Verdu *et*
170 *al.* 2009; Bertorelle *et al.* 2010; Lombaert *et al.* 2010; Robert *et al.* 2011).

171 In this paper, we describe a methodological innovation to process more efficient
172 ABC scenario probability estimation using linear discriminant analysis (LDA)
173 transformations on Ss before computing logistic regression. We first describe the principle
174 and goals of the method. We then use simulated *pods* to assess its main features (precision
175 and computation time) in comparison to probability estimation using logistic regression on
176 raw (i.e. not LDA-transformed) Ss. Finally, we illustrate the method on real microsatellite
177 datasets produced by Lombaert *et al.* (2011) to make inferences about the worldwide
178 routes of invasion of the coccinelid *Harmonia axyridis*.

179

180 **Materials and Methods**

181

182 *Linear discriminant analysis*

183
184 The Linear Discriminant Analysis (LDA) is a standard technique for supervised
185 classification. For a modern and comprehensive presentation of LDA, we invite readers to
186 refer to classical textbooks such as Ripley (1996), McLachlan (2004) or Hastie *et al.*
187 (2009). The LDA dates back to Fisher (1936) who proposed the dimension reduction
188 technique that contributed to the popularity of LDA. Actually, the classifier estimated with
189 the LDA depends only on some linear projection of the dataset onto a linear subspace
190 whose dimension is smaller than the number of groups, denoted by K . It is not our purpose
191 here to explain how this low-dimensional projection of the data can further leads to a LDA
192 classifier which provides automatic rules to classify a new data point to the class with the
193 largest posterior probability. As a matter of fact, we are here only interested in the
194 dimension reduction part of LDA and hence in the construction of the $(K - 1)$ discriminant
195 variables. Those discriminant variables are non-correlated, linear combinations of the
196 original variables that maximise the between-class variance relative to the within-class
197 variance, which is assumed identical among the different classes. This minimizes the
198 overlap between the classes when projected on the discriminant subspace if the within-
199 class distribution were Gaussian. Note that the discriminant variables are ordered with
200 respect to their ability to move the classes further apart.

201 In the methodological framework considered here (i.e. that of computing posterior
202 probabilities of scenarios using ABC), we used LDA to transform the set of usually large
203 number J of summary statistics (Ss) into $(K - 1)$ independent variables maximizing the
204 differences among the K compared scenarios (assuming $K < J$). The goal was to reduce the
205 dimension of the set of explanatory variables from J non-independent to $(K - 1)$
206 independent variables, and this whatever the value of J . Certainly, variance of the Ss varies
207 among the different scenarios. Even in that case, however, the projection onto the

208 discriminant subspace was proved relevant as a dimension reduction technique; see the
209 classical textbooks cited above. It is worth noting that we also weighted the simulated
210 datasets to give more importance to the ones that are closer to the observed dataset. The
211 LDA functions were used to transform both the (raw) simulated and observed Ss. Details
212 on LDA computations and transformation of Ss are given in the Appendix S1

213 Let us recapitulate how computation of the discriminant variables was included in
214 practice as a single additional step of the ABC process allowing the computation of
215 posterior probabilities of scenarios.

216 Step 1: We selected a subset of x % (typically 1 %) best simulations in a standard
217 ABC reference table (i.e. the table where parameter values drawn from priors and
218 corresponding simulated Ss have been recorded) usually including 10^6 simulations for each
219 of the K compared scenarios. This selection was based on the standard normalized
220 Euclidian distance computed between the observed and simulated “raw” (i.e. not
221 transformed) Ss (e.g. Beaumont *et al.* 2002) and hence corresponded to the x %
222 simulations with the smallest Euclidian distances.

223 Step 2 (LDA step; see Appendix S1 for details): we used LDA to transform the raw
224 Ss of this subset of x % best simulations into $(K - 1)$ discriminant variables maximizing
225 the differences among the K compared scenarios. When computing LDA functions, we
226 weighted the simulated datasets with the Epanechnikov kernel commonly used in the local
227 regression (equation 5 in Beaumont *et al.* 2002).

228 Step 3: We estimated the posterior probabilities of each competing scenario by
229 polychotomous logistic regression (Cornuet *et al.* 2008) on the x % best simulated datasets
230 now summarized by $(K - 1)$ discriminant variables instead of J non-independent variables
231 (i.e. raw Ss statistics). Confidence intervals (i.e. 95% CI) were computed for each posterior
232 probability using the $(s - 1)$ independent variables following Cornuet *et al.* (2008).

233 Hence, our proposal included only a single additional step (i.e. Step 2) when
234 compared to the computation traditionally proposed by different authors (e.g. Beaumont
235 2008; Fagundes *et al.* 2007; Cornuet *et al.* 2008 & 2010). Processing Step 2 substantially
236 decreases the number of explanatory variables through the production of LDA variables
237 maximizing the differences among the compared scenarios. This provides three main
238 advantages. First, computation of scenario probabilities using the polychotomous
239 regression of Step 3 becomes (much) faster and sometimes simply feasible. Second, a
240 lower number of explanatory variables may also improve the accuracy of the ABC
241 approximation, particularly when the number of simulations is not large enough to offset
242 the number of Ss. Finally, using LDA-transformed Ss avoids correlations among
243 explanatory variables.

244

245 *Tests on simulated datasets*

246 Pseudo-observed datasets (*pods*) were simulated from a set of known scenarios and prior
247 distributions to compare posterior probabilities obtained through the logistic regression
248 performed on both LDA-transformed and raw Ss. The *pods* were defined to mimic the real
249 microsatellite dataset of the ABC analysis 1 processed by Lombaert *et al.* (2011) on the
250 invasive coccinellid *Harmonia axyridis*. The *pods* hence included 18 microsatellites
251 genotyped in five population samples (18 to 35 individuals per population samples). This
252 dataset was produced to make inferences about the origin of the invasive *H. axyridis*
253 population established in Eastern North America in 1988 (ENA), considering altogether
254 two populations from the native range, two strains used for biocontrol release and one
255 (target) population from the introduction range (ENA). In this analysis, Lombaert *et al.*
256 (2011) defined ten competing scenarios considering a native or biocontrol population as a

257 source for ENA or admixture between them (see Lombaert *et al.* 2010 and 2011 for
258 details).

259 As in analysis 1 of Lombaert *et al.* (2011), genetic variation within and between
260 populations was summarized in the *Pods* using a set of (raw) statistics traditionally
261 employed in ABC (Cornuet *et al.* 2008 & 2010; Guillemaud *et al.* 2010). For each
262 population and each population pair we used the mean number of alleles per locus, the
263 mean expected heterozygosity and the mean allelic size variance. The other statistics used
264 were the mean ratio of the number of alleles over the range of allele sizes, pairwise F_{ST}
265 values, mean individual assignment likelihoods of population i assigned to population j and
266 the maximum likelihood estimate of admixture proportion. The total number of Ss was 86.

267 We choose this particular scenarios-priors-Ss setting because it had the potential to
268 fairly illustrate our new methodological developments based on LDA-transformed Ss. This
269 setting was characterized by relatively high (mean) type I error rates (ca. 0.40, due to the
270 large prior parameter space used to generate *Pods*, this space including “areas” for which
271 the discrimination among scenarios was difficult) and relatively small (mean) type II error
272 rates (ca. 0.07). High type I error rates corresponds to situations where probability values
273 of the target scenario can be small to high depending on the parameter values of the
274 analysed *pod*, hence virtually including the all spectrum of probabilities between 0 and 1.
275 This allows a better (and fairer) comparison of results between raw and LDA-transformed
276 Ss (cf. it is difficult to compare probability estimations when all values are between say
277 0.95 and 1.0). Moreover, this particular setting was chosen because it corresponded to
278 complex evolutionary models and large datasets that nevertheless could be analyzed for a
279 large number of *Pods* using logistic regression on both LDA-transformed and raw Ss. More
280 complex data and scenario settings (with larger number of scenarios and/or raw Ss) were

281 computationally too heavy to obtain probability estimations on a large enough number of
282 *Pods* in a manageable time using logistic regression on raw Ss “(i.e. < 15 min per *pod* on a
283 single standard biprocessor computer; see below). The results presented here were however
284 qualitatively similar to those obtained considering various alternative settings (with smaller
285 or larger numbers of scenarios and/or raw Ss) that we have also tested with our
286 methodological innovation (results not shown).

287 The ABC analyses of the *Pods* were performed using parameter values drawn from
288 the prior distributions described in Table S1 and by simulating 10^6 datasets for each of the
289 ten competing scenarios. For each *pod* we estimated the posterior probabilities of the
290 scenarios using a polychotomous logistic regression on the 1% of simulated datasets
291 closest to the observed dataset, considering either LDA-transformed or raw Ss.

292 We produced a first set of 500 *Pods* under scenario 5 (the scenario selected after
293 ABC treatment by Lombaert *et al.* 2011), drawing parameters values into the distributions
294 described in table S1. This scenario 5 is presented graphically in figure S1; the nine other
295 competing scenarios correspond to alternative source(s) of the target introduced population
296 (see Lombaert *et al.* 2011 for details). For each *pod*, we used the logistic regression on
297 either the 9 LDA-transformed or the 86 raw Ss to estimate the posterior probability and
298 95% CI of scenario 5 relatively to the set of ten compared scenarios. The number of
299 iterations of the Newton-Raphson algorithm used by the logistic regression computations
300 and the mean time of each iteration were also recorded for each *pod*.

301 We then produced a second set of 1,000 *Pods* including 10 subsets of 100 *Pods*
302 simulated under each of the ten compared scenarios, drawing parameter values from the
303 same distributions (Table S1). Each *Pods* subset was used to estimate type I and type II

304 errors on scenario choice using either the 9 LDA-transformed or the 86 raw Ss. Type I
305 error of a given scenario is the proportion of *pods* simulated from this scenario for which
306 this scenario does not have the highest posterior probability. Type II error is the proportion
307 of *pods* for which the scenario with the highest posterior probability is not the given true
308 one.

309 Finally, we evaluated the impact of the dimensionality of the simulated dataset (i.e.
310 the “curse of dimensionality” mentioned in the *Introduction* section), using either the 9
311 LDA-transformed or the 86 raw Ss. For different amount of simulated datasets, we
312 estimated the type I and II error rates from 500 *pods* simulated under scenario 5 (type I
313 error for scenario 5) and 500 *pods* simulated under scenario 1 (type II error for scenario 5
314 which in this case corresponds to the proportion of times that scenario 5 was selected when
315 *pods* have been produced under scenario 1). Scenario 1 was chosen to evaluate type II
316 errors because this scenario has shown the largest type II errors in the abovementioned
317 analyses. To consider different dimensionalities of simulated datasets, we decreased the
318 number of datasets simulated for each of the ten compared scenarios from 10^6 to 10^4 ,
319 keeping the proportions of datasets closest to the observed dataset selected for the logistic
320 regression at 1% of the total number of simulated datasets

321 All analyses were processed on a 2 CPU Intel Xeon X5472 computer (Windows XP
322 platform, 32 bits system, 4 Go of RAM) using a modified version of the package DIYABC
323 V1. This modified version is available under request from AE. LDA-transformation of Ss
324 before logistic regression will be implemented in a new multiplatform version of DIYABC
325 that will be freely available later in 2012.

326

327 *Tests on real datasets*

328 We used the real microsatellite datasets of Lombaert *et al.* (2011) to compare scenario
329 choice and probability estimation computing logistic regression on both LDA-transformed
330 and raw Ss. These datasets, which included 18 microsatellites genotyped on five to eight
331 population samples (18 to 42 individuals per population samples), were used to make five
332 consecutive ABC analyses about the worldwide routes of invasion of the coccinellid *H.*
333 *axyridis*, considering altogether populations from the native range, the introduction range
334 and biocontrol release actions, with potential admixture between them (see Lombaert *et al.*
335 2010 and 2011 for details).

336 We used prior distributions and Ss identical to those described in the previous
337 section (*Tests on simulated datasets*; Table S1). Following Lombaert *et al.* (2010 and
338 2011), we performed five consecutive ABC analyses of invasion scenarios involving
339 successive *H. axyridis* outbreaks that were successively recorded in the invaded range. As
340 previously detailed, analysis 1 dealt with the introduction pathway for the first recorded
341 outbreak in eastern North America in 1988, defining ten competing scenarios. Analysis 2
342 dealt with the second outbreak recorded in western North America in 1991, taking into
343 account the scenario selected in analysis 1, hence defining 15 competing scenarios. The
344 European and South American outbreaks in 2001 were addressed in analyses 3 and 4,
345 respectively (15 scenarios for each outbreak), taking into account the scenario selected in
346 analysis 1 and 2. Finally, the African outbreak in 2004 was considered in analysis 5 (28
347 scenarios), taking into account the scenarios selected in analyses 1, 2, 3 and 4. The total
348 number of raw Ss varied from 86 (analysis 1) to 223 (analysis 5), whereas the total number
349 of LDA-transformed Ss varied from 9 (analysis 1) to 27 (analysis 5).

350 The ABC analyses were performed by simulating 10^6 microsatellite datasets for
351 each competing scenario in the first four analyses and 5×10^5 datasets per scenario in
352 analysis 5 because of the high number of scenarios (28) and raw summary statistics (223)
353 which made a larger analysis computationally too heavy, even when using LDA-
354 transformed Ss. For each of the five analyses, we estimated the posterior probabilities of
355 the competing scenarios using a polychotomous logistic regression on the 1% of simulated
356 datasets closest to the observed dataset, considering either LDA-transformed or raw Ss.
357 Computation times were also recorded to illustrate the gain obtained in computation speed
358 when using LDA-transformed Ss.

359 Finally, we evaluated the impact of the number of simulated datasets recorded in
360 the reference table for analysis 1 on the estimation of the probability of scenario 5 using
361 either LDA-transformed or raw Ss. To this aim, we decreased the number of datasets
362 simulated for each of the ten compared scenarios from 10^6 to 10^4 , keeping the proportions
363 of datasets closest to the observed dataset selected for the logistic regression at 1% of the
364 total number of simulated datasets.

365 All analyses were processed on a 2 CPU Intel Xeon E5540 computer (Windows XP
366 platform, 32 bits system, 4 Go of RAM) using a modified version of the package DIYABC
367 V1 (available under request from AE).

368

369 **Results**

370

371 *Tests on simulated datasets*

372 Figure 1A illustrates the strong correlation between the probability values of scenario 5
373 obtained from *Pods* computing logistic regression on LDA-transformed Ss and raw Ss
374 (Pearson's correlation coefficient = 0.940). One can see, however, a trend for a globally
375 slightly lower scenario probability with LDA-transformed Ss (see linear regression
376 equation in the legend of Figure 1A). Figure 1B shows that 95% CI are almost always
377 smaller with LDA-transformed Ss.

378 Figure 2 summarizes the type I and II error rates obtained with LDA-transformed
379 and raw Ss. We found that these error rates substantially varied among scenarios but were
380 to a large extent similar for both methods for a given scenario. P-values computed using
381 Fisher exact test were higher than 0.6 for all scenarios for mean type II errors and were
382 lower than 5% for a single scenario for type I errors ($p = 0.047$ for scenario 7; p -value non
383 significant after applying the false discovery rate correction method of Benjamini &
384 Hochberg 1995).

385 The gain in computation time with LDA-transformed Ss was high. First, the
386 number of iterations needed to reach convergence during the logistic regression analysis
387 was lower with LDA-transformed Ss (mean = 7.320, SD = 1.420) than with raw Ss (mean
388 = 9.190, SD = 2.250). Second, the mean time of each such iteration was considerably
389 smaller with LDA-transformed Ss (mean = 7.034 sec, SD = 0.791) than with raw Ss (mean
390 = 888.146 sec, SD = 65.374). This translated into a computation speed increase by a mean
391 factor 128.128 (SD = 19.482) per iteration and 163.601 (SD = 46.456) for a completed
392 logistic regression analysis. The computation time for the LDA-transformation of raw Ss
393 before the regression was negligible.

394 Results summarized in Table 1 indicate that we did not face the curse of
395 dimensionality problem (see definition in the *Introduction* section) at least in the present
396 setting. Even for a large number of Ss and a strongly degraded number of simulated
397 datasets including only 10^4 datasets per scenario (total of 10^5 datasets for the ten compared
398 scenarios in this case), the error rates did not dramatically increase. The increase of type I
399 and II error rates with smaller datasets is (only) slightly faster for raw SS than for LDA-
400 transformed Ss.

401

402 *Tests on real datasets*

403 As will be further illustrated below on real datasets, our methodological innovation
404 is particularly attractive when practitioners have to deal with a large number of complex
405 scenarios involving a large number of Ss. Table 2 summarizes our results on scenario
406 choice and probability estimation computing logistic regression on both LDA-transformed
407 and raw Ss obtained on the real microsatellite datasets of Lombaert *et al.* (2011). For each
408 of the five consecutive analyses, the same scenario had the highest probability and was
409 hence selected using either LDA-transformed or raw Ss. The probabilities of the most
410 likely scenarios were slightly smaller with LDA-transformed Ss for analyses 1, 3 and 4,
411 and slightly larger for analysis 2. In contrast to computation based on LDA-transformed
412 Ss, analysis 5 could not be processed with raw Ss due to computer memory overflow. In all
413 analyses the 95% CI of the most likely scenario never overlapped those of competing
414 scenarios. As found with simulated *pods*, 95% CI with LDA-transformed Ss were smaller
415 than those with raw Ss.

416 In agreement with *pods* analyses, the gain in computation time with LDA-
417 transformed Ss was substantial. For all analyses, both the number of iterations needed to
418 reach convergence during the logistic regression and the mean computation time for each
419 such iteration was smaller with LDA-transformed Ss. This translated into a computation
420 speed increase by a factor 72 to 101 per iteration and 93 to 159 for a completed logistic
421 regression analysis.

422 Figure 3 indicates that analysis 1, processed either on LDA-transformed or raw Ss,
423 is rather robust to the potential difficulties associated with the curse of dimensionality.
424 Estimations of the probability of scenario 5 start to fluctuate substantially and 95% CIs to
425 increase considerably for simulation efforts including less than 2×10^5 datasets per
426 scenarios. No obvious differences could be observed between LDA-transformed and raw
427 Ss.

428

429 **Discussion**

430

431 Model comparison is an active research field among the widespread developments
432 currently undergone in ABC (e.g. Beaumont *et al.* 2009; Bertorelle *et al.* 2010; Csilléry *et*
433 *al.* 2010; Beaumont 2010; Robert *et al.* 2011). Here, we propose a methodological
434 innovation to deal with the discrimination among a large set of complex scenarios through
435 more efficient ABC probability computation using a linear discriminant analysis (LDA) on
436 Ss before the logistic regression analysis. Statistical methods to select appropriate Ss to
437 optimize model selection are still under development and discussed (see for instance

438 Fearnhead & Prangle 2012 and associated discussions). Our LDA-based transformation of
439 Ss represents a practical and straightforward way to tackle this question.

440 We show, using both simulated and real datasets, that posterior probabilities of
441 scenarios computed from LDA-transformed and raw Ss are strongly correlated. LDA-
442 transformed Ss tend, however, to provide slightly lower probability values and hence to be
443 somewhat conservative with respect to scenario discrimination. On the other hand, model
444 probabilities estimated from LDA-transformed Ss are characterized by smaller 95% CI.
445 The later feature is expected to decrease the number of inconclusive results if non-
446 overlapping of CI is taken as a criterion to select a scenario. When scenario selection is
447 made on the basis of the highest probability, type I and II errors were nevertheless similar
448 for both methods. The lower number of LDA variables used for the logistic regression
449 analysis (e.g. 9 LDA-transformed Ss versus 86 raw Ss in the *pods* we analyzed) is likely to
450 explain, to a large extent, both the smaller 95% CIs of probability estimates and the smaller
451 number of iterations needed to reach convergence during the regression.

452 A major practical advantage of using LDA-transformed Ss is that it substantially
453 decreases the dimension of explanatory variables making computation of scenario
454 probability (much) faster and sometimes simply feasible when the memory space is not
455 large enough to compute the matrix of second partial derivatives of the likelihood (p1 of
456 Supplementary material in Cornuet *et al.* 2008), as in Analysis 5 using the real dataset of
457 Lombaert *et al.* (2011). This allows larger data-scenarios settings to be analyzed. It is
458 worth stressing, however, that because LDA-transformation only plays on the number of
459 Ss and not on the number of parameters of the models, such transformation should not
460 motivate ABC practitioners to over-parameterize their models.

461 Faster probability computation increases the ability of ABC practitioners to analyze
462 large numbers of *pods* (for instance using the option “Evaluate confidence in scenario
463 choice” in the package DIYABC). It hence makes it easier to process a manageable
464 empirical evaluation of the power to discriminate among a given set of scenarios by
465 computing type I and II errors from sufficiently large number of *pods*, especially for large
466 sets of complex scenarios (see e.g. Robert *et al.* 2011 for theoretical arguments in favor of
467 such experimental explorations). Several authors have suggested to use scenario
468 probabilities computed from *pods* to evaluate type I and II errors to estimate the posterior
469 probability of a model among a set of k models given the observed posterior probability of
470 a real dataset, $P(M_k \text{ is the true model} \mid \text{observed estimated posterior probability} = x)$. Such
471 computation can then be used to adjust the posterior probabilities estimated from the real
472 dataset, taking part of the errors associated with ABC into account (see Fagundes *et al.*
473 2007; Lombaert *et al.* 2011).

474 Other potential advantages of LDA-transformation of raw Ss include reducing the
475 difficulties associated with the curse of dimensionality and avoiding correlation among
476 explanatory variables (i.e. multi-co-linearity) during the regression step. At least
477 theoretically the dimensionality issue might be offset by increasing the number of
478 simulations, but the amount of time then needed for concrete implementation might be
479 unreasonable. It is worth stressing, however, that the actual impact of such potential issues
480 remains difficult to assess in a generic manner as it probably differs depending on the
481 analyzed observed dataset, as well as on the Ss and/or scenario settings. Table 1 and Figure
482 3 both indicate a good robustness to the numbers of simulated datasets, as a substantial
483 effect could be observed only for particularly low (and in practice rarely used) number of
484 simulated datasets. Analyses carried on *pods* suggest a slightly better robustness when

485 using LDA-transformed rather than raw Ss, at least when using type I and II error rates as
486 criterion (cf. the slightly smaller increase of errors with smaller datasets for LDA-
487 transformed than raw Ss). It is difficult to know, however, to which extent this result
488 reflects the lower number of LDA variables used for the regression and/or the fact that a
489 substantial number of raw Ss are non-independent variables.

490 We believe that our LDA-based methodological innovation will usefully enlarge
491 the tool box available to biologists to make ABC inferences on more complex and hence
492 more realistic demographic processes that have acted on natural populations.

493

494 **Acknowledgements**

495 This research was financially supported by the French Agence Nationale de la Recherche
496 grants ANR-09-BLAN-0145-01 - EMILE to all authors. We thank the EMILE working
497 group and Eric Bazin for useful discussions.

498

499 **References**

500

501 Ascunce MS, Yang CC, Oakey J, Calcaterra L, Wu WJ, Shih CJ, Goudet J, Ross KG, and
502 Shoemaker D (2011) Global invasion history of the fire ant *Solenopsis invicta*.
503 *Science*, **331**, 1066-1068.

504 Bazin E, Dawson KJ, Beaumont MA (2010) Likelihood-Free inference of population
505 structure and local adaptation in a Bayesian hierarchical model. *Genetics*, **185**, 587-
506 602.

507 Beaumont MA (2008) Joint determination of topology, divergence time and immigration in
508 population trees. In: Simulation, Genetics and Human Prehistory (eds Matsamura S,
509 Forster P, Rrenfrew C), pp. 135–154. McDonald Institute for Archaeological
510 Research, Cambridge, UK.

511 Beaumont M (2010) Approximate Bayesian Computation in Evolution and Ecology.
512 *Annual Review of Ecology and Evolutionar. Systematics*, **41**, 379-406.

513 Beaumont M, Rannala B (2004) The Bayesian revolution in genetics. *Nature Review*
514 *Genetics*, **5**, 251-261.

515 Beaumont MA, Zhang WY, Balding DJ (2002) Approximate Bayesian computation in
516 population genetics. *Genetics*, **162**, 2025-2035.

517 Beaumont MA, Cornuet J-M, Marin J-M, Robert CP (2009) Adaptivity for ABC
518 algorithms: the ABC-PMC scheme. *Biometrika*, **96**, 983–990.

519 Benjamini Y, Hochberg Y (1995) Controlling the false discovery rate: a practical and
520 powerful approach to multiple testing. *Journal of the Royal Statistical Society B*,
521 **57**, 289-300.

522 Bertorelle G, Benazzo A, Mona S (2010) ABC as a flexible framework to estimate
523 demography over space and time: some cons, many pros. *Molecular Ecology*, **19**,
524 2609-2625.

525 Besley DA, Kuh E, Welsch RE (2004) Regression diagnostics: identifying influential data
526 and sources of collinearity. Publish by Wiley & Sons, Inc., Hoboken, New Jersey.

527 Blum MGB, François O (2009) Non-linear regression models for Approximate Bayesian
528 Computation. *Statistics and Computing*, **20**, 63-73.

529 Cornuet J-M, Santos F, Beaumont MA, Robert CP, Marin J-M, Balding DJ, Guillemaud T,
530 Estoup A. (2008) Inferring population history with DIY ABC: a user-friendly
531 approach to approximate Bayesian computation. *Bioinformatics*, **24**, 2713-2719.

532 Cornuet J-M, Ravigne V, Estoup A (2010) Inference on population history and model
533 checking using DNA sequence and microsatellite data with the software DIYABC
534 (v1.0). *BMC Bioinformatics* **11**, 401doi:10.1186/1471-2105-11-401.

535 Csilléry K, Blum M, Gaggiotti O, François O (2010) Approximate Bayesian Computation
536 (ABC) in practice. *Trends in Ecology and Evolution*, **25**, 410-418.

537 Estoup A, Beaumont M, Sennedot F, Moritz C, Cornuet J-M (2004) Genetic analysis of
538 complex demographic scenarios: spatially expanding populations of the cane toad,
539 *Bufo marinus*. *Evolution*, **58**, 2021-2036.

540 Estoup A, Guillemaud T (2010) Reconstructing routes of invasion using genetic data: why,
541 how and so what? *Molecular Ecology*, **19**, 4113-4130.

542 Excoffier L, Heckel G (2006) Computer programs for population genetics data analysis: a
543 survival guide. *Nature Review Genetics*, **7**, 745-758.

544 Fagundes NJR, Ray N, Beaumont M, Neuenschwander S, Salzano FM, Bonatto SL,
545 Excoffier L (2007) Statistical evaluation of alternative models of human evolution.

546 *Proceedings of the National Academy of Sciences of the United States of America*,
547 **104**, 17614-17619.

548 Fearnhead P, Prangle D (2012) Constructing summary statistics for approximate Bayesian
549 computation: Semi-automatic approximate Bayesian computation, *Journal of the*
550 *Royal Statistical Society B*, **74**, 1-28.

551 Fisher RA (1936) The use of multiple measurements in taxonomic problems. *Annals of*
552 *Eugenics*, **7**, 179-188.

553 Guillemaud T, Beaumont M, Ciosi M, Cornuet J-M, Estoup A (2010) Inferring
554 introduction routes of invasive species using approximate Bayesian computation on
555 microsatellite data. *Heredity*, **104**, 88-99.

556 Hastie T, Tibshirani R, Friedman J (2009) The Elements of Statistical Learning Data
557 Mining, Inference, and Prediction, Second Edition Springer Series in Statistics,
558 Springer-Verlag, New York

559 Jakobsson M, Hagenblad J, Tavaré S, Säll T, Halldén C, *et al.* (2006) A recent unique
560 origin of the allotetraploid species *Arabidopsis suecica*: evidence from nuclear
561 DNA markers. *Molecular and Biological Evolution*, **23**, 1217-1231.

562 Leuenberger C, Wegmann D (2010) Bayesian Computation and model selection without
563 likelihoods. *Genetics*, **184**, 243-252

564 Lombaert E, Guillemaud T, Cornuet J-M, Malausa T, Facon B, *et al* (2010) Bridgehead
565 effect in the worldwide invasion of the biocontrol harlequin ladybird. *PLoS ONE*,
566 **5**, e9743, doi:10.1371/journal.pone.0009743.

567 Lombaert E, Guillemaud T, Thomas C, Lawson Handley L-J, Li J, *et al.* (2011) Inferring
568 the origin of populations introduced from a genetically structured native range by
569 approximate Bayesian computation: case study of the invasive ladybird *Harmonia*
570 *axyridis*. *Molecular Ecology*, **20**, 4654–4670.

571 Luciania F, Sisson SA, Jiang H, Francis AR, Tanaka MM (2009) The epidemiological
572 fitness cost of drug resistance in *Mycobacterium tuberculosis*. *Proceedings of the*
573 *National Academy of Sciences of the United States of America*, **106**, 14711-14715.

574 McLachlan GJ (2004) *Discriminant analysis and statistical pattern recognition*. Wiley
575 Interscience.

576 Neuenschwander S, Largiader CR, Ray N, Currat M, Vonlanthen P, Excoffier L (2008)
577 Colonization history of the Swiss Rhine basin by the bullhead (*Cottus gobio*):
578 Inference under a Bayesian spatially explicit framework. *Molecular Ecology*, **17**,
579 757-772.

580 Pritchard JK, Seielstad MT, Perez-Lezaun A, Feldman MW (1999) Population growth of
581 human Y chromosomes: A study of Y chromosome microsatellites. *Molecular*
582 *Biology and Evolution*, **16**, 1791-1798.

583 Ratmann O, Andrieu C, Wiuf C, Richardson S (2009) Model criticism based on likelihood-
584 free inference, with an application to protein network evolution. *Proceedings of the*
585 *National Academy of Sciences of the United States of America*, **106**, 10576-10581.

586 Ripley BD (1996) *Pattern Recognition and Neural Networks*, Cambridge University Press.

587 Robert CP, Cornuet J-M, Marin J-M, Pillai NS. (2011) Lack of confidence in
588 approximate Bayesian computation model choice. *Proceedings of the National*
589 *Academy of Sciences of the United States of America*, **108**, 15112-15117.

590 Rosenblum EB, Hickerson MJ, Moritz C (2007) A multilocus perspective on colonization
591 accompanied by selection and gene flow. *Evolution*, **61**, 2971-2985.

592 Toni T, Welch D, Strelkowa N, Ipsen A, Stumpf MPH (2009) Approximate Bayesian
593 computation scheme for parameter inference and model selection in dynamical
594 systems. *Journal of the Royal Society Interface*, **6**, 187-202.

595 Verdu P, Austerlitz F, Estoup A, Vitalis R, Georges M, *et al.* (2009) Origins and genetic
596 diversity of pygmy hunter-gatherers from Western Central Africa. *Current Biology*,
597 **19**, 312-318.

598 Wegmann D, Leuenberger C, Excoffier L (2009) Efficient approximate Bayesian
599 computation coupled with Markov chain Monte Carlo without likelihood. *Genetics*,
600 **182**, 1207– 1218.

601 Wegmann, D, Leuenberger C, Neuenschwander S, Excoffier L (2010) ABCtoolbox: a
602 versatile toolkit for approximate Bayesian computations. *BMC Bioinformatics*, **11**,
603 116.

604 Weiss G, von Haeseler A (1998) Inference of population history using a likelihood
605 approach. *Genetics*, **149**, 1539–1546.

606

607 **Figures legends**

608

609 **Figure 1. Probability estimations of scenario 5 computed using LDA-transformed or**
610 **raw summary statistics for 500 pods simulated under scenario 5 (ten scenarios**
611 **compared).**

612 Note: (A) Pearson's correlation coefficient between probability estimations = 0.940 (95%
613 CI = [0.928, 0.949]). Solid line: $y = x$; dotted line: linear regression line $y = 0.818436x +$
614 0.004878 . (B) 95% CIs (i.e. 2.5% and 97.5% quantiles) for each probability values
615 obtained from either LDA-transformed summary statistics (black lines) or raw summary
616 statistics (grey lines).

617

618 **Figure 2. Confidence in discriminating scenarios using LDA-transformed or raw**
619 **summary statistics.**

620 Note: Type I error: exclude scenario x when it is actually scenario x . Type II error: choose
621 scenario x when it is not scenario x . Results are based on 100 *pods* per scenario (total of
622 ten compared scenarios). The compared scenarios correspond to variants of the scenario 5,
623 the latter being detailed in Figure S1.

624

625 **Figure 3. Probabilities of scenario 5 computed from the real dataset of Lombaert et**
626 **al. (2011) for different numbers of simulated datasets.**

627 Note Black = LDA-transformed summary statistics. Grey = raw summary statistics. Plain
628 and dotted lines are for probability estimations and 95% CIs, respectively. Probabilities of
629 scenario 5 were estimated for number of datasets simulated for each of the ten compared
630 scenarios decreasing from 10^6 to 10^4 , keeping the proportions of datasets closest to the
631 observed dataset selected for the logistic regression at 1% of the total number of simulated
632 datasets.

633 **Table 1. Type I and II error rates estimated for different numbers of simulated**
 634 **datasets.**

635
 636

| | | Number of simulated datasets for each of the 10 compared scenarios | | | | |
|------------------|--------------------|-----------------------------------------------------------------------|--------|-----------------|-----------------|--------|
| | | 10^6 | 10^5 | 5×10^4 | 2×10^4 | 10^4 |
| Type I error | LDA-transformed Ss | 0.560 | 0.556 | 0.584 | 0.592 | 0.622 |
| | Raw Ss | 0.450 | 0.492 | 0.530 | 0.536 | 0.624 |
| Type II error | LDA-transformed Ss | 0.056 | 0.056 | 0.052 | 0.062 | 0.080 |
| | Raw Ss | 0.060 | 0.062 | 0.072 | 0.088 | 0.116 |

637
 638
 639
 640
 641
 642
 643

Note: Type I error rates were estimated for scenario 5 from 500 *pods*. Type II errors were estimated for scenario 5 when simulating 500 *pods* under scenario 1. The number of datasets simulated for each of the ten compared scenarios decreased from 10^6 to 10^4 , keeping the proportions of datasets closest to the observed dataset selected for the logistic regression at 1% of the total number of simulated datasets

Table 2. Scenario choice and posterior probability estimated from either LDA-transformed or raw summary statistics when considering the real microsatellite datasets of Lombaert et al. (2011).

| | Logistic regression on raw summary statistics | | | | Logistic regression on LDA-transformed summary statistics | | | | | | Speed gain |
|-------------------------------------------------------------------|-----------------------------------------------|---------------------|--------------------------------|----------------------------|-----------------------------------------------------------|-------------|---------------------|--------------------------------|----------------------------|---------------------|-------------------------------------|
| | Nb of stats | Selected scenario # | Posterior probability [95% CI] | Mean time per NR iteration | Nb of NR iterations | Nb of stats | Selected scenario # | Posterior probability [95% CI] | Mean time per NR iteration | Nb of NR iterations | Per iteration \ over all iterations |
| Consecutive ABC analyses (nb simulations per scen. \ nb of scen.) | | | | | | | | | | | |
| Analysis 1 (10 ⁶ \ 10 scenarios) | 86 | 5 | 0.6242 [0.5767, 0.6717] | 5' 05" | 11 | 9 | 5 | 0.5420 [0.5325, 0.5516] | 3" | 7 | 101.7 \ 159.8 |
| Analysis 2 (10 ⁶ \ 15 scenarios) | 124 | 1 | 0.4425 [0.3746, 0.5105] | 38' 45" | 11 | 14 | 1 | 0.5767 [0.5559, 0.5976] | 31" | 7 | 75.0 \ 117.9 |
| Analysis 3 (10 ⁶ \ 15 scenarios) | 124 | 13 | 0.8134 [0.7107, 0.9160] | 38' 38" | 9 | 14 | 13 | 0.7487 [0.7214, 0.7760] | 32" | 6 | 72.4 \ 93.1 |
| Analysis 4 (10 ⁶ \ 15 scenarios) | 124 | 4 | 0.9489 [0.9315, 0.9663] | 33' 36" | 9 | 14 | 4 | 0.9227 [0.9139, 0.9315] | 27" | 7 | 74.7 \ 96.0 |
| Analysis 5 (5x10 ⁵ \ 28 scenarios) | 223 | NC* | NC | NC | NC | 27 | 4 | 0.6864 [0.6456, 0.7272] | 6' 11" | 7 | NC |

Note: The probabilities of the competing scenarios were computed using a logistic regression on the 1% of simulated datasets closest to the real *Harmonia axyridis* datasets. NR iterations = Newton-Raphson iterations (Cornuet *et al.* 2008). NC: not computable. * Because the full computation of analysis 5 was not feasible (Lombaert *et al.* 2011), an alternative method was used to compare scenarios by first setting aside 11 scenarios using the direct approach (Cornuet *et al.* 2008) on the 0.01% datasets closest to the observed dataset. The scenario 4 was then selected due to its highest posterior probability in a subsequent analysis (using polychotomous logistic regression and raw Ss) performed on the 19 remaining scenarios (see Lombaert *et al.* 2011 for details).

Fig. 1

(A)

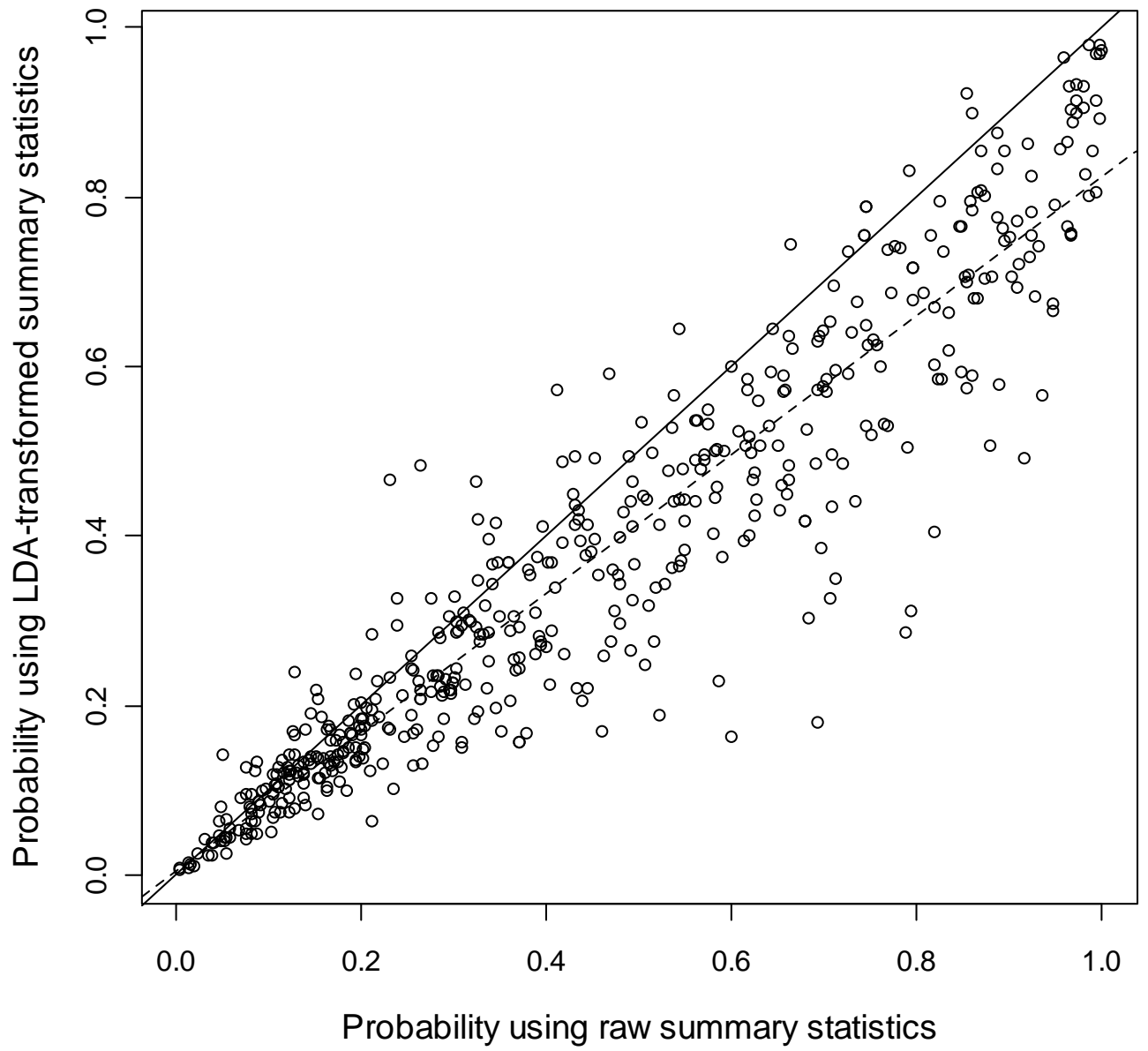


Fig. 1 continued

(B)

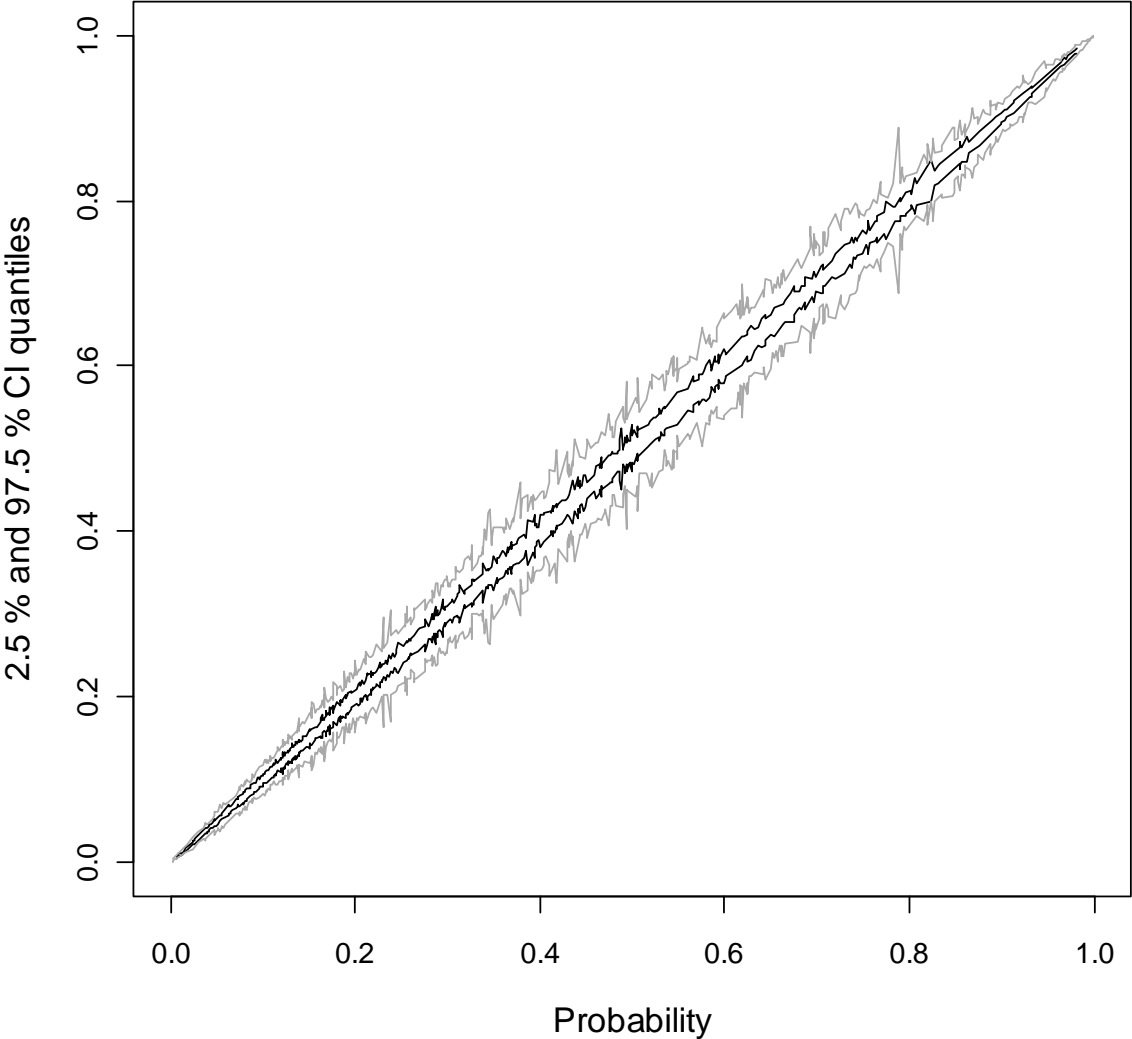


Fig. 2

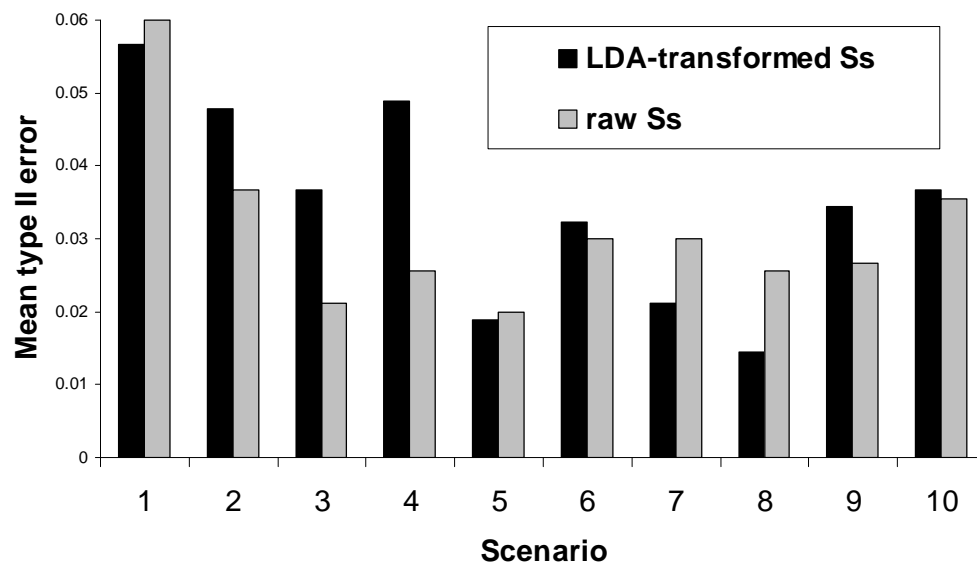
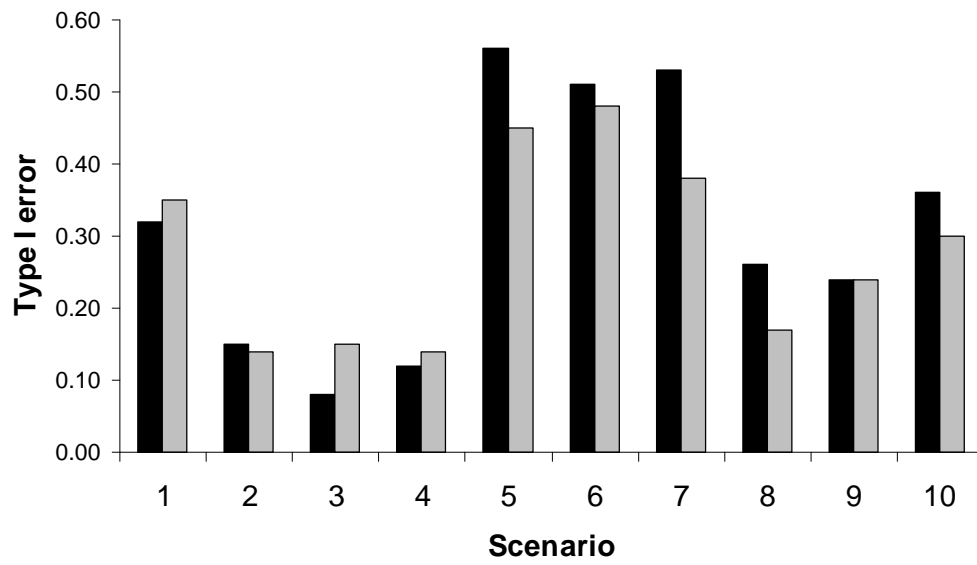


Fig. 3

