

1

DIYABC

2

version 2.0

3

A user-friendly software
for inferring population history through
Approximate Bayesian Computations

7

J.M. Cornuet, P. Pudlo, J. Veyssié,

6

A. Dehne-Garcia, A. Estoup and V. Ravigné
Centre de Biologie et de Gestion des Populations

8

Institut National de la Recherche Agronomique
Campus International de Baillarguet, CS 30016 Montferrier-sur-Lez
34988 Saint-Gély-du-Fesc Cedex, France
(diyabc@cbgp.supagro.fr)

December 18, 2012

Contents

2	1.	Preface	3
3	1.1	Acknowledgements	4
4	1.2	References to cite	4
5	1.3	Web site	4
6	2.	Methodology	5
7	2.1	Basic notions on ABC	5
8	2.2	Historical model parameterization	5
9	2.3	Mutation model parameterization (microsatellite and DNA sequence loci)	10
10	2.4	SNPs do not require mutation model parameterization	11
11	2.5	Prior distributions	11
12	2.6	Algorithms for data simulation : main features	11
13	2.7	Summary statistics	12
14	2.8	Pre-evaluation of scenarios and prior distributions	14
15	2.9	Estimation of posterior distributions of parameters	14
16	2.10	Model checking	15
17	2.11	Measures of performances	15
18	2.12	Comparison of scenarios	16
19	3.	The Graphic User Interface	18
20	3.1	What is a <i>DIYABC</i> Project ?	18
21	3.2	Options of the home screen	18
22	3.3	Defining a new project	19
23	3.4	Building the reference table	28
24	3.5	Performing analyses	29
25	3.6	Simulating data sets	46
26	3.7	The <i>Settings</i> option of the File menu	59
27	4.	Implementation details	65
28	4.1	Software design	65

1	4.2	Files	65
2	4.3	Missing data	68
3	4.4	Data files	68
4	5.	Cluster version	71
5	5.1	A note about the random number generator used in DIYABC	73

1. Preface

In less than 10 years, Approximate Bayesian Computations (ABC) have developed in the Population Genetics community as a new tool for inference on the past history of populations and species. Compared to other approaches based on the computation of the likelihood which are still restrained to a very narrow range of evolutionary scenarios and mutation models, the ABC approach has demonstrated its ability to stick to biological situations that are much more complex and hence realistic. However, this approach still requires numerous computations to be performed so that it has been used mostly by specialists (i.e. statisticians and programmers). This has almost certainly restrained the possible impact of ABC in population genetic studies. We believe that this situation must be improved and therefore we have developed a computer program for the large community of experimental biologists. We therefore designed DIYABC as a user-friendly program allowing non specialist biologists to achieve their own analysis.

The first version (*DIYABCv0.x*) had been written especially for microsatellite data. There were at least two reasons for that. The first one is that we have been among the first to develop and use this class of markers in population genetic studies (e.g. Estoup *et al.*, 1993). Since then, we have developed microsatellites in numerous species as well as we have published theoretical studies and reviews on these markers (e.g. Estoup *et al.*, 2002). The second reason is that microsatellites have been and still are very popular markers in the population geneticist community and there is now a large quantity of data that might benefit of an ABC approach.

The second version of our software (*DIYABCv1.x*) has been designed to make use of DNA sequence data. This has several immediate consequences. For instance, the standard Genepop data file format has been extended to incorporate sequence data. This has been done in collaboration with the authors of *Genepop* and explained in subsection 4.1.1. In this version, sequence loci are considered in the same way as microsatellite loci, i.e. they are considered as genetically independent and intra-locus recombination is not (yet) available. Concerning mutation models for DNA sequences, we used the same philosophy as for microsatellites, i.e. the program considers only simple and widely used models, keeping in mind that a higher-dimensional parameter space will be less well explored than a lower-dimensional space. Note that none of these mutation models includes insertion-deletions. Also five categories of loci (either microsatellites or DNA sequences) were considered in this second version : autosomal diploid, autosomal haploid, X-linked, Y-linked and mitochondrial. Note that X-linked loci can be used for an haplo-diploid species in which both sexes have been sampled. If non-autosomal loci have been typed in population samples, the sex-ratio of the species will have to be provided (see subsection 4.1.1).

Other improvements over version 0 included :

1. the use of multithread technology in order to exploit multicore/multiprocessor computers. This is especially useful when building the reference table and for several other intensive computation steps, such as the multinomial logistic regression,
2. a new option which helps the detection of "bad" prior modelisation of the data,
3. another new option which helps evaluate the goodness of fit of a given model-parameter posterior combination (i.e. Model checking),
4. many new screens implemented not only to treat sequence data, but also to cope with the new options described above, as well as to offer useful complementary information on the current run.

The third version of *DIYABC* (*DIYABCv2.x*) has been entirely recoded in order to be used under the usual three OS (Windows, Mac and Linux). Also the code for computations has been separated from that of the graphic user interface (GUI). The former has been rewritten in C++ and the latter is a mixture of Python and Qt (PyQt). The user can then launch computations with or without using the GUI. The GUI's uses are :

1. the management of projects
2. the input of the historical and genetical models
3. the parameterization of analyses
4. the launch of computations of the reference table and of the various required analyses
5. the visualization of results

1 Also, as DNA sequences have been added in the second version, a new category of markers has been
 2 added to the third version : Single Nucleotide Polymorphisms (SNPs). Instead of extending once more
 3 the Genepop format, a new data simple format has been designed for these markers.

4 This version includes all improvements of version 1.x and a few new improvements such as :

- 5 • loci of the same type (i.e. microsatellites on one hand or DNA sequences on the other hand) can
 6 be associated in one or more groups. This allows for instance to define different mutation models
 7 for microsatellites with motifs of different lengths.
- 8 • the model checking option is now presented as a direct option (not a suboption of the ABC esti-
 9 mation of parameters) which largely simplifies its use.
- 10 • the logistic regression can be performed on factorial discriminant analysis components instead of
 11 all summary statistics. This reduces the number of dependent variables, thus allowing to run large
 12 "confidence in scenario choice" analyses including many summary statistics and scenarios.
- 13 • ascertainment bias in the design of SNPs can be tentatively corrected by considering "reference"
 14 samples in which the loci need to be polymorphic in order to belong to the SNP set. These reference
 15 samples are not necessarily included of the actual samples.

16 1.1 Acknowledgements

17 We thank Mark Beaumont who has been at the origin of our interest for ABC. He offered us constant
 18 help and inspiration since the beginning. We also thank David Balding who welcomed one of us (JMC)
 19 in his team during the whole writing of the program and who organized several workshops on ABC
 20 during the same period. We are indebted to Christian Robert, Jean-Michel Marin, Stuart Baird, Thomas
 21 Guillemaud, Renaud Vitalis, Gael Kergoat, Gilles Guillot and David Welsh with whom we discussed
 22 many theoretical and practical aspect of DIYABC in the numerous meetings financed by a grant from
 23 the French Research National Agency (project *MISGEPOP* ANR-05-BLAN-196). The same grant is
 24 also aknowledged for having paid for the 2-year salary of FS. This research was also supported by an EU
 25 grant awarded to JMC as an EIF Marie-Curie fellowship (project *StatInfPopGen*) and which allowed
 26 him to come to David Balding's place at Imperial College (London, UK). Current and future developments
 27 of DIYABC are financed by a new grant from the French Research National Agency (project *EMILE*
 28 ANR-09-BLAN-0145) awarded in september 2009.

29 1.2 References to cite

- 30 • **version 0** : Cornuet J.M., F. Santos, M.A. Beaumont, C.P. Robert, J.M. Marin, D.J. Balding, T.
 31 Guillemaud and A. Estoup. Inferring population history with DIYABC: a user-friendly approach
 32 to Approximate Bayesian Computations (2008) *Bioinformatics*, **24** (23), 2713-2719.
- 33 • **version 1** : Cornuet J.M., V. Ravigné and A. Estoup, 2010. Inference on population history
 34 and model checking using DNA sequence and microsatellite data with the sofware DIYABC (v1.0)
 35 (2010) *BMC Bioinformatics*, **11**, 401.
- 36 • **version 2** : Cornuet J.M., Pudlo P., Veyssier J., Dehne-Garcia A., V. Ravigné and A. Estoup.
 37 Improved inference on population history using SNP's and DIYABC (v2). *submitted*

38 1.3 Web site

39 <http://www.montpellier.inra.fr/CBGP/diyabc>

40

¹ 2. Methodology

² 2.1 Basic notions on ABC

³ Approximate Bayesian Computation or ABC is a bayesian approach in which the posterior distributions
⁴ of the model parameters are determined by replacing the computation of the likelihood (probability of
⁵ observed data given the values of the model parameters) by a measure of similarity between observed
⁶ and simulated data. The posterior distributions are estimated from parameter values providing simulated
⁷ data that are the most similar to observed data. Historically, different ways of estimating this similarity
⁸ have been proposed, but all have been based on statistics summarizing information conveyed by the data
⁹ set. In population genetics, data most often relate to individuals that have been genotyped at a given set
¹⁰ of loci, these individuals being representative of the studied populations. The summary statistics are for
¹¹ instance the mean number of alleles per population or genetic distances between pairs of populations. It is
¹² much easier to measure the similarity between small sets of summary statistics than between large sets of
¹³ multilocus genotype data. When the number of summary statistics is low, it is possible to select simulated
¹⁴ data for which *all* the summary statistics are close to those of the observed data (Pritchard *et al.*, 1999;
¹⁵ Estoup *et al.*, 2001; Estoup and Clegg, 2003). However, for more complex scenarios necessitating a larger
¹⁶ number of summary statistics, it becomes almost impossible to find such simulated data sets. Beaumont
¹⁷ *et al.* (2002) have hence proposed to measure similarity through the Euclidian distance between observed
¹⁸ and simulated summary statistics, after normalization by standard deviations of simulated statistics. In
¹⁹ addition, these authors introduced a step of local linear regression aimed at favoring simulated data sets
²⁰ that are closer to the observed one.

²¹ In practice, the ABC approach can be summarized in three successive steps (Excoffier *et al.*, 2005) :
²² i) generating simulated data sets, ii) selecting simulated data sets closest to observed data set and iii)
²³ estimating posterior distributions of parameters through a local linear regression procedure.

²⁴ In addition, this approach provides a way of comparing different scenarios that can explain observed
²⁵ data. Two measures of posterior probabilities of scenarios are proposed. The first measure is simply the
²⁶ relative proportion of each scenario in the simulated data sets closest to observed data sets (Miller *et*
²⁷ *al.*, 2005; Pascual *et al.*, 2007). The second measure is obtained by a logistic regression of each scenario
²⁸ probability on the deviations between simulated and observed summary statistics (Fagundes *et al.*, 2007;
²⁹ Beaumont, 2008).

³⁰

³¹ In order to simulate data, one has first to define one (or possibly several) model(s). Each model
³² includes a historical model describing how the sampled populations are connected to their common an-
³³ cestor and a mutational model describing how allelic states of the studied genes are changing along their
³⁴ genealogical trees.

³⁶ 2.2 Historical model parameterization

³⁷ The evolutionary scenario, which is quantified by the historical model, can be described as a succession in
³⁸ time of "events" and "inter event periods". The events considered in the program are a restricted set of
³⁹ possible events affecting populations evolution. In the current version of the program, we consider only 4
⁴⁰ categories of events : population divergence, discrete change of effective population size, admixture and
⁴¹ sampling (the last one has been added to allow considering samples taken at different times). Between two
⁴² successive events affecting a population, we assume that populations evolve independently (e.g. without
⁴³ migration) and with a fixed effective size. The usual parameters of the historical model are the times
⁴⁴ of occurrence of the various events (counted in generations), the effective sizes of populations and the
⁴⁵ admixture rates. When writing the scenario, events are provided sequentially backward in time. Although
⁴⁶ this choice may not be natural at first sight, it is coherent with coalescence theory on which are based all
⁴⁷ data simulations in the program. For that reason, the keywords for a divergence or an admixture event
⁴⁸ are `merge` and `split`, respectively. Two other keywords, `varNe` and `sample`, correspond to a discrete
⁴⁹ change in effective population size and a gene sampling, respectively.

⁵⁰ A scenario takes the form of a succession of lines (one line per event), each line starting with the time of
⁵¹ the event, then the nature of the event, and ending with several other data depending on the nature of
⁵² the event. Following is the syntax used for each category of event :

⁵³ **population sample** : $\langle \text{time} \rangle \text{ sample } \langle \text{pop} \rangle [\text{nmales } \text{n females}]$

⁵⁴ $\langle \text{time} \rangle$ is the time (always counted in number of generations) at which the sample was taken and

`<pop>` is the population number from which is taken the sample. It is worth stressing here that **samples are considered in the same order as they appear in the data file**.

`[nmales nfemales]` is only used for SNP loci to indicate the number of males (respectively females) from the sample that have been used to detect SNPs. These males and females appear in the corresponding sample of the data file.

6 **population size variation** : `<time> varNe <pop> <Ne>`

From time `<time>`, looking backward in time, population `<pop>` will have an effective size `<Ne>`.

8 **population divergence** : `<time> merge <pop0> <pop1>`

At time `<time>`, looking backward in time, population `<pop1>` "merges" with population `<pop0>`. Hereafter, only `<pop0>` "survives".

11 **population admixture** : `<time> split <pop0> <pop1> <pop2> <rate>`

At time `<time>`, looking backward in time, population `<pop0>` "splits" between populations `<pop1>` and `<pop2>`. A gene lineage from population `<pop0>` joins population `<pop1>` (respectively `<pop2>`) with probability `<rate>` (respectively `1-<rate>`).

A scenario is a succession of lines as described above. However, in order to cope with special situations (see explanations in Note 9 below), we added a first line giving the effective sizes of sampled populations before the first event described, looking backward in time. Expressions between arrows, other than population numbers, can be either a numeric value (e.g. 25) or a character string (e.g. t0). In the latter case, it is considered as a parameter of the model. So the only possible parameters of the historical model are times of events, effective population sizes and admixture rates.

The program offers the possibility to add or remove scenarios, by just clicking on the corresponding buttons. The usual shortcuts (CTRL+C, CTRL+V and CTRL+X) can be used to edit the different scenarios. Some or all parameters can be in common among scenarios.

25 Notes

26 1. There are two ways of giving a fixed value to effective population sizes, times and admixture rates.
 27 Either the fixed value appears as a numeric value in the scenario windows or it is given as a string
 28 value like any parameter. In the latter case, one gives this parameter a fixed value by choosing a
 29 Unifom distribution and setting the minimum and maximum to that value in the prior setting (see
 30 section 2.4).

31 2. All expressions must be separated by at least one space.

32 3. All expressions relative to parameters can include sums or differences. For instance, it is possible
 33 to write :

```
t0 merge 2 3
t0+t1 merge 1 2
```

This means that `t1` is the time elapsed between the two events. By imposing `t1>0` (as explained in section **prior and posterior distributions**), this implies that the divergence of populations 1 and 2 is always more ancient than the divergence of populations 2 and 3. However, one cannot mix a parameter and a numeric value (e.g. `t1+150` will result in an error). This can be done by writing `t1+t2` and fixing `t2` by choosing a uniform distribution with lower and upper bounds both equal to 150.

42 4. Time is always given in generations. Since we look backward, time increases towards past.

43 5. Negative times are allowed (e.g. the example given in section 3), but not recommended.

44 6. Population numbers must be consecutive natural integers starting at 1. The number of population
 45 can exceed the number of samples and vice versa : in other words, unsampled populations can be
 46 considered in the scenario on one hand, and the same population can be sampled more than once
 47 on the other hand.

48 7. Multi-furcating population trees can be considered, by writing several divergence events occurring
 49 at the same time. However, one has to be careful to the order of the `merge` events. For instance,
 50 the following piece of scenario will fail :

```
100 merge 1 2
```

1 100 merge 2 3

2 This is because, after the first line, population 2, which has merged with population 1, does not
 3 "exist" anymore (the surviving population is population 1). So, it cannot receive lineages of popu-
 4 lation 3 as it should as a result of the second line. The correct ways are either to put line 2 before
 5 line 1, or to change line 2 to :

6 100 merge 1 3.

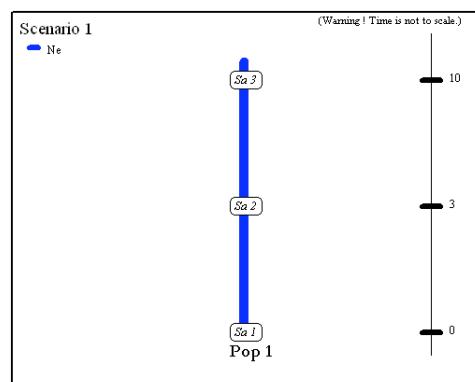
- 7 8. Since times of events can be parameters, the order of events can change according to the values
 8 taken by the time parameters. In any case, before simulating a data set, the program sorts out
 9 events by increasing times ¹. If two or more events occur at the same time, the order is that of the
 10 scenario as it is written by the the user.

- 11 9. Most scenarios begin with sampling events. We then need to know the effective size of the popu-
 12 lations to perform the simulation of coalescences until the next event concerning each population.
 13 One way would have been to provide the population size on the same line of the scenario description.
 14 However, in some scenarios with varying population sizes, it can not be determined what is the
 15 effective size at the sampling time before the set of time parameter values is generated. For that
 16 reason, we decided to provide the effective size and the sampling description on two distinct lines.

17 **Examples** Below are some usual scenarios with increasing complexity. Each scenario is coded on the
 18 left side and a graphic representation given by DIYABC is printed on the right side

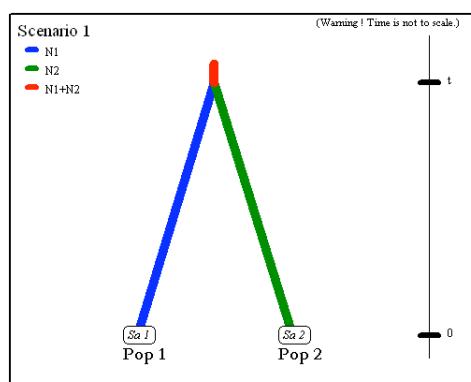
- 19 1. One population from which several samples have been taken at various generations : 0, 3 and 10.
 20 The only unknown parameter of the scenario² is the effective population size.

22 **Ne**
 0 sample 1
 3 sample 1
 10 sample 1



- 23 2. Two populations of size N1 and N2 have diverged t generations in the past from an ancestral pop-
 24 ulation of size N1+N2.

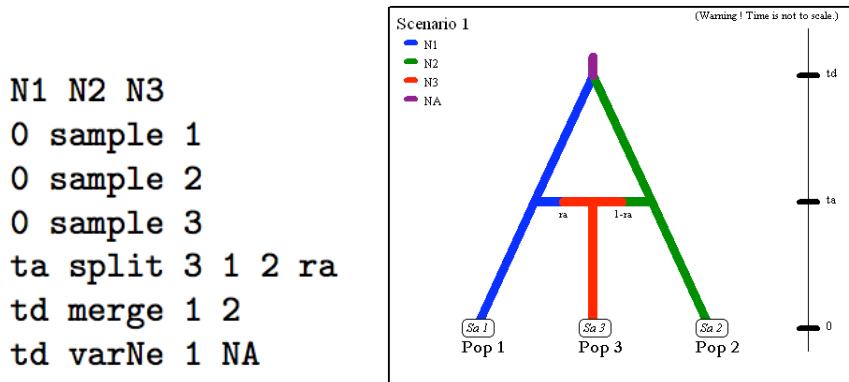
26 **N1 N2**
 0 sample 1
 0 sample 2
 t merge 1 2
 t varNe 1 N1+N2



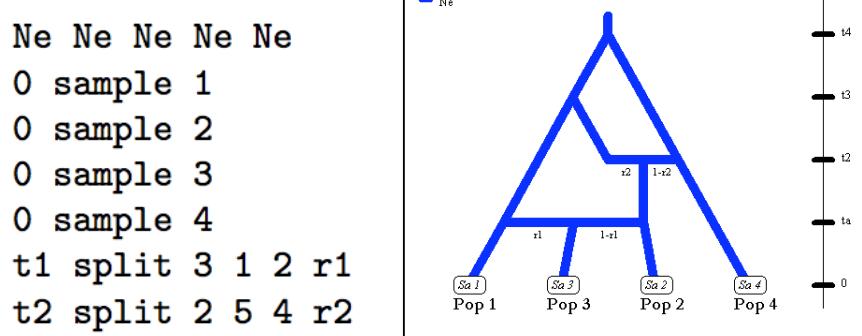
¹Sorting events by increasing times can only be done when all time values are known, i.e. when simulating datasets. When checking scenarios, all time values are not yet defined, so that when visualizing a scenario, events are represented in the same order as they appear in the window used to define the scenario.

²Of course, there are also one or more parameter(s) for the mutation model.

- 1 3. Two parental populations (1 and 2) with constant effective populations sizes N_1 and N_2 have diverged at time t_d from an ancestral population of size N_A . At time t_a , there has been an admixture event between the two populations giving birth to an admixed population (3) with effective size N_3 and with an admixture rate r_a relative to population 1.
- 2
- 3
- 4
- 5



- 6
- 7 4. The next scenario is slightly more complicated. It includes four population samples and two admixture events. For simplicity sake, all populations are assumed to have identical effective sizes (N_e).
- 8
- 9



11 Note that although there are only four samples, the scenario includes a fifth population. This population which diverged from population 1 at time t_3 was a parent in the admixture event occurring at time t_2 . Note also that the first line must include the effective sizes of the five populations.

12

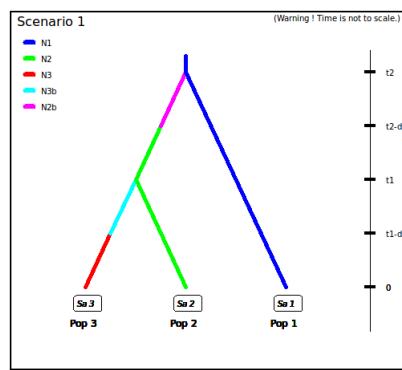
13

14

- 15 5. The following three scenarii correspond to a classic invasion history from an ancestral population (population 1). In scenario 1, population 3 is derived from population 2, itself derived from population 1. In scenario 2, population 2 derived from population 3, itself derived from population 1. In scenario 3, both populations 2 and 3 derived independently from population 1. The same trio of scenarii will be taken later in a fully described example. Note that when a new population is created from its ancestral population, there is an initial size reduction (noted here N_{2b} for population 2 and N_{3b} for population 3) since the invasive population generally starts with a few immigrants.
- 16
- 17
- 18
- 19
- 20
- 21

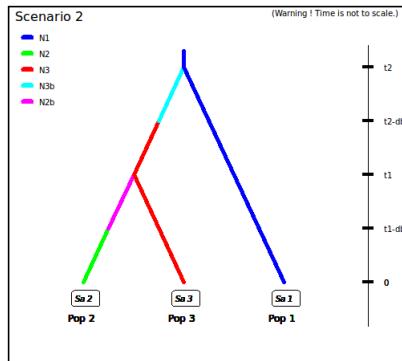
22 Scenario 1

N1 N2 N3
 0 sample 1
 0 sample 2
 0 sample 3
 t1-db VarNe 3 N3b
 t1 merge 2 3
 t2-db VarNe 2 N2b
 t2 merge 1 2



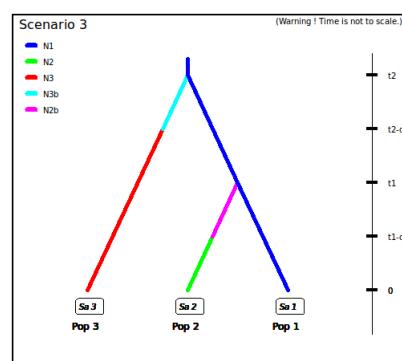
Scenario 2

N1 N2 N3
 0 sample 1
 0 sample 2
 0 sample 3
 t1-db VarNe 2 N2b
 t1 merge 3 2
 t2-db VarNe 3 N3b
 t2 merge 1 3



Scenario 3

N1 N2 N3
 0 sample 1
 0 sample 2
 0 sample 3
 t1-db VarNe 2 N2b
 t1 merge 1 2
 t2-db VarNe 3 N3b
 t2 merge 1 3



1 2.3 Mutation model parameterization (microsatellite and DNA sequence loci)

2 The program can analyse microsatellite data and DNA sequence data altogether as well as separately.
3 In the current version, there are still two restrictions. First, all loci in an analysis must be genetically
4 independent. Second, for DNA sequence loci, intralocus recombination is not considered.

5 Loci are grouped by the user according to its needs (this an improvement of the current version which
6 imposed all loci of a given category to follow the same mutation model). A different mutation model can
7 be defined for each group. For instance, one group can include all microsatellites with motifs that are 2
8 bp long and another group those with a 4 bp long motif. Also, with DNA sequence loci, nuclear loci can
9 be grouped together and a mitochondrial locus form a separate group.
10

11
12 The parameterization of the two categories of markers is now described below.

13 2.3.1 Microsatellite loci

14 Although a variety of mutation models have been proposed for microsatellite loci (Whittaker *et al.*,
15 2003), it is usually sufficient to consider only the simplest models (Cornuet *et al.*, 2006). This has the
16 non-negligible advantage of reducing the number of parameters, which can be a real issue when complex
17 scenarios are considered. This is why we chose the Generalized Stepwise Mutation model (Estoup *et al.*,
18 2002). Under this model, a mutation increases or decreases the length of the microsatellite by a number
19 of repeated motifs following a geometric distribution. This model necessitates only two parameters :
20 the mutation rate (μ) and the parameter of the geometric distribution (P). The same mutation model
21 is imposed to all loci of a given group. However, each locus has its own parameters (μ_i and P_i) and,
22 following a hierarchical scheme, each locus parameter is drawn from a gamma distribution with mean the
23 mean parameter. Note also that :

- 24** 1. individual loci parameters (μ_i and P_i) are considered as nuisance parameters and hence are never
25 recorded. Only mean parameters are recorded.
- 26** 2. The variance or shape parameter of the gamma distributions are set by the user and are NOT
27 considered as parameters.
- 28** 3. The SMM or Stepwise Mutation Model is a special case of the GSM in which the number of repeats
29 involved in a mutation is always one. Such a model can be easily achieved by setting the maximum
30 value of mean P (\bar{P}) to 0. In this case, all loci have their P_i set equal to 0 whatever the shape of
31 the gamma distribution.
- 32** 4. All loci can be given the same value of a parameter by setting the shape of the corresponding
33 gamma distribution to 0 (this is NOT a limiting case of the gamma, but only a way of telling the
34 program).

35 Eventually, to give more flexibility to the mutation model, the program offers the possibility to consider
36 mutations that insert or delete a single nucleotide to the microsatellite sequence. In the previous version,
37 this option was considered as marginal, and was not treated in the same way as the motif size stepwise
38 mutational process, i.e. there was no associated parameter that could be adjusted to the data. This has
39 been changed in this version : it is now possible to use a mean parameter (named $\mu_{(SNI)}$) with a prior
40 to be defined and individual loci having either values identical to the mean parameter or drawn from a
41 Gamma distribution.

42 2.3.2 DNA sequence loci

43 Note first that this version of the program does not consider insertion-deletion mutations, mainly because
44 there does not seem to be much consensus on this topic. Concerning substitutions, only the simplest
45 models are considered. We chose the Jukes-Cantor (1969) one parameter model, the Kimura (1980) two
46 parameter model, the Hasegawa-Kishino-Yano (1985) and the Tamura-Nei (1993) models. The last two
47 models include the ratios of each nucleotide as parameters. However, in order to reduce the number
48 of parameters, these ratios have been fixed to the values calculated from the observed data set for each
49 DNA sequence locus. Consequently, this leaves two and three parameters for the Hasegawa-Kishino-Yano
50 (HKY) and Tamura-Nei (TN), respectively. Also, two adjustments are possible : one can fix the fraction
51 of constant sites (those that cannot mutate) on the one hand and the shape of the Gamma distribution

- ¹ of mutations among sites on the other hand.
² As for microsatellites, all sequence loci of the same group are given the same mutation model with
³ mean parameter(s) drawn from priors and each locus has its own parameter(s) drawn from a Gamma
⁴ distribution (same hierarchical scheme). Notes 1, 2 and 4 of previous subsection (2.3.1) apply also for
⁵ sequence loci.

⁶ 2.4 SNPs do not require mutation model parameterization

- ⁷ SNPs have two characteristics that allow to get rid of mutation models : they are polymorphic and they
⁸ present only two allelic (ancestral and derived) states. In order to be sure that all analyzed SNP loci
⁹ have the two characteristics, non polymorphic loci are disregarded right from the beginning of analyses.
¹⁰ Consequently, no matter *how* it occurred, we can assume that there occurred one and only one mutation
¹¹ in the coalescence tree of sampled genes. We will see below that this largely simplifies (and speeds
¹² up) SNP data simulation. Also, this advantageously reduces the dimension of the parameter space (no
¹³ mutation parameters which are often considered as nuisance parameters). There is however a potential
¹⁴ drawback which is the absence of any calibration generally brought by priors on mutation parameters :
¹⁵ only (time/effective size) ratios will be informative.

¹⁶ 2.5 Prior distributions

- ¹⁷ The Bayesian aspect of the ABC approach implies that parameter estimations are based on prior knowl-
¹⁸ edge about these parameters. This translates into prior distributions of parameters. The program offers
¹⁹ a choice among usual probability distributions, i.e. Uniform, Log-Uniform, Normal or Log-Normal for
²⁰ historical parameters and Uniform, Log-Uniform or Gamma for mutation parameters. Extremum values
²¹ and other parameters (e. g. mean and standard deviation) must be filled in by the user.
²² In addition, one can impose some simple conditions on historical parameters. For instance, there can
²³ be two times parameters with overlapping prior distributions. However, we want that the first one, say
²⁴ t_1 , to always be larger than the second one , say t_2 . For that, we just need to set $t_1 > t_2$ in the
²⁵ corresponding edit-windows. Such a condition needs to be between two parameters (not a parameter and
²⁶ a number, though this can be set up by giving a minimum and a maximum to the prior distribution) and
²⁷ more precisely between two parameters of the same category (i.e. two effective sizes, two times or two
²⁸ admixture rates). The limit to the number of conditions is imposed by the logics, not by the program.
²⁹ The only binary relationships accepted here are $>$, $<$, \geq and \leq .

³⁰ 2.6 Algorithms for data simulation : main features

- ³¹ Data simulation is based on the Wright-Fisher model. It consists in generating the genealogy of all
³² sampled genes until their most recent common ancestor using coalescence theory.
³³ This begins by randomly drawing a complete set of parameters from their own prior distributions and
³⁴ that satisfy all imposed conditions. Then, once events have been ordered by increasing times, a sequence
³⁵ of *actions* is constructed. If there are more than one locus, the same sequence of actions is used for all
³⁶ successive loci. Possible *actions* fall into four categories :

³⁷ adding a sample to a population :

³⁸ Add as many gene lineages to the population as there are genes in the sample.

³⁹ merge two populations :

⁴⁰ Move the lineages of the second population into the first population.

⁴¹ split between two populations :

⁴² Distribute the lineages of the admixed populations between the two parental population according
⁴³ to the admixture rate.

⁴⁴ coalesce and mutate lineages within a population :

⁴⁵ There are two possibilities here, depending on whether the population is *terminal* or not. We call
⁴⁶ *terminal* the population including the most recent common ancestor of the whole genealogy. In
⁴⁷ a terminal population, coalescences and mutations stop when the MRCA is reached whereas in a
⁴⁸ non terminal population, coalescence and mutations stop when the upper (most ancient) limit is
⁴⁹ reached. In the latter case, coalescences can stop before the upper limit is reached because there
⁵⁰ remains a single lineage, but this single remaining lineage can still mutate.

Two different algorithms are implemented : a generation by generation simulation or a continuous time simulation. The choice, automatically performed by the program, is based on an empirical criterion which ensures that the (approximate³) continuous time algorithm is chosen whenever it is faster than the (exact³) generation by generation while keeping the relative error on the coalescence rate below 5% (see Cornuet *et al.* (2008) for a description of this criterion).

In any case, a coalescent tree is generated over all sampled genes.

Then the simulation process diverges depending on the type of markers : for microsatellite or DNA sequence loci, mutations are distributed over the branches according to a Poisson process whereas for SNP loci, one mutation is applied to a single branch of the coalescent tree, this branch being drawn at random with probability proportional to its length.

Eventually, starting from an ancestral allelic state (established as explained below), all allelic states of the genealogy are deduced forward in time according to the mutation process. For microsatellite loci, the ancestral allelic state is taken at random in the stationary distribution of the mutation model (not considering potential single nucleotide indel mutations). For DNA sequence loci, the procedure is slightly more complicated. First, the total number of mutations over the entire tree is evaluated. Then according to the proportion of constant sites and the gamma distribution of individual site mutation rates, the number and position of mutated sites are generated. Finally, these mutated sites are given 'A', 'T', 'G' or 'C' states according to the selected mutation model. For SNP loci, the ancestral allelic state is arbitrarily set to 0 and it becomes equal to 1 after the mutation.

Each category of loci has its own coalescence rate deduced from male and female effective population sizes . In order to combine different categories (e.g. autosomal and mitochondrial), we have to take into account the relationships among the corresponding effective population sizes. This can be achieved by linking the different effective population sizes to the effective number of males (N_M) and females (N_F) through the sum $N_T = N_F + N_M$ and the ratio $r = N_M/(N_F + N_M)$. We use the following formulae for the probability of coalescence of two lineages within this population :

$$\text{autosomal diploid loci} : p = \frac{1}{8r(1-r)N_T}$$

$$\text{autosomal haploid loci} : p = \frac{1}{4r(1-r)N_T}$$

$$\text{X-linked loci / haplo-diploid loci} : p = \frac{1+r}{9r(1-r)N_T}$$

$$\text{Y-linked loci} : p = \frac{1}{rN_T}$$

$$\text{Mitochondrial loci} : p = \frac{1}{(1-r)N_T}$$

Users have to provide a (total) effective size N_T (on which inferences will be made) and a sex-ratio r . If no sex ratio is provided, the default value of r is taken as 0.5.

2.7 Summary statistics

For each category (microsatellite, DNA sequences or SNP) of loci, the program proposes a series of summary statistics among those used by population geneticists. These summary statistics are mean values or variances over loci of the same group and characterize a single, a pair or a trio of population samples. These are :

2.7.1 for microsatellite loci

Single sample statistics :

1. mean number of alleles across loci
2. mean gene diversity across loci (Nei, 1987)
3. mean allele size variance across loci
4. mean M index across loci (Garza and Williamson, 2001; Excoffier *et al.*, 2005)

Two sample statistics :

1. F_{ST} between two samples (Weir and Cockerham, 1984)

³The terms *approximate* and *exact* are relative to the basic assumptions of the Wright-Fisher model, not to the biological reality of the process.

2. mean index of classification (two samples) (Rannala and Mountain, 1997; Pascual *et al.*, 2007)
3. $(\delta\mu)^2$ distance between two samples (Golstein *et al.*, 1995)
4. mean number of alleles across loci (two samples)
5. mean gene diversity across loci (two samples)
6. mean allele size variance across loci (two samples)
7. shared allele distance between two samples (Chakraborty and Jin, 1993)

7 **Three sample statistics :**

- 8 1. Maximum likelihood coefficient of admixture (Choisy *et al.*, 2004)

9 **2.7.2 for DNA sequence loci**

10 **Single sample statistics :**

- 11 1. number of distinct haplotypes
- 12 2. number of segregating sites
- 13 3. mean pairwise difference
- 14 4. variance of the number of pairwise differences
- 15 5. Tajima's D statistics (Tajima, 1989)
- 16 6. Number of private segregating sites (=number of segregating sites if there is only one sample)
- 17 7. Mean of the numbers of the rarest nucleotide at segregating sites⁴
- 18 8. Variance of the numbers of the rarest nucleotide at segregating sites

19 **Two sample statistics :**

- 20 1. number of distinct haplotypes in the pooled sample
- 21 2. number of segregating sites in the pooled sample
- 22 3. mean of within sample pairwise differences
- 23 4. mean of between sample pairwise differences
- 24 5. F_{ST} between two samples (Hudson *et al.*, 1992)

25 **Three sample statistics :**

- 26 1. Maximum likelihood coefficient of admixture (adapted from Choisy *et al.*, 2004)

27 **2.7.3 for SNP loci**

28 **Single sample statistics :**

- 29 1. proportion of loci with null gene diversity (= proportion of monomorphic loci)
- 30 2. mean gene diversity across polymorphic loci (Nei, 1987)
- 31 3. variance of gene diversity across polymorphic loci
- 32 4. mean gene diversity across all loci (Garza and Williamson, 2001; Excoffier *et al.*, 2005)

33 **Two sample statistics :**

- 34 1. proportion of loci with null Nei's distance between the two samples (Nei, 1972)
- 35 2. mean across loci of non null Nei's distances between the two samples
- 36 3. variance across loci of non null Nei's distances between the two samples
- 37 4. mean across loci of Nei's distances between the two samples

⁴This statistics can provide information in case of recent demographic variation : a recent expansion increases the number of singletons (nucleotides occurring just once at a segregating site) resulting in a low value of this statistics, whereas a recent decline will produce an opposite result.

5. proportion of loci with null F_{ST} distance between the two samples (Weir and Cockerham,
1984)
6. mean across loci of non null F_{ST} distances between the two samples
7. variance across loci of non null F_{ST} distances between the two samples
8. mean across loci of F_{ST} distances between the two samples

6 **Three sample statistics :**

- 7 1. Maximum likelihood coefficient of admixture (Choisy *et al.*, 2004)

8 **2.8 Pre-evaluation of scenarios and prior distributions**

9 This option is proposed to users since version 1.0. The purpose is to check that at least one combination
10 of scenarios and priors can produce simulated data sets that are close enough to the observed data set.
11 This is performed through two kinds of analyses. In the first one, a principal component analysis is
12 performed in the space of summary statistics on at most 100,000 simulated data set and the observed
13 data is added on each plane of the analysis in order to evaluate how the latter is surrounded by simulated
14 data sets. In addition to this global approach, there is a second one in which each summary statistic of
15 the observed data set is ranked against those of the simulated data set. This second analysis helps finding
16 which aspects of the model (including prior) have been mistated. For instance, a grossly overestimated
17 genetic distance (in simulated data sets compared to the observed one) may suggest a misspecification of
18 the prior distribution of the time of divergence of the two involved populations or of the mean mutation
19 rate of the markers. Using this new option before running a full ABC treatment is a convenient way to
20 reveal misspecification of models (scenarios) and/or prior distributions of parameters (see Cornuet *et al.*,
21 2010, for an illustration)

22 **2.9 Estimation of posterior distributions of parameters**

23 Several steps are necessary to get posterior distributions of parameters. First, the normalized Euclidian
24 distance between the observed data set and each simulated data set is computed as the sum of squared
25 differences of summary statistics weighted by the inverse of their variance in the entire set of simulated
26 data. For the i -th data set, the distance is :

$$d_i = \sqrt{\sum_{j=1}^{nstat} \frac{(s_{ij} - s_j^{obs})^2}{V_j}} \quad (1)$$

27 in which s_{ij} is the j -th summary statistics from the i -th data set, s_j^{obs} is the j -th summary statistics
28 from the observed data set and V_j is the variance of the the j -th summary statistics across all simulated
29 data sets. Only the closest data sets are selected for further treatments. The latter includes a local linear
30 regression step aimed at improving the posterior distributions of the parameters (Beaumont *et al.*, 2002).
31 Basically, a multiple linear regression is performed in which summary statistics are the independent
32 variables and parameters the dependent variables. But this regression is also *local* in the sense that more
33 weight in the regression is given to data sets that are closest to the observed data set. This is performed
34 by using a kernel function (the Epanechnikov kernel following Beaumont *et al.* (2002) :

$$K_\delta(d) = \begin{cases} (1.5/\delta)(1 - (d/\delta)^2), & t \leq \delta \\ 0, & t > \delta \end{cases} \quad (2)$$

35 Eventually, parameters are adjusted through this process as :

$$\phi_{ik}^* = \phi_{ik} - (\mathbf{s}_i - \mathbf{s}^{obs})\boldsymbol{\beta}_k \quad (3)$$

36 in which ϕ_{ik} is the k -th parameter of the i -th selected data set, ϕ_{ik}^* is the adjusted corresponding pa-
37 rameter, \mathbf{s}_i is the row vector of summary statistics of the i -th selected data set, \mathbf{s}^{obs} is the row vector of
38 summary statistics of the observed data set and $\boldsymbol{\beta}_k$ is the transposed k -th row vector of the regression
39 coefficient matrix.

40 The adjusted ϕ_{ik}^* of the selected data sets are an approximate sample of the posterior distribution of
41 parameters (Beaumont *et al.*, 2002).

2.10 Model checking

Checking the model is crucial to statistical analysis (p161 in Gelman *et al.*, 1995). Model checking (i.e. the assessment of the goodness-of-fit of a model parameter posterior combination) is a facet of ABC analysis that has been so far neglected (but see Ingvarsson, 2008). Following Gelman *et al.* (1995; pp 159-163), we already implemented this option in DIYABCv1.0, to measure the discrepancy between a model parameter posterior combination and a real data set by considering various sets of test quantities. These test quantities can be chosen among the large set of ABC summary statistics proposed in the program. This option is based on the same kinds of analysis as section 2.7. The main difference is the set of simulated data. Whereas in section 2.7, prior distributions of parameters have been used to simulate data sets, here we use posterior distributions of the same parameters, hence simulating data from the *posterior predictive distribution*.

The first analysis is a principal component analysis in the space of summary statistics using data sets simulated with the **prior** distributions of parameters (exactly as in section 2.7) and the observed data as well as **data sets from the posterior predictive distribution** are represented on each plane of the PCA. If the model fits well the data, one should see on each PCA plane a wide cloud of data sets simulated from the prior, with the observed data set in the middle of a small cluster of datasets from the posterior predictive distribution.

In the second analysis, each summary statistics of the observed data set is ranked against the distribution of the corresponding summary statistics from the posterior predictive distribution. Summary statistics play here the role of *test statistics* (p169 in Gelman *et al.*, 1995).

Since summary statistics are generally not sufficient, it is advised to use different sets of summary statistics to compute the posterior distribution of parameters on one hand and to check the model on the other hand (see Cornuet *et al.*, 2010). This has been implemented in DIYABC.

2.11 Measures of performances

As stressed in previous studies (e.g. Excoffier *et al.*, 2005), the ABC approach provides an efficient way of assessing its own performances for estimating posterior distributions of parameters. The reference table, the building of which represents generally 95 to 99% of the computing time, can be reused with pseudo-observed (test) data sets which in fact have been obtained through simulation with known values of parameters. It is then rather quick and easy to evaluate the performance of the method for parameter estimation by computing statistics such as estimation biases or mean square errors.

These measures of performance have been fully integrated into DIYABC. The performance measures computed by DIYABC are :

the average relative bias : the difference between the point estimate (e) and the true value (v) divided by the true value, $\frac{1}{n} \sum_{i=1}^n \frac{e_i - v_i}{v_i}$, averaged over the n test data sets,

the square Root of the Relative Mean Square Error (RRMSE) : the square root of the average square difference between the point estimate and the true value, divided by the true value, $\sqrt{\frac{1}{n} \sum_{i=1}^n (\frac{e_i - v_i}{v_i})^2}$

the square Root of the Relative Mean Integrated Square Error (RRMISE) : the square root of the average (over test data sets) of the integrated square error (measured on each test data set) divided by the true value, $\sqrt{\frac{1}{n} \sum_{i=1}^n (\frac{\sum_{j=1}^{m_i} (x_{ij} - v_i)^2}{m_i v_i^2})}$, x_{ij} and m_i being the sampled values and the sample size of the posterior distribution in the i -th test data set, respectively.

the Relative Mean Absolute Deviation (RMAD) : the average (over test data sets) of the mean absolute deviation (measured on each data set), divided by the true value, $\frac{1}{n} \sum_{i=1}^n (\frac{\sum_{j=1}^{m_i} |x_{ij} - v_i|}{m_i |v_i|})$

the 50% and 95% coverages : the proportion of test data sets for which the 50% and 95% credibility intervals respectively include the true value.

1 **the factor 2** :the proportion of test data sets for which the point estimate is at least half and at most
 2 twice the true value.

3 **the Relative Median Bias (RMB)** : the 50% quantile of the bias (measured on each data set) divided
 4 by the true value. The bias is computed respectively for each point estimate

5 **the Relative Median Absolute Deviation (RMedAD)** : the 50% quantile (over test data sets) of
 6 the median (over each data set) of the absolute difference between each value of the posterior
 7 distribution sample and the true value divided by the true value.

8 **the Relative Median of the Absolute Error (RMAE)** : the 50% quantile (over test data sets) of
 9 the absolute value of the difference between the point estimate (in each data set) and the true value
 10 divided by the true value.

11 *DIYABC* considers the following three point estimates : mean, median and mode of the ϕ_{ik}^* (sample
 12 of the posterior distribution of each parameter), as defined in subsection 1.7.

13 Concerning the true value (v) appearing in the above formulae, DIYABC offers two possibilities :

14 1. All values v are fixed by the user. If any one of these values is outside the limits given to the
 15 prior for the corresponding parameter, a warning message is issued but the analysis can proceed if
 16 needed.

17 2. All values v are drawn from distributions. These distributions can be different from those of priors.
 18 They may even not be overlapping (no warning message is issued whatever the user's choice).

19 If you want to fix some parameter values and draw the other from distributions, choose the second option
 20 and give the same desired values as minimum and maximum for those fixed parameter values.

21 In order to better assess the information brought by genetic data, *DIYABC* provides a double
 22 estimate of all these bias/precision statistics. As expected, the first one is based on genetic data given
 23 in the data file. The second one is computed as if there was no genetic information, *i.e.* estimates are
 24 based only on parameter priors. Technically, a sample of parameter values is drawn at random from the
 25 reference table. This sample of the same size of the sample of posterior values is used in place of the
 26 latter in all computations.

28 2.12 Comparison of scenarios

29 The ABC approach can also be used to compare possible scenarios for the same data file through the
 30 computation of the posterior probabilities of each scenario and this option is naturally implemented in
 31 DIYABC.

32

33 2.12.1 Reference table

34 First, the reference table can include as many scenarios as desired. By default, the prior probability of
 35 each scenario is uniform, that is each scenario will have approximately the same number of simulated
 36 data sets. But, if for any reason, one wants a different prior probability for each scenario, there is the
 37 possibility to do so.

38

39 Scenarios are drawn according to their own prior probability and then only parameters that are defined
 40 for the drawn scenario are generated from their respective prior distribution. Scenarios may or may not
 41 share parameters.

42 When conditions apply to some parameters (see subsection 2.4), the program provides the possibility of
 43 choosing between two options :

44 1. parameter sets are drawn in their respective prior distributions until all conditions are fulfilled.
 45 2. a single parameter set is drawn and only if all condition are fulfilled, the simulation is performed
 46 and the data set is recorded in the reference table.

47 When there is only one scenario, both options are equivalent, although in the latter option, there might
 48 be less simulated data sets that are recorded than one asked. When there is more than one scenario,
 49 the second option can be viewed as a way to set prior probabilities on scenario that result from imposed
 50 conditions on parameters (see Miller *et al.* (2005) for an example).

¹ **2.12.2 Posterior probability of scenarios**

² The program *DIYABC* provides two estimates of the posterior probability of each scenario :
³ **a direct estimate :** This is simply the number of times that a given scenario is found in the first n_δ
⁴ simulated data sets once the latter, produced under several scenarios, have been sorted by ascending
⁵ distances to the observed data set.

⁶

⁷ **a logistic regression estimate :** Following M.A. Beaumont's suggestion (Fagundes *et al.*, 2007; Beaumont,
⁸ 2008), a polychotomous weighted logistic regression is performed on the first n_δ data sets with
⁹ the proportion of the scenario as the dependent variable and the differences between observed and
¹⁰ simulated data set summary statistics as the independent variables. The intercept of the regression
¹¹ (corresponding to an identity between simulated and observed summary statistics) is taken as the
¹² point estimate. In addition, 95% confidence intervals are computed (Cornuet *et al.*, 2008).

¹³ Since both estimates are dependent upon the chosen threshold (δ), the program provides a range of
¹⁴ 100 estimates for the direct approach (for each one 100-th of n_δ between 0 and n_δ) and up to 10 estimates
¹⁵ for the logistic regression estimates (e.g. one estimate for $kn_\delta/10$ with $k \in [1, 2, \dots, 10]$ when the number
¹⁶ of analyses is set to 10). These estimates are represented in two graphs, one for each kind of estimate.
¹⁷ These two graphs can be printed and/or saved (in *svg*, *jpg*, *png* or *pdf* format). Values can also be output
¹⁸ as a text file.

¹⁹ **2.12.3 Confidence in scenario choice**

²⁰ The program DIYABC offers a last option that allows one to evaluate the confidence in a scenario choice.
²¹ Suppose that we compare 3 scenarios for a given data set and that e.g. scenario 2 had maximum posterior
²² probability. By using this option, we can estimate type I and type II errors when choosing scenario 2 as
²³ the true scenario. To do so, we simulate a given number of data sets according to scenario 1, 2 and 3.
²⁴ Then we count the proportion of times that scenario 2 has not the highest posterior probability among
²⁵ the three competing scenarios when it is the true scenario (type I error, estimated from test data sets
²⁶ simulated under scenario 2) or the proportion of times that scenario 2 has highest posterior probability
²⁷ when it *not* the true scenario (type II error, estimated from test data sets simulated under scenarios 1
²⁸ and 3).

²⁹

³⁰ In *DIYABCv2.0*, a new possibility is offered to the user that may be useful when dealing with many
³¹ summary statistics and many scenarios. In this particular case, the logistic regression has to deal with
³² large matrices and the amount of needed memory on one hand and the computation time on the other
³³ hand can become problematically large. An approximate solution is to replace summary statistics by the
³⁴ components of a factorial discriminant analysis which reduces the number of independent variables to
³⁵ the smallest of number of summary statistics and scenarios. Although the result is only approximate, it
³⁶ can be a useful guide in some specific cases. The gain in time can be large. For instance, the time can
³⁷ be reduced by a 30X factor.

³⁸

³⁹ As for the bias/precision analysis, parameter values can be fixed to given values or drawn from given
⁴⁰ distributions (not necessarily the same as those used as priors for the reference table).

¹ 3. The Graphic User Interface

² When launching the GUI, the home screen appears like this :



⁴

⁵ You can already notice that *DIYABC* works with projects. This notion is new to version 2 of
⁶ *DIYABC*. It is explained in subsection 3.1.

⁷ 3.1 What is a *DIYABC* Project ?

⁸ A *DIYABC* project is a unit of work materialized by a specific and unique directory. A project is defined
⁹ by at least one observed data set and one reference table header file. These files are located in the *Project*
¹⁰ *directory* which name includes an identifier, the date of creation and a number (between 1 and 100).

¹¹

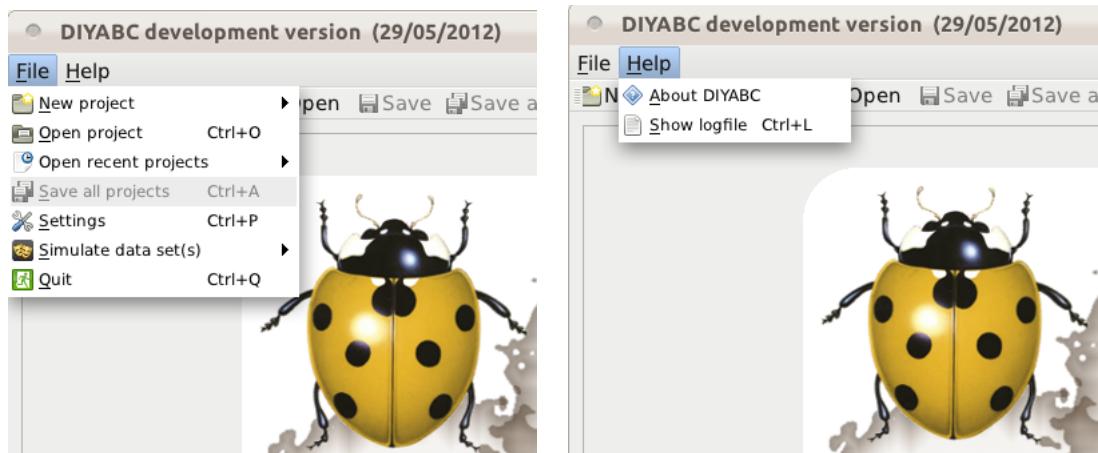
¹² The header file, always named `header.txt`, contains all information necessary to compute a reference
¹³ table associated with the data : i.e. the scenarios, the scenario parameter priors, the characteristics of
¹⁴ loci, the loci parameter priors and the summary statistics to compute. As soon as the first records of
¹⁵ the reference table have been saved in the reference table file, always named `reftable.bin` and also
¹⁶ included in the project directory, the project is "locked". This means that the header file can not be
¹⁷ changed anymore. If one needs to change a scenario or a parameter prior, or a summary statistics, a new
¹⁸ project needs to be defined. This is to guarantee that all subsequent actions performed on the project
¹⁹ are in coherence with the current data and header files. It is of course strongly advised NOT to move
²⁰ files among projects. Incidentally, the `header.txt` file is only built when the project has been saved, the
²¹ information progressively input by the user being saved in a series of temporary files.

²²

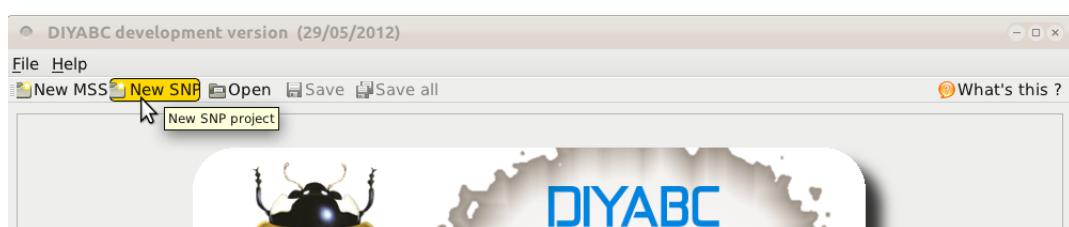
²³ Once a sufficiently large reference table has been simulated, analyses can be performed. Their different
²⁴ output files are copied to the *analysis* directory included in the project directory, and containing as many
²⁵ directories as analyses performed. Hence, it is now much easier to know with certainty the conditions of
²⁶ each analysis.

²⁷ 3.2 Options of the home screen

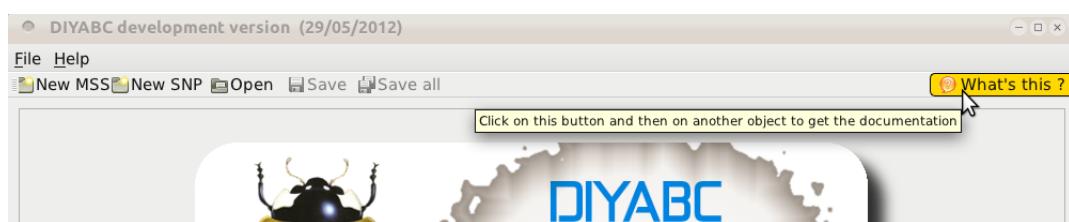
²⁸ The home screen above has two menus and several buttons.
²⁹ Let's start with the menus. Below are shown all submenus :



- 1 The File menu has seven options, namely New project, Open project, Open recent projects, Save all projects, Settings, Simulate data set(s) and Quit. All are self explanatory.
- 2 The Help menu has two options : About DIYABC which opens up a small window providing the names and address of the authors and Show logfile which gives access to a logfile viewer in which are recorded all actions and messages about the execution of the GUI.
- 3 Just below the menu are five shortcuts to main File menu options.



4 On the right, the field What's this ? is an another way to get help on a specific GUI object :



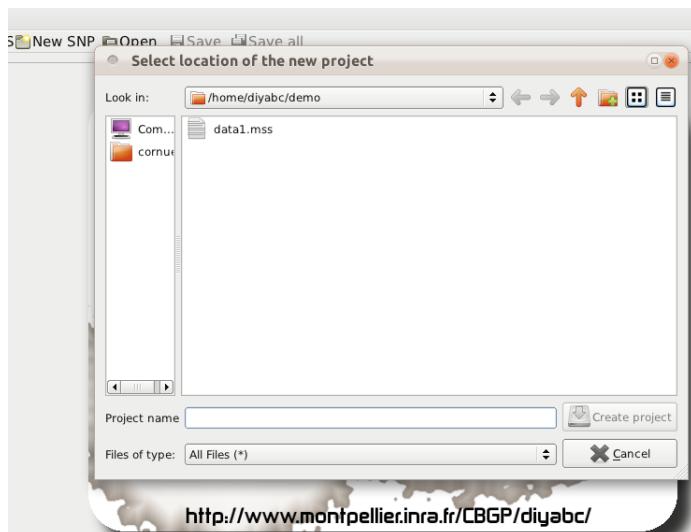
5 Eventually, below the logo, there are three buttons which are duplicate shortcuts :



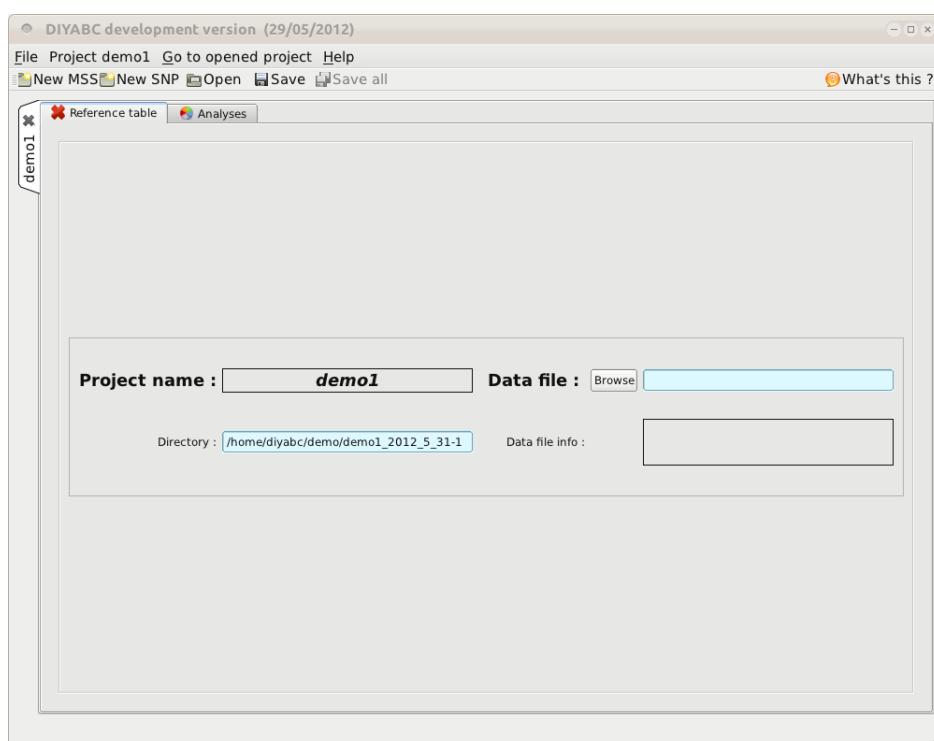
19 3.3 Defining a new project

- 20 Defining a new project requires different steps which are not the same whether the data are SNPs or microsatellites/DNA sequences (MSS). Let start with an MSS project : click on one of the following :
- 21 • File menu > New project > Microsatellites and/or sequences
- 22 • the menu shortcut New MSS

- 1 • the bottom left button **New Microsat/Sequence project**
- 2 or press simultaneously the **Control** and **M** keys.
- 3 A new window appears in which the user can choose a location and a name for the new project as shown below :
- 4
- 5



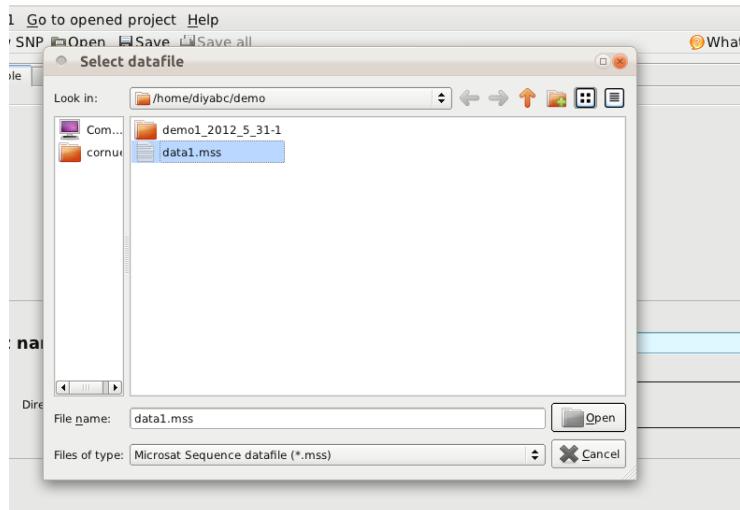
- 6
- 7 Let's enter **demo1** as the project name and click on the **Create project** button.
- 8 The following screen appears :
- 9



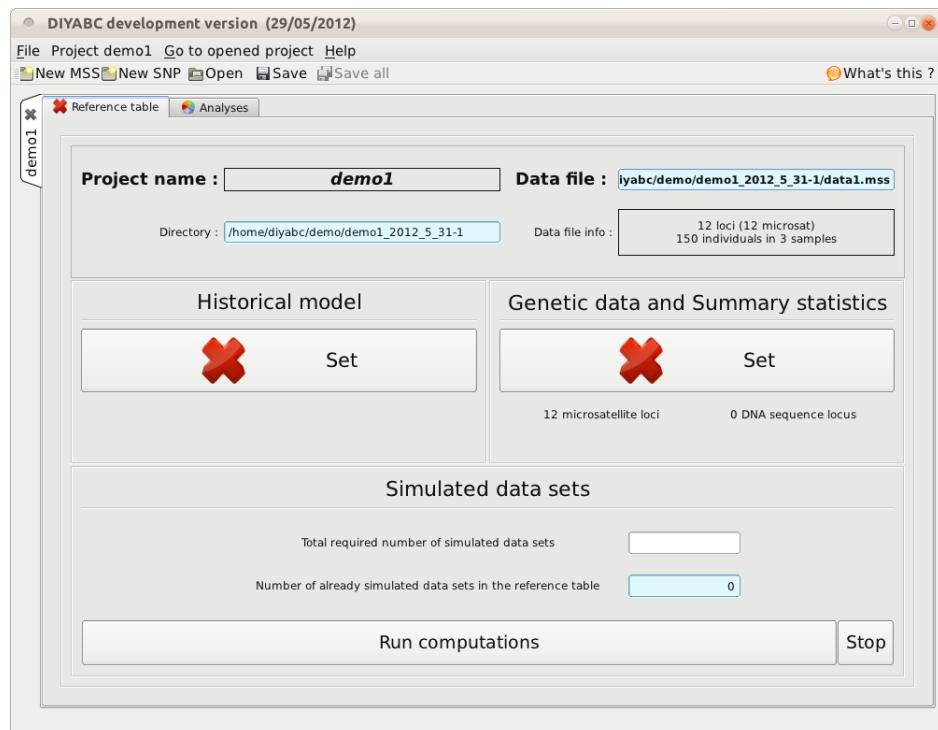
- 10
- 11 The **demo1** project and all its future files will be located in the directory **demo1_2012_5_31-1**.

12 **3.3.1 Step 1 : choosing the data file**

- 13 We next need to choose the data file of the project. This is performed by clicking on the corresponding **Browse** button (previous screen). The usual file browsing screen appears (below) and one has to select a Genepop format data file, here **data1.mss**.
- 14
- 15
- 16



- 1 Go to opened project Help
 2 SNP Open Save Save all
 3 Select datafile
 4 Look in: /home/diyabc/demo
 5 demo1_2012_5_31-1
 6 cornut
 7 data1.mss
 8
 9
- 1 Clicking on the **Open** button leads to the following screen with the edit field filled with the name of the data file and some characteristics of this data file appearing on the screen (number of loci, individuals and samples).
 2 Below these fields are two panels indicating that we need to provide information about the Historical model (left panel) and about the Genetic data and associated Summary statistics (right panel). The red crosses on both panels will change to green checks once the corresponding information will be completed.

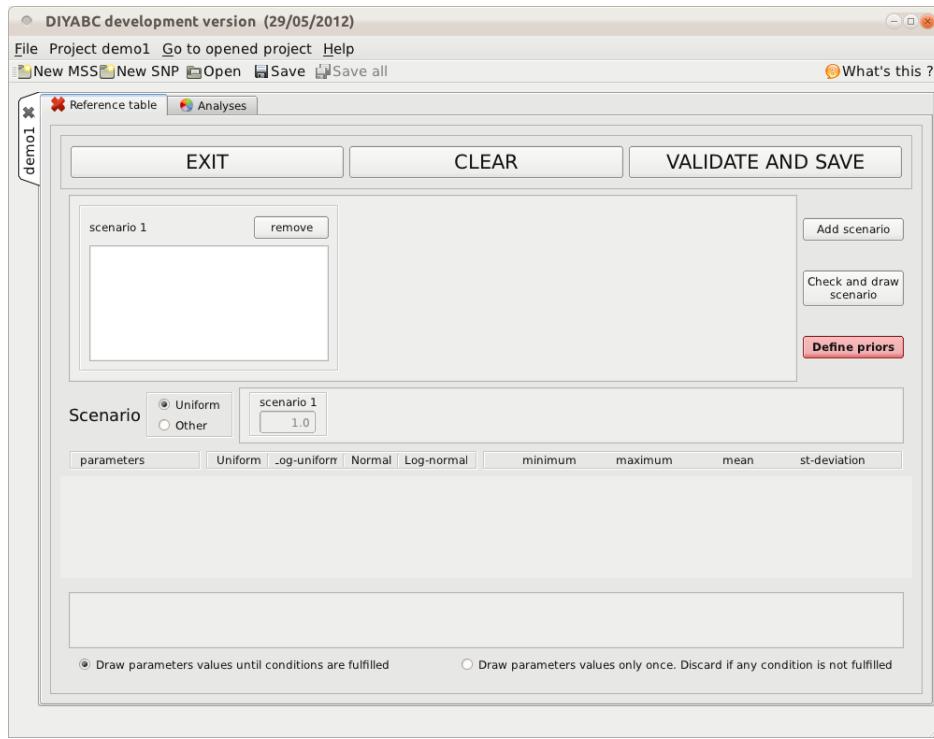


10

1 3.3.2 Inform the Historical model

- 2 Click on the corresponding **Set** button. The following screen, familiar to users of previous versions,
3 appears:

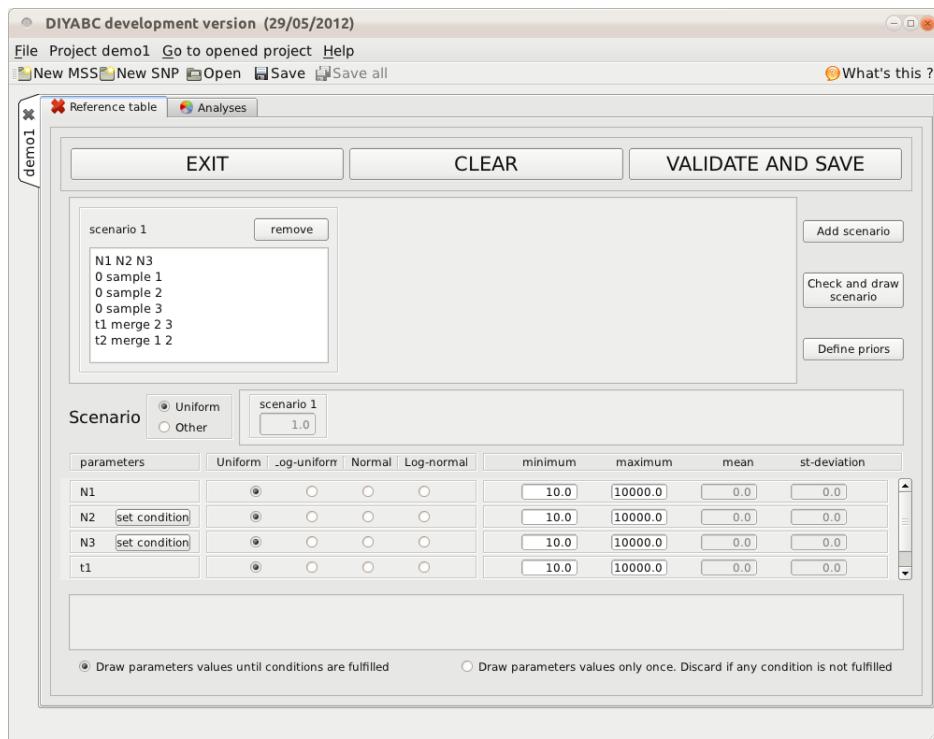
4



5

- 6 Let's enter a simple scenario in scenario 1 edit window and click on the **Define priors** button. We get
7 this :

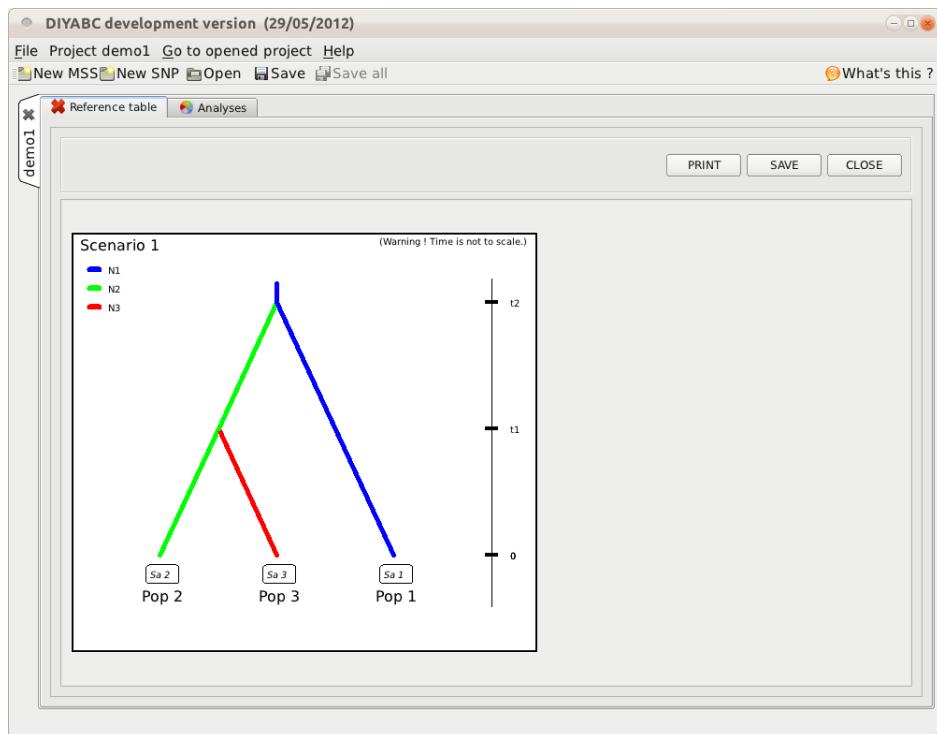
8



9

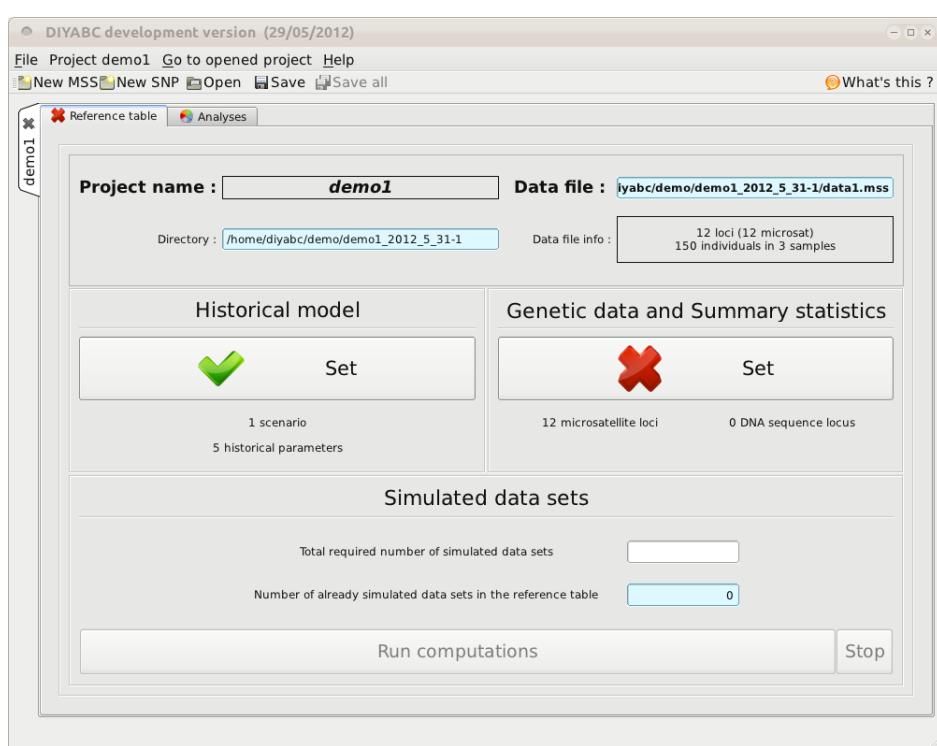
- 10 The parameter prior frame allows to choose the prior density of each parameter. A parameter is
11 anything in the scenario that is not a keyword (here `sample` and `merge`), nor a numeric value. In our
12 little scenario, parameters are hence : N1, N2, N3, t1 and t2. In our example, we need to set the priors
13 on t1 and t2 such that t2>t1. We can do it either by using the `set condition` button or by playing
14 with the minimum and maximum values of the two parameters.

1
2 If we click on the **Check scenario** button, the logic of the scenario is checked and if it is found OK,
3 and if the scenario is drawable, the drawing appears on a new frame :
4



5
6 The scenario can be saved by clicking on the **SAVE** button. The frame can be close by clicking on
7 the **CLOSE** button.

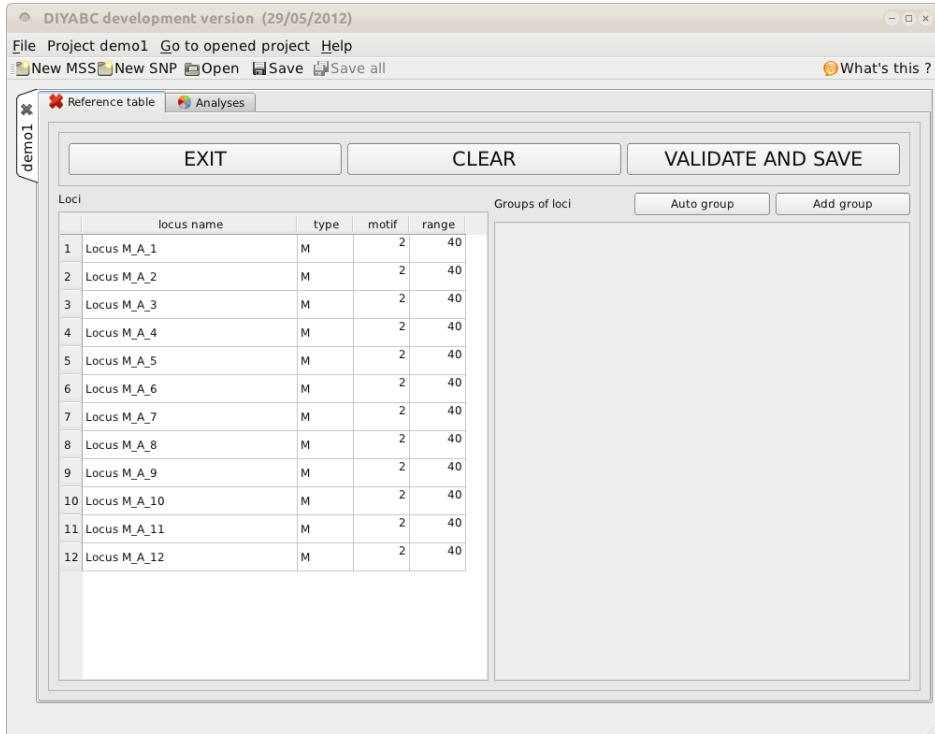
8
9 Since the scenario has been checked, we can validate and save the historical model by clicking on the
10 **VALIDATE AND SAVE** button (bottom screen of p 21). We then go back to the project screen in which
11 the historical model has now received the green check sign.



1 3.3.3 Inform the Genetic model

- 2 Click on the corresponding **Set** button. We get the following screen :

3

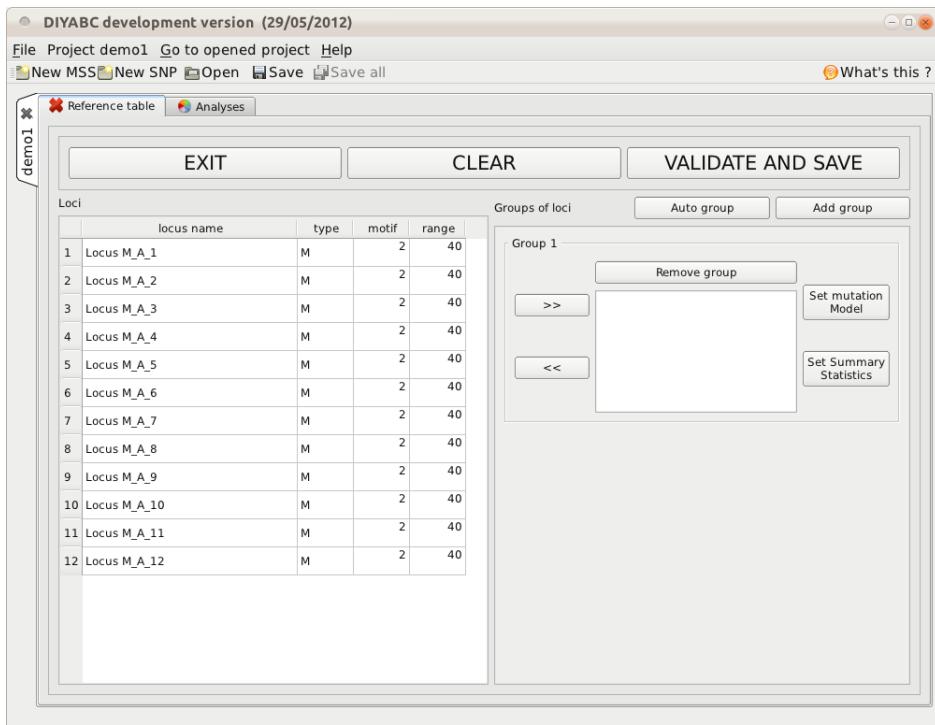


4

5 On the left part of the screen, there is the list of loci, with their type (M for microsatellites or S for DNA sequences) and the motif size and range for microsatellite loci only. Actually, the values for motif size and range are just default values and do not necessarily correspond to the actual data. The user who knows the real values for its data is required to set the correct values at this stage. If the range is too short to include all observed values, a message appears in a box asking to enlarge the corresponding range. Note that the range is measured in number of motifs, so that a range of 40 for a motif length of 2 bp means that the difference between the smallest and the longest alleles should not exceed 80 bp.

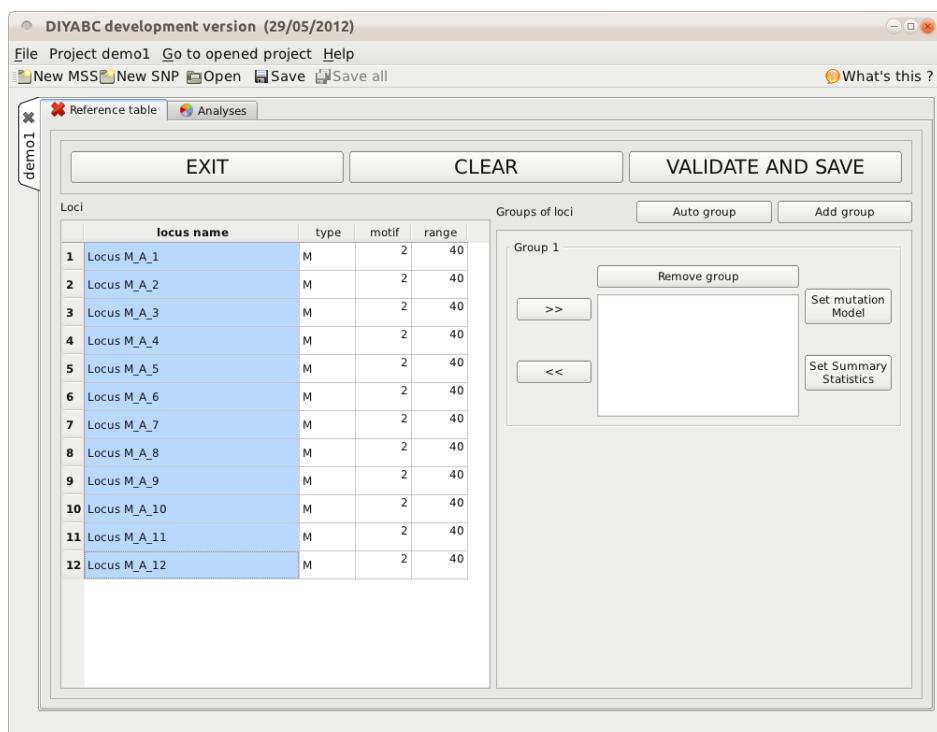
- 12 We then need to define at least one group of loci by clicking on the **Add group** button. We get this :

13

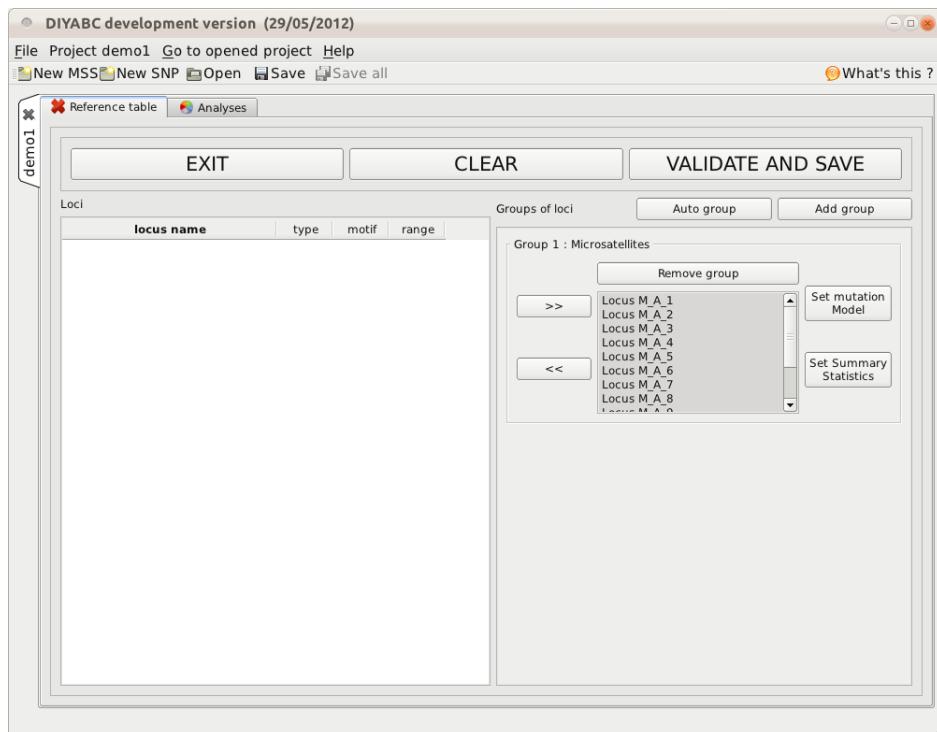


14

1 Suppose we want the three loci in the same group. We select them like in any table, extending the
 2 selection with the Shift and Control keys (see below) :



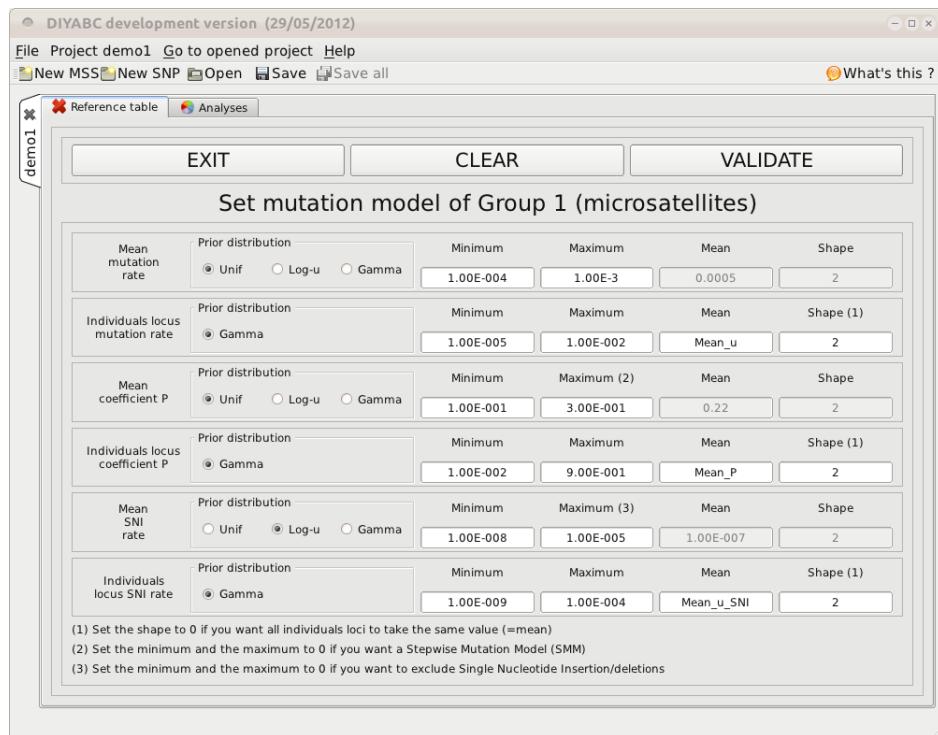
4 and then pressing the **>>** button :



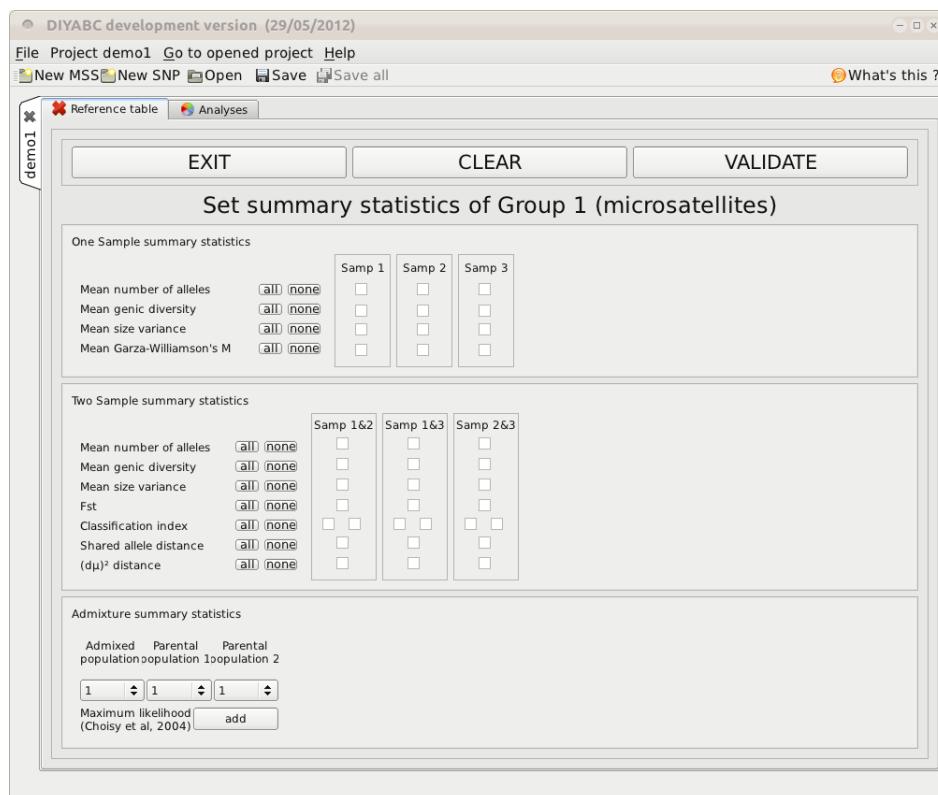
7 Note that the **Auto group** button would have produced the same result of putting all the microsatellite loci in the same group.

10 We then need to define the mutation model and the summary statistics of the locus group. Clicking
 11 on the **Set mutation model** button, the following screen appears :

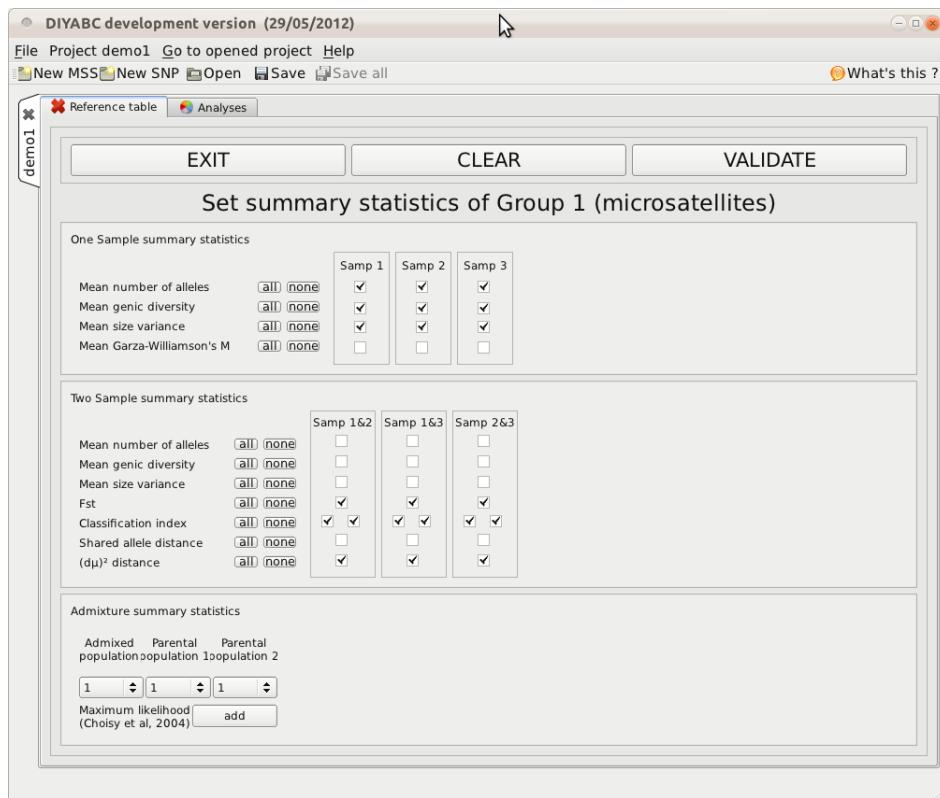
13



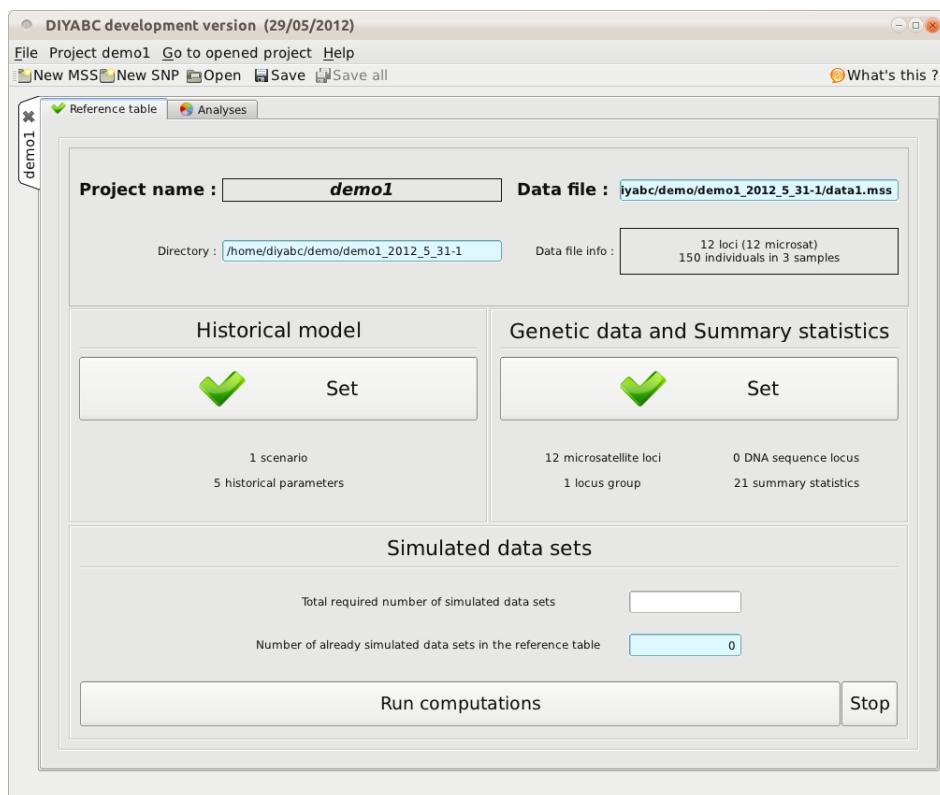
- Once the mutation model of Group 1 is defined, we click on the **VALIDATE** button to go back to the previous screen. Clicking on the **Set Summary statistics** button, we get the following screen :



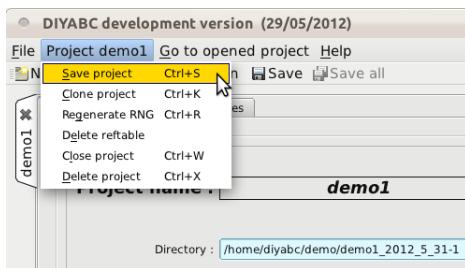
- We define summary statistics by checking the corresponding boxes :



Once finished, we click on the **VALIDATE** button to go back to the screen of p24. Now, we can validate also this screen which brings us back to the screen of p22. The latter looks now like this :

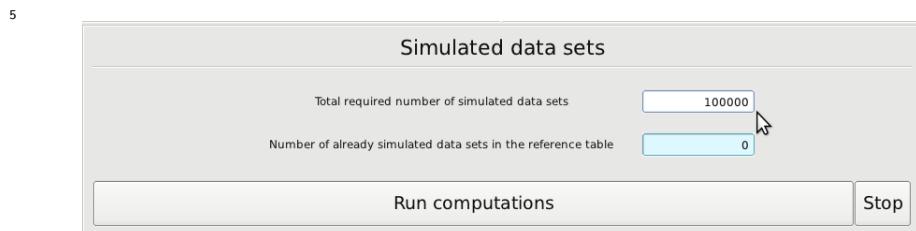


At that moment, the project directory includes the following files : a copy of the data file, and four configuration files : `conf.analysis`, `conf.gen.tmp`, `conf.hist.tmp`, `conf.tmp`. Note that the project is not yet saved. To save the project, we need either to save it explicitly by using the **File** menu (see below) or to start simulating data sets (next section). Saving the project results in saving the `header.txt` file in the project directory.

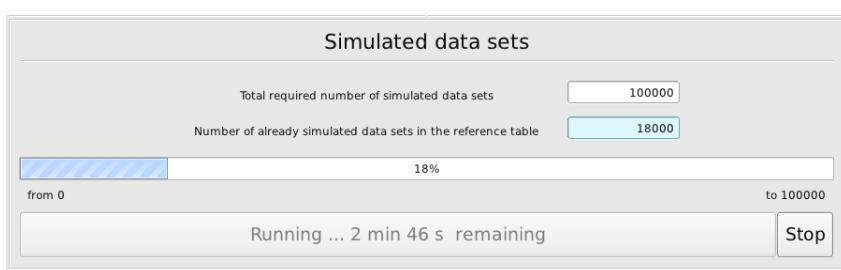


3.4 Building the reference table

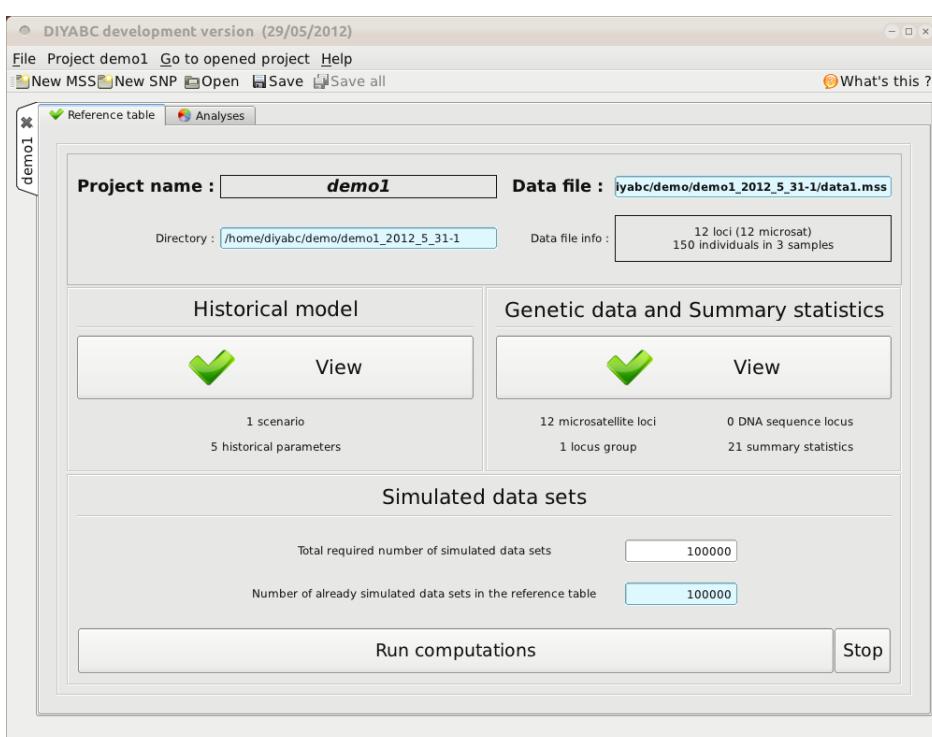
- 3 Keeping on the current screen, indicate the required number of data sets to simulate for the reference table :



Then click on the **Run computations** button. If things go well, you will soon see the progress both into the edit window "Number of simulated data sets in the reference table" and in the progress bar below. Also, you have an estimate of the remaining time (at the left of the **Run computations** button):

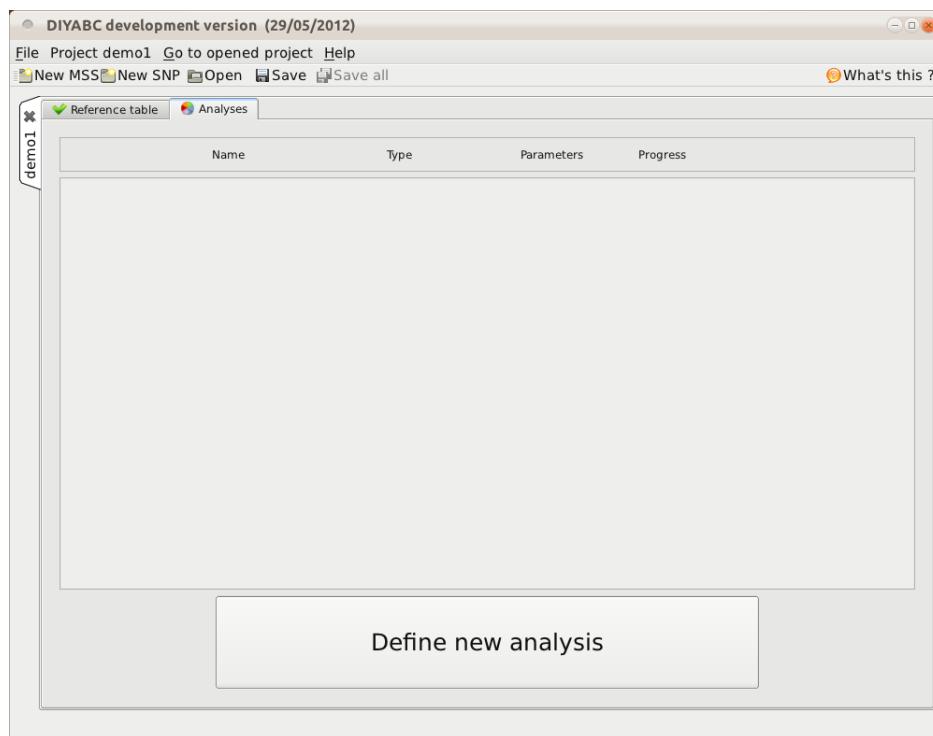


When the computation is finished, the screen looks like this :

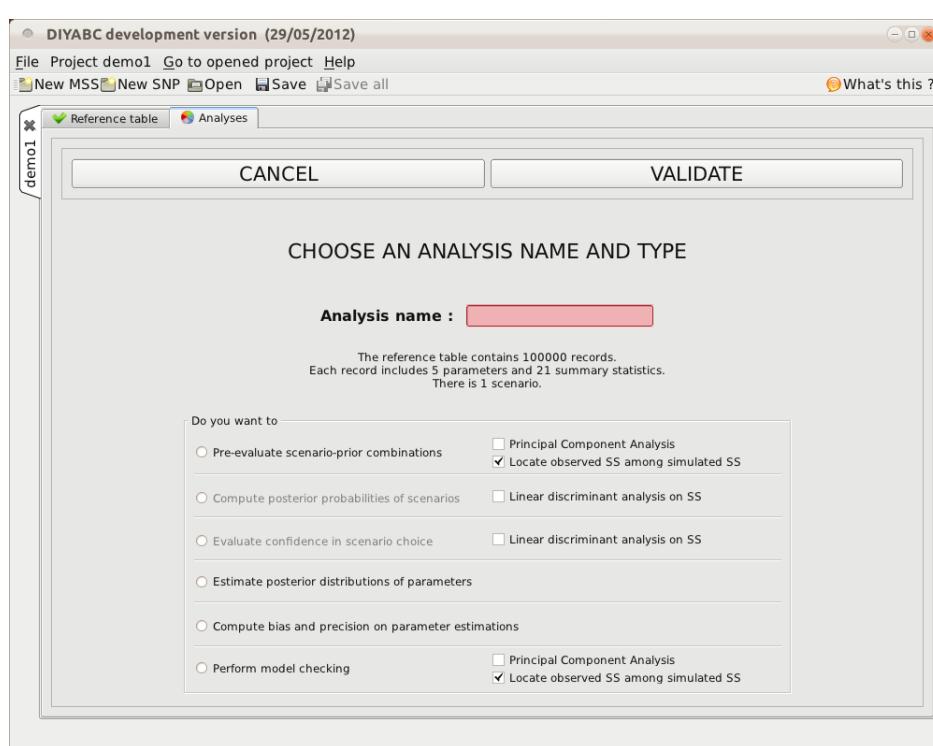


1 3.5 Performing analyses

2 We have now everything necessary to perform analyses. The current screen shows two tabs : Reference
3 table and Analyses. Let's click on the Analyses tab. We get this new screen :



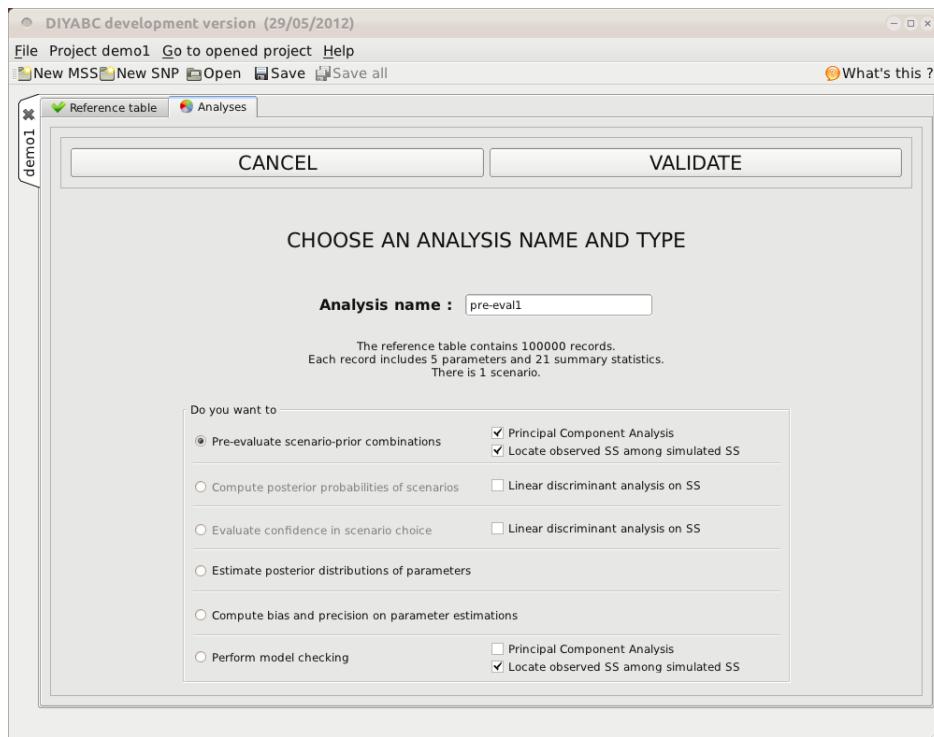
5
6 First, we need to define the analysis we want to perform. So we click on the **Define new analysis**
7 button and get this new screen :



9
10 We need to choose among the six possible types of analyses (actually, only four of them are possible,
11 since the reference table includes a single scenario). We decide to first check whether the model (scenario
12 and parameter prior definition) is off the target or not. This can be appreciated through the analysis
13 denominated **Pre-evaluate scenario prior combination**. To illustrate the result, we also ask for a
14 principal component analysis by checking the corresponding square. Eventually, we give the name of

1 pre-eval1 to this first analysis. The screen now looks like this :

2

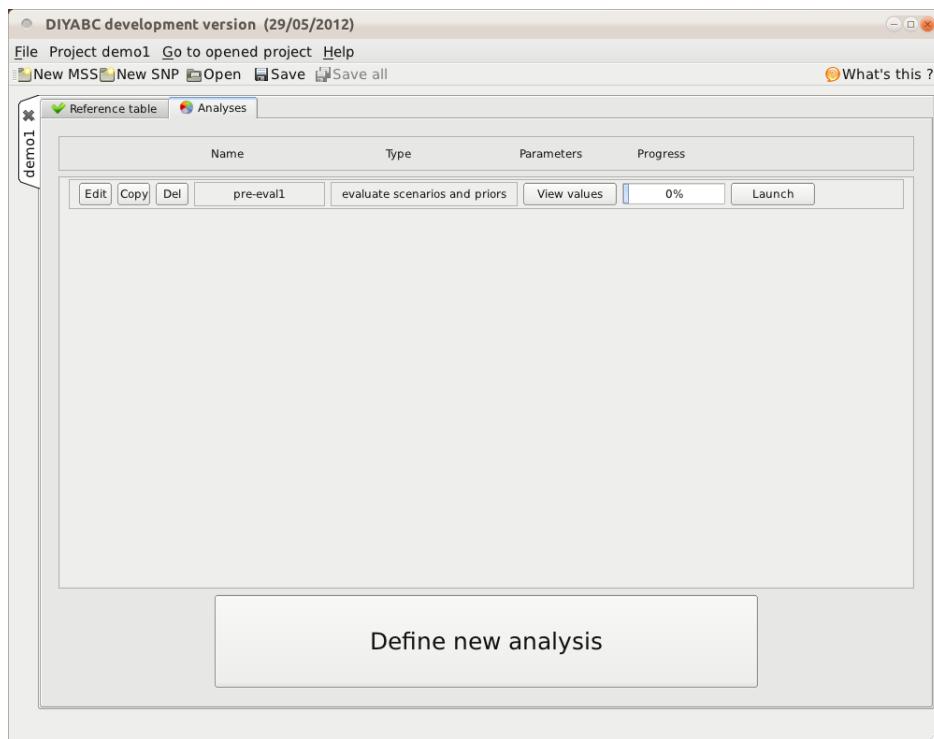


3

4

5 After clicking on the **VALIDATE** button, we go back to the previous screen. However, the new analysis now appears on top of the analysis panel. For each analysis, this panel provides its name and type, the list of parameters that will be transmitted (in a coded way) to the computation program, a progress bar that approximates the progress of the analysis run, and four buttons. The right button has to be clicked to launch the analysis. The three left buttons provide a way to copy an analysis (**Copy** button), to make some modifications (**Edit** button) before launching it or to delete the analysis (**Del** button).

11



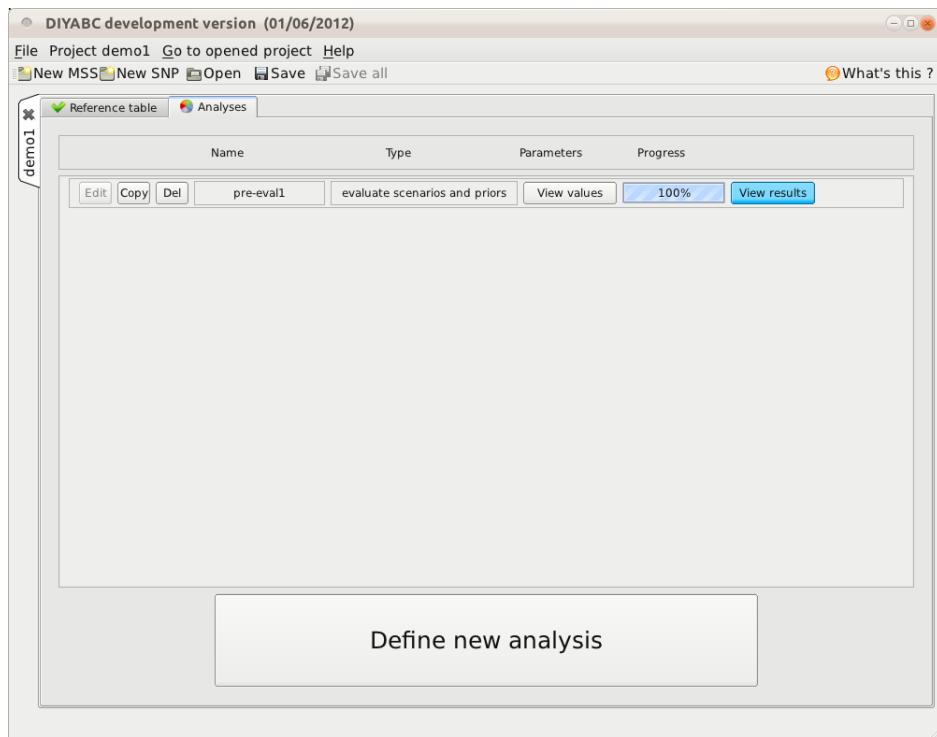
12

13

14 Let's click on the **Launch** button. This analysis is very fast (ca 1 second) so that the progress bar

1 shows almost immediately a 100% value :

2

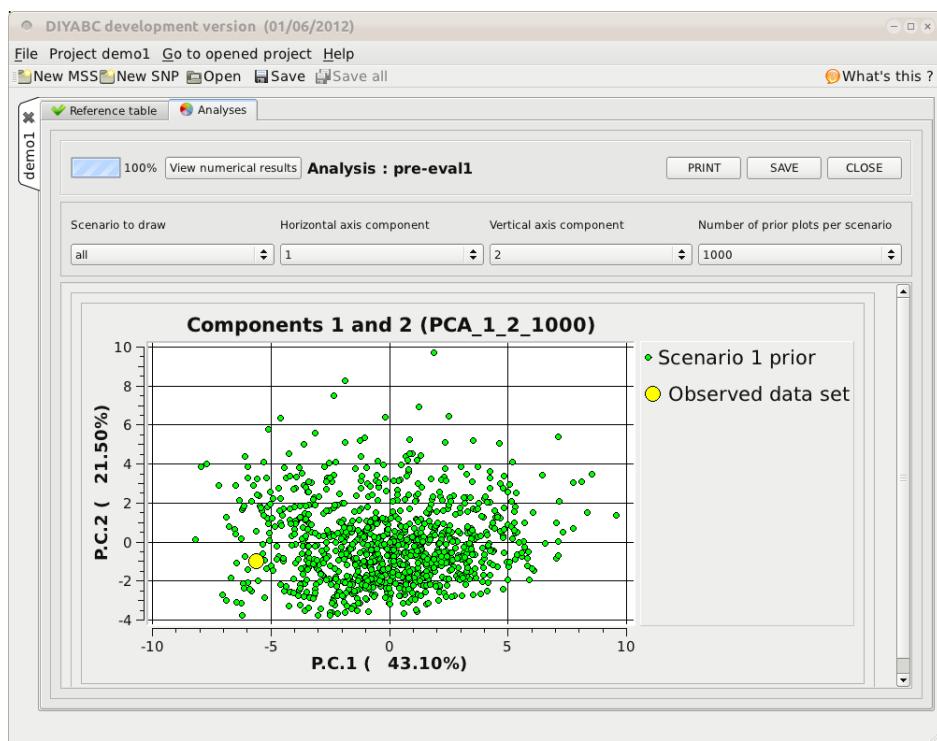


3

4

5 To view results, just click on the **View results** button. After some seconds (while the program reads
6 the PCA result file), we can see this :

7

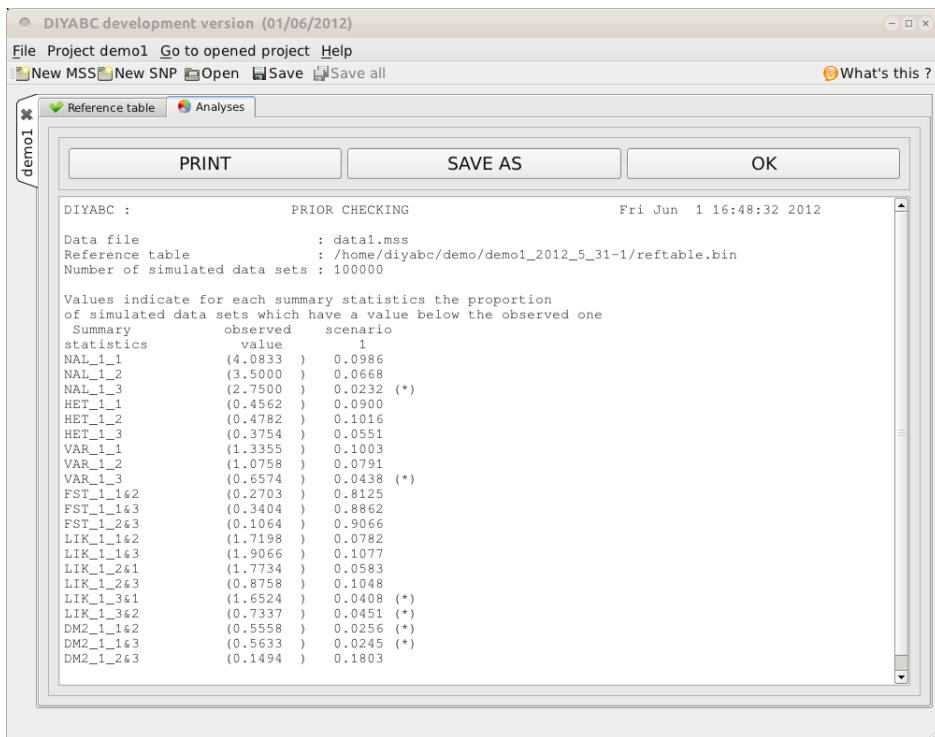


8

9

10 The results are shown PCA plane by PCA plane. Each (small) dot represents a simulated dataset from
11 the reference table and the large yellow dot represents the observed data set. The initial components of
12 datasets are the values of the summary statistics from which are computed the principal components. The
13 four drop-lists (**Scenario to draw**, **Horizontal axis component**, **Vertical axis component**, **Number**
14 **of prior plots per scenario**) can be used to explore further the results of the PCA.

- 1 The graphic can be printed or saved (**PRINT** and **SAVE** buttons, respectively). Clicking on the **CLOSE**
 2 button closes the result window. Eventually, clicking on the **View numerical results** opens up another
 3 screen as shown below :



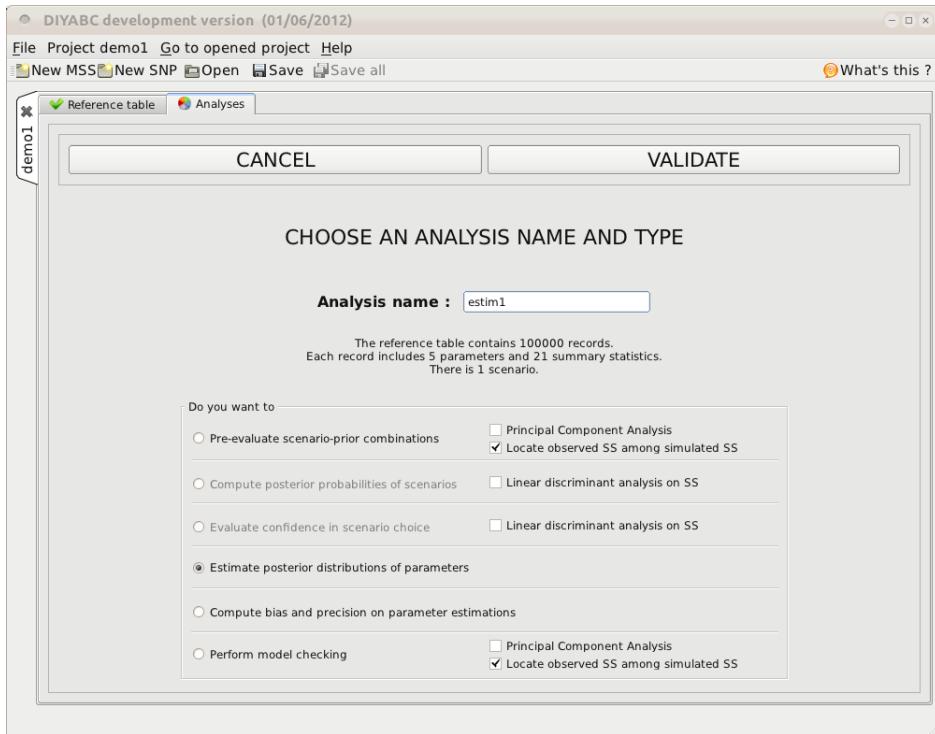
- 5
 6 This screen is obtained by computing for each summary statistics the proportion of simulated data
 7 (considering the total reference table) that have a value below the value of the observed dataset. A star
 8 indicates proportions lower than 5% or greater than 95% (two stars, <1% or >1%; three stars, <0.1% or
 9 >0.1%).
 10

- 11
 12 As usual, results can be printed (**PRINT**) and/or saved (**SAVE**). Click on **OK** to leave this screen.
 13
 14 Although we get one star for a few summary statistics, we conclude that our model is suitable enough
 15 to proceed to other ABC analyses.

16 **3.5.1 ABC parameter estimation**

- 17 Back on the screen of page 30, we click on the **Define new analysis** button. We choose the **Estimate**
 18 **posterior distribution of parameters** option and we call `estim1` this second analysis :

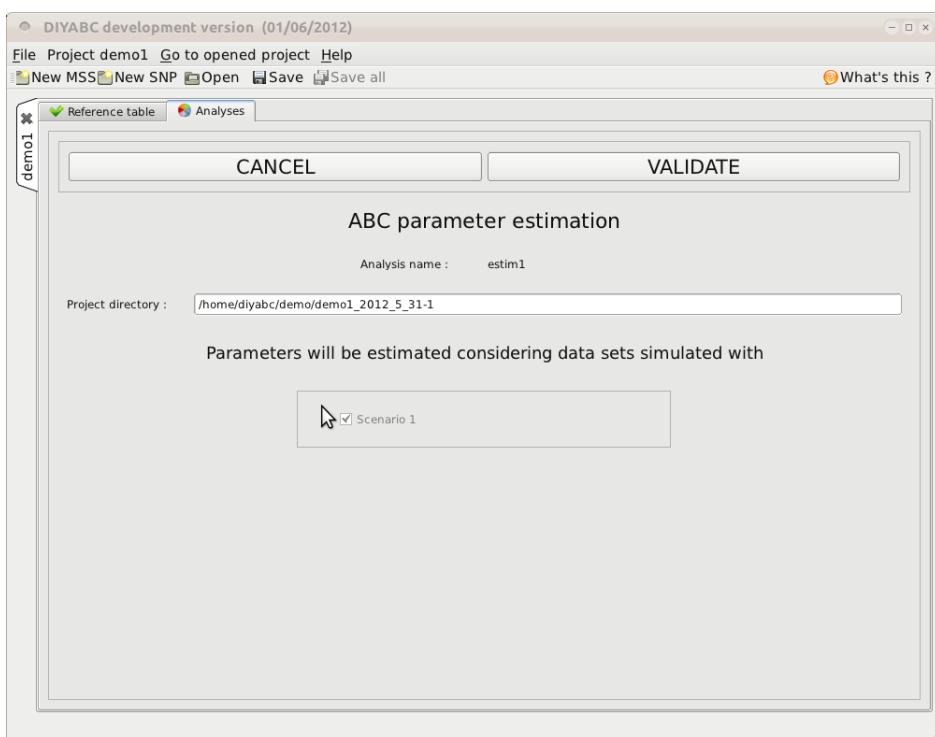
19



1

2

3 We click on the **VALIDATE** button and get the following screen in which we can choose the scenario
4 to use for this estimation. Since a single scenario has been defined, there is nothing else to do than to
5 click on the **VALIDATE** button :



6

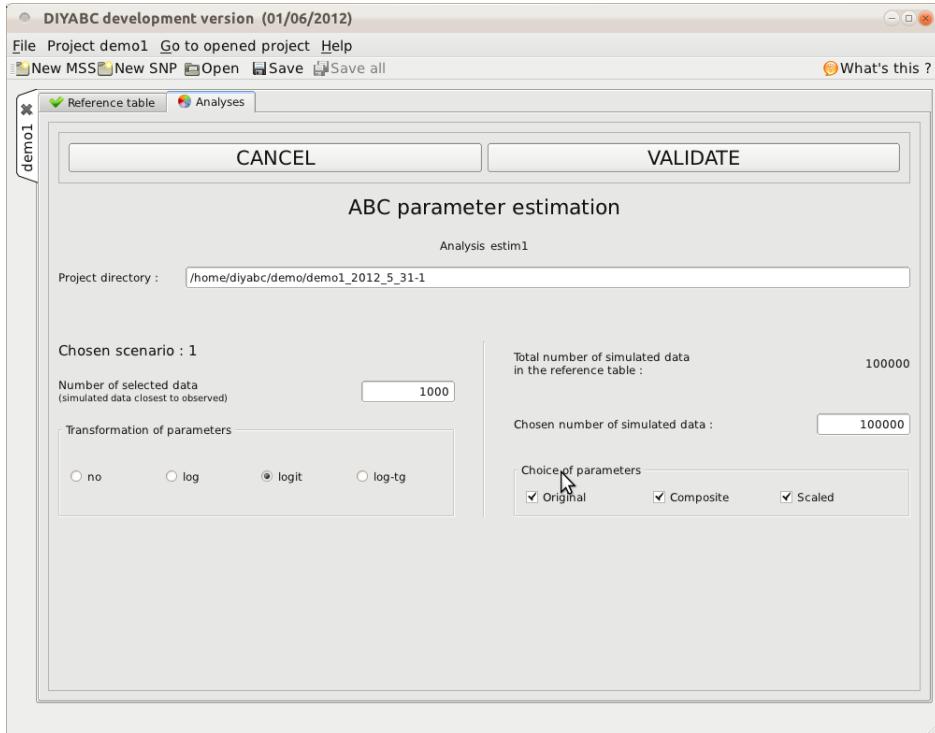
7

8

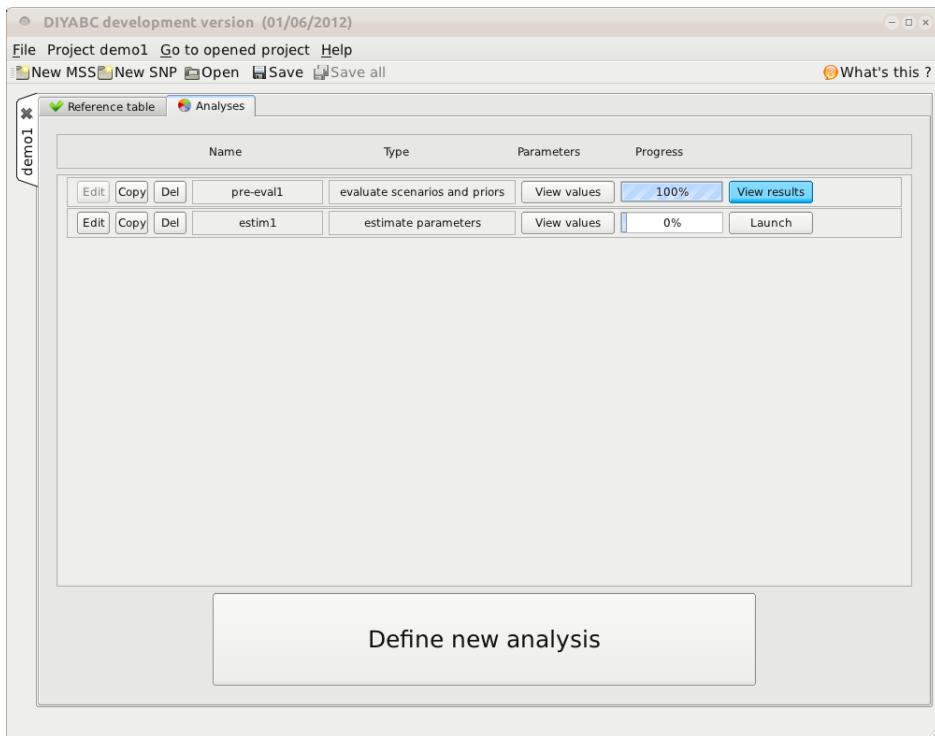
9 We get then the following screen in which we can make several choices :

- 10 • on the left hand side, we can choose the number of simulated datasets that will be used for the
11 local linear regression (cf section 2.1).
- 12 • below, we can select the transformation of parameter values that can generally improve the results
13 (default = logit transform).

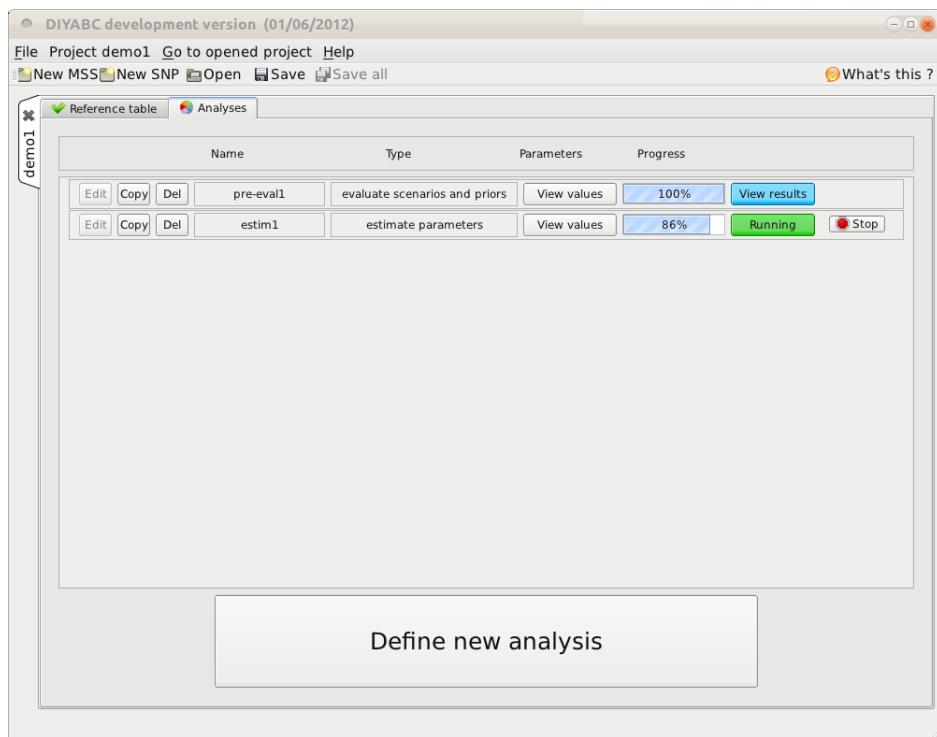
- 1 • on the right hand, we can truncate the reference table to a specified number of datasets.
- 2 • eventually, estimations can be performed either on raw parameters, and/or combinations of parame-
- 3 ters that are generally more estimable. *Composite* parameters are products of effective population
- 4 sizes or times by mean mutation rate whereas *Scaled* parameters are ratios of effective population
- 5 sizes or times by mean effective population size.



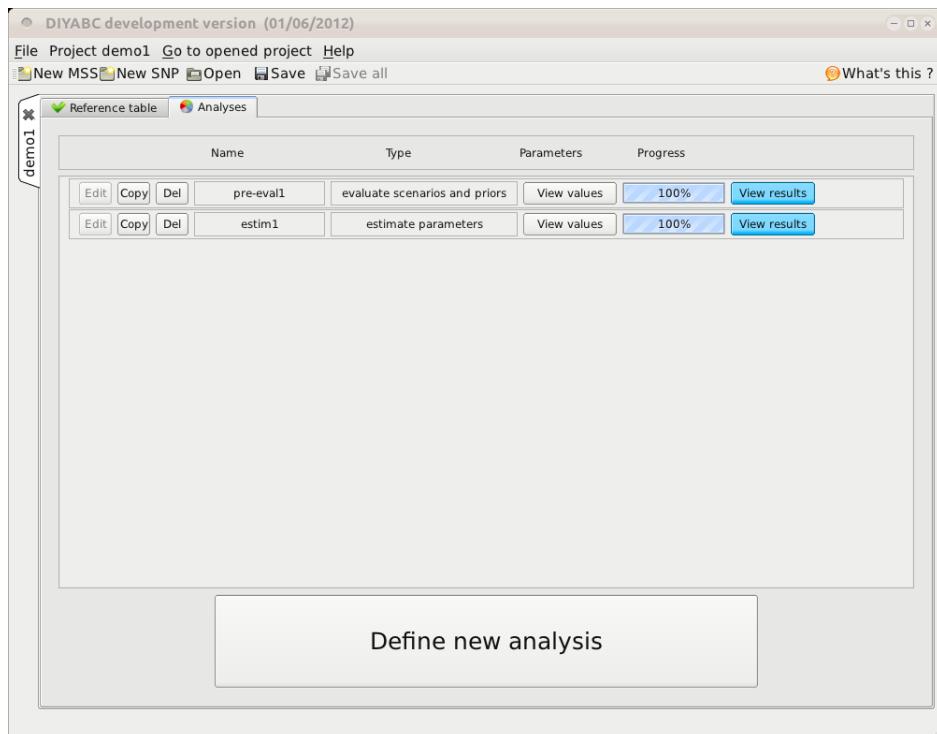
6
7 Apart from the number of closest datasets that we set at 10,000, we keep all other default values and
8 click on the **VALIDATE** button. We get back to the Analysis control panel which now looks like this:



1 We click on the **Launch** button. The analysis progress is now visible :

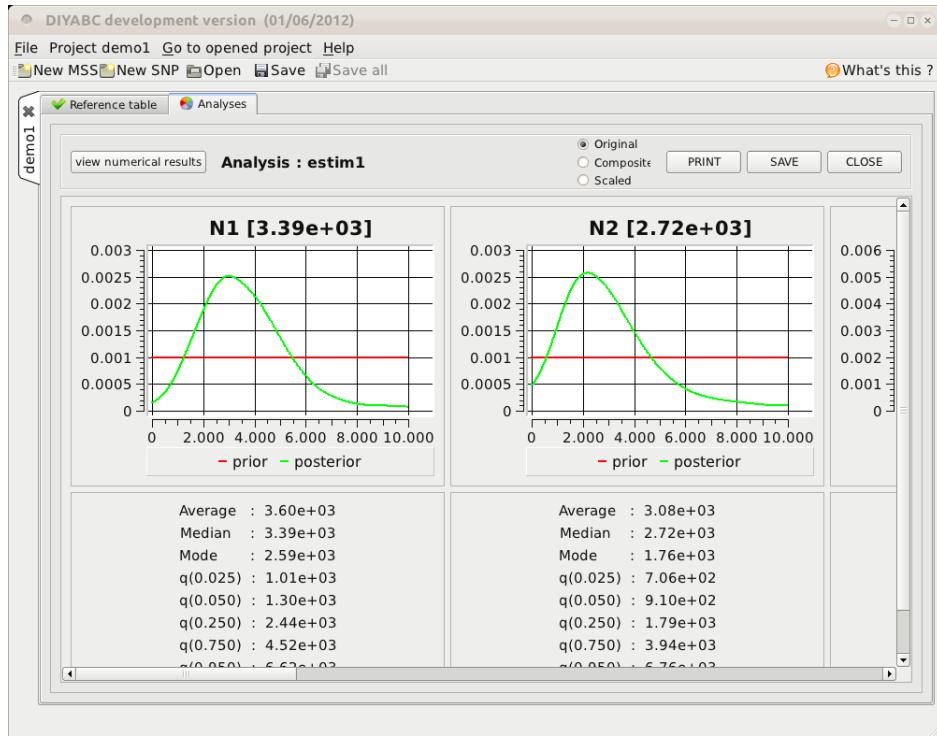


5 As long as the analysis is not terminated, we could stop it by clicking on the **Stop** button. Once this
6 second analysis is finished, we can view its results by clicking on the **View results** button :



10 Let's have a look :

11

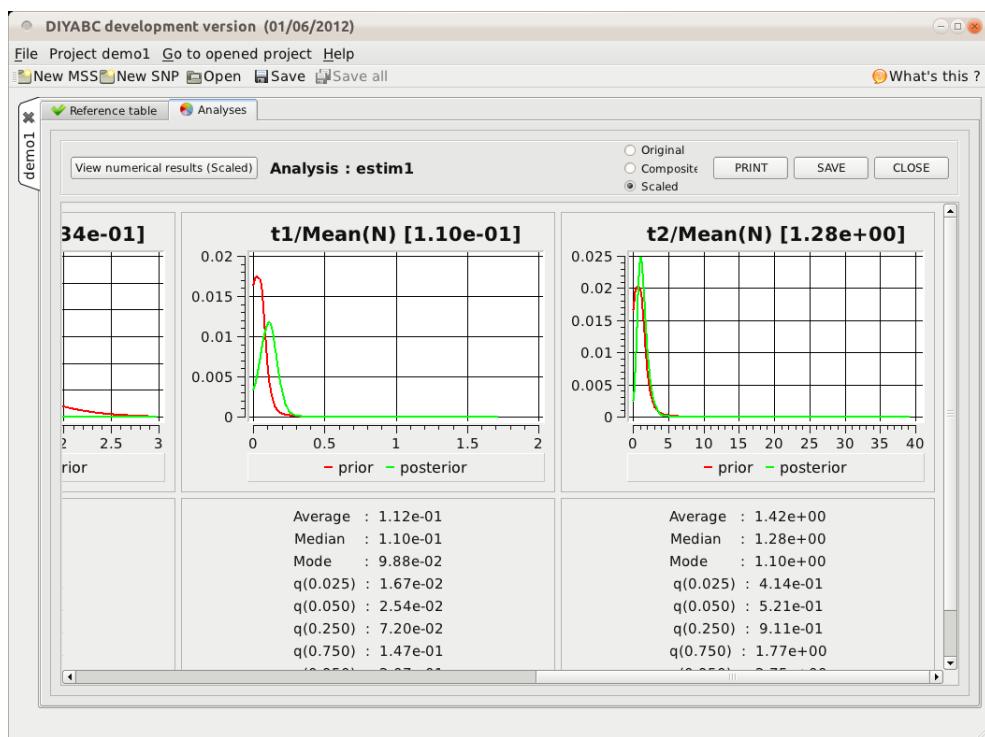


In the scrolling window, we get graphics showing the prior (red curve) and posterior (green curve) distributions of all parameters. Below each graphics are statistics (mean, median, mode and quantiles) of the posterior distribution. The latter are grouped in a table that appears when clicking on the upper left **view numerical results** button, showing this :

Parameter	mean	median	mode	q025	q050	q250	q750	q950	q975
N1	3.60e+03	3.39e+03	2.59e+03	1.01e+03	1.30e+03	2.44e+03	4.52e+03	6.62e+03	7.55e+03
N2	3.08e+03	2.72e+03	1.76e+03	7.06e+02	9.10e+02	1.79e+03	3.94e+03	6.76e+03	7.94e+03
N3	1.39e+03	8.33e+02	4.00e+02	1.44e+02	1.92e+02	4.56e+02	1.59e+03	4.78e+03	6.58e+03
t1	2.84e+02	2.91e+02	2.94e+02	4.29e+01	6.80e+01	1.93e+02	3.82e+02	4.70e+02	4.85e+02
t2	3.70e+03	3.30e+03	2.22e+03	8.78e+02	1.07e+03	2.12e+03	4.86e+03	7.86e+03	8.71e+03
umic_1	1.54e-04	1.31e-04	1.07e-04	9.90e-05	1.01e-04	1.13e-04	1.65e-04	2.81e-04	3.56e-04
Pmic_1	1.60e-01	1.49e-01	1.06e-01	1.00e-01	1.02e-01	1.21e-01	1.90e-01	2.53e-01	2.68e-01
snimic_1	8.42e-07	1.69e-07	1.10e-08	1.08e-08	1.19e-08	3.58e-08	8.24e-07	4.31e-06	6.12e-06

We go back to the previous screen by clicking the **OK** button.

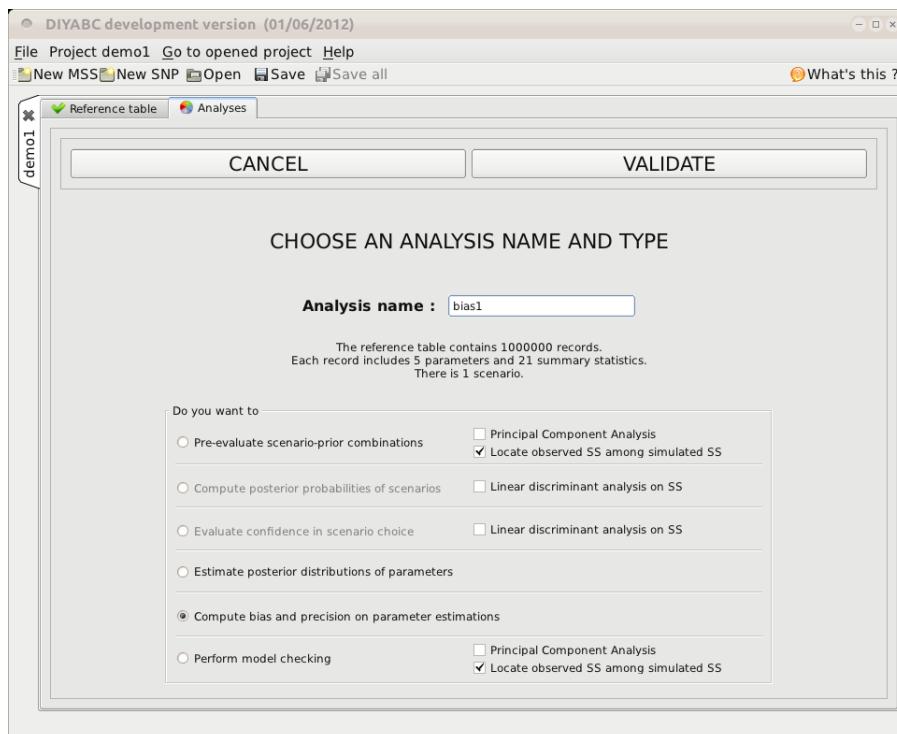
We can also have results for *Composite* or *Scaled* parameters. Below is an example of *Scaled* time parameters obtained by clicking on the *Scaled* radio button and scrolling the graphs window to the right:



4
5

6 3.5.2 Bias and precision

7 Let's define a new analysis (click on the **Define new analysis** button) and choose the option **Compute bias and precision on parameter estimations**. We give it the name **bias1** :



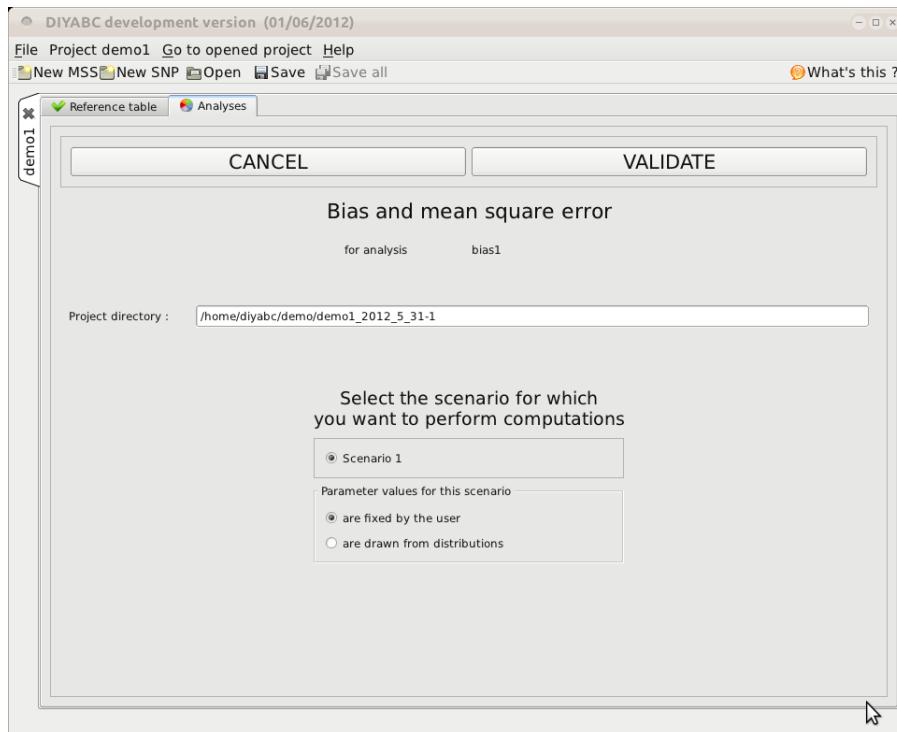
10
11
12

In this kind of analysis, pseudo-observed data are simulated with known values of parameters copying

- 1 the exact configuration of the observed dataset in terms of sample sizes (taking into account missing data)
 2 and are submitted to the same ABC estimation process. If we assume that the evolutionary scenario
 3 is correct, the comparison of real and estimated values of parameters provide some information of the
 4 precision of the estimation process.

5 We validate and get this screen :

6

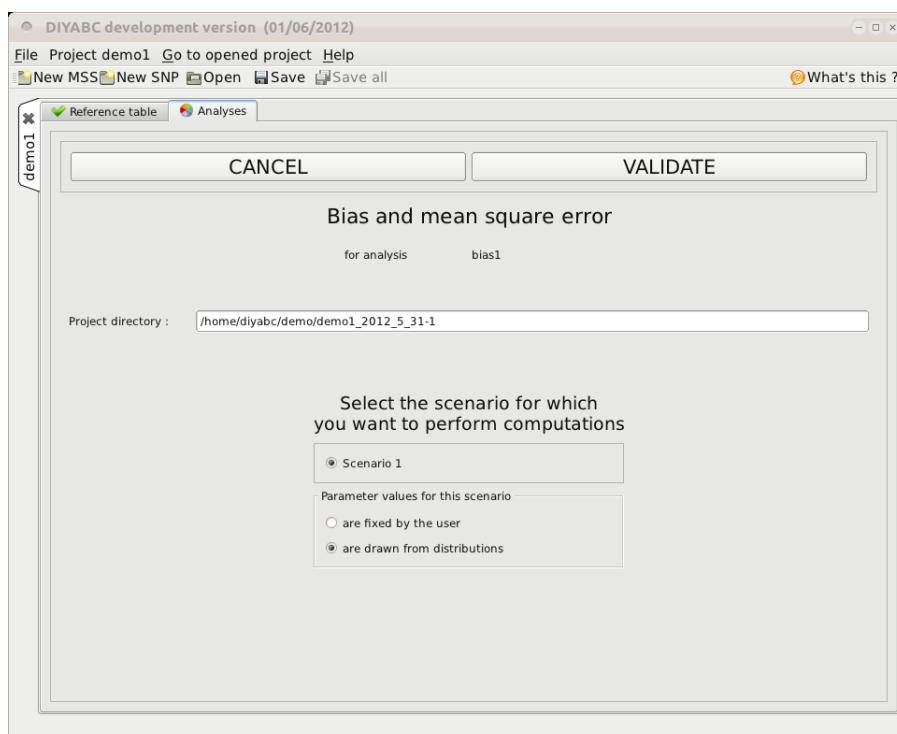


7

8

9 We choose to draw parameter values from distributions :

10

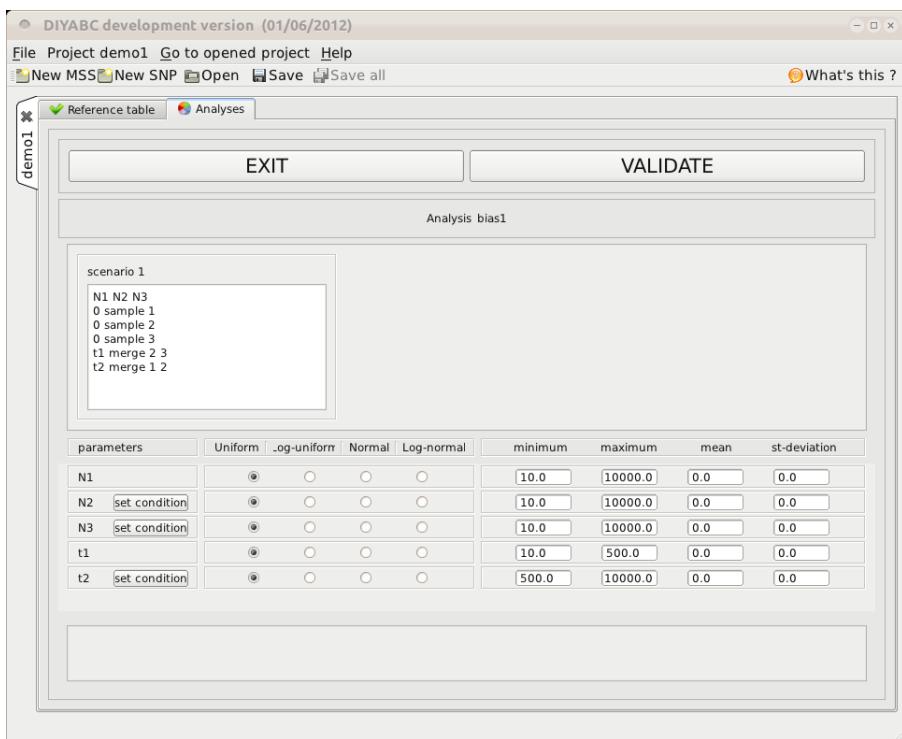


11

12

13 Clicking on the **VALIDATE** button, we get this screen which allows us to choose distributions.

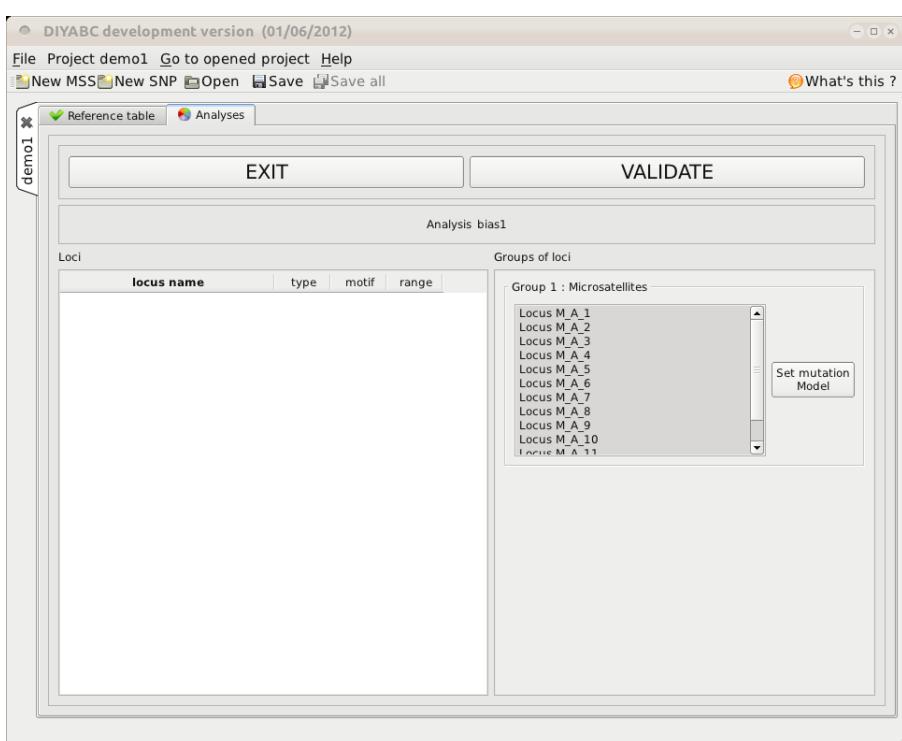
14



1

2

3 By default, the screen suggests the prior distributions that have been used to build the reference
 4 table. However, these distributions can be edited if necessary. We decide not to change them and click
 5 on **VALIDATE** which brings us to the following screen :

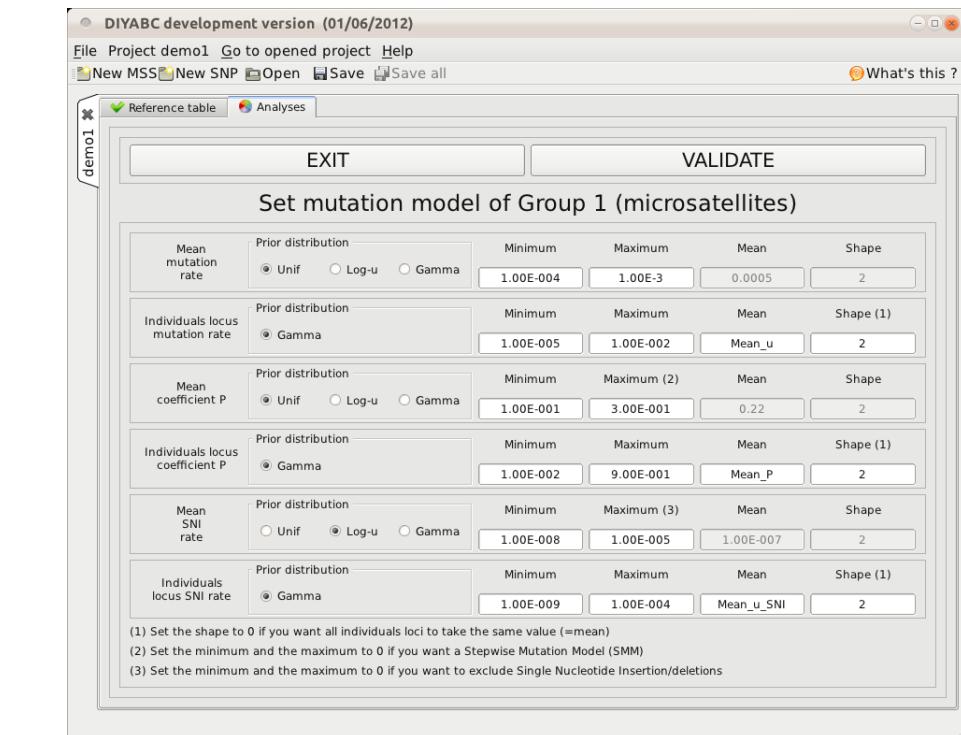


7

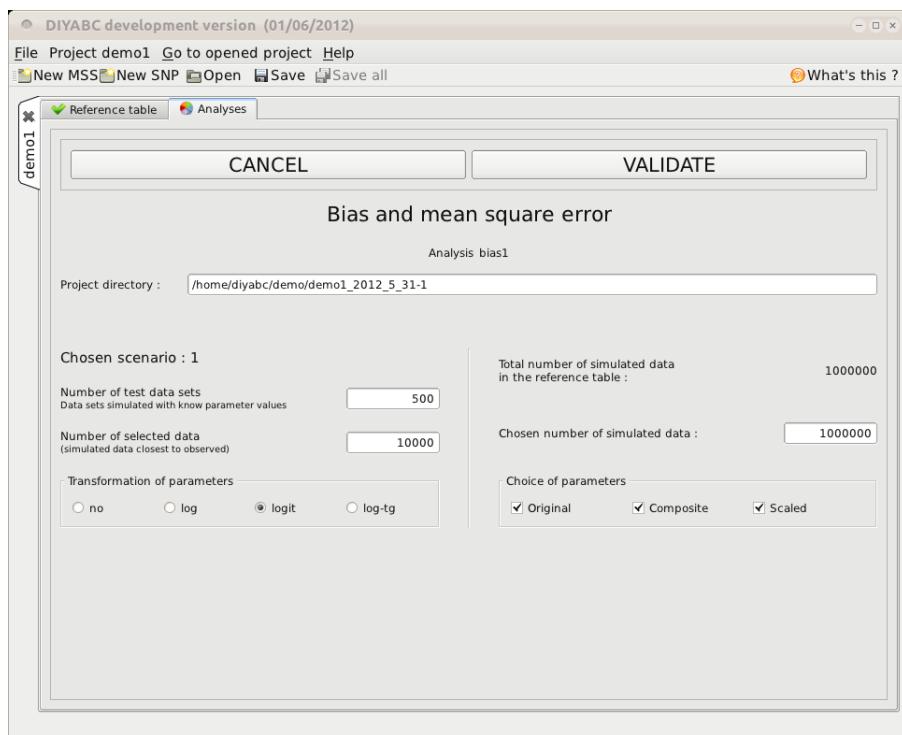
8

9 If we want to keep the same distributions for mutation parameters as when building the reference
 10 table, we just click on **VALIDATE**. If we need to change them, we click on **Set mutation model** which
 11 would bring the following screen :

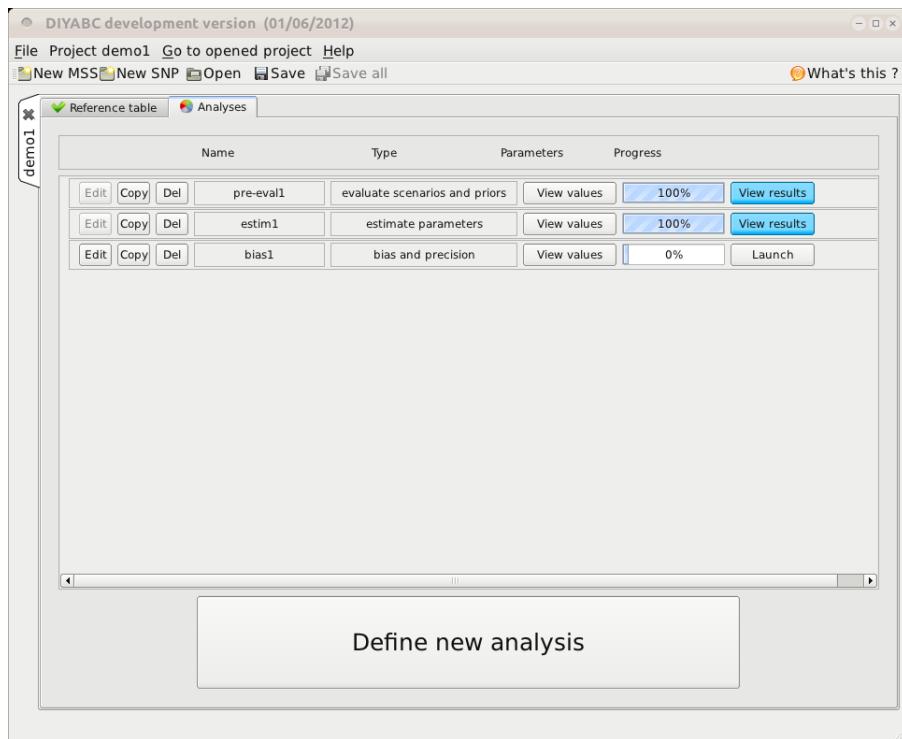
12

1
2
3
4

After validating twice, we get the last screen necessary to define this kind of analysis :

5
6
7
8
9

This screen is similar to that for parameter estimation (see section 3.5.1). After validating, we get back to the analysis panel with a third analysis defined :

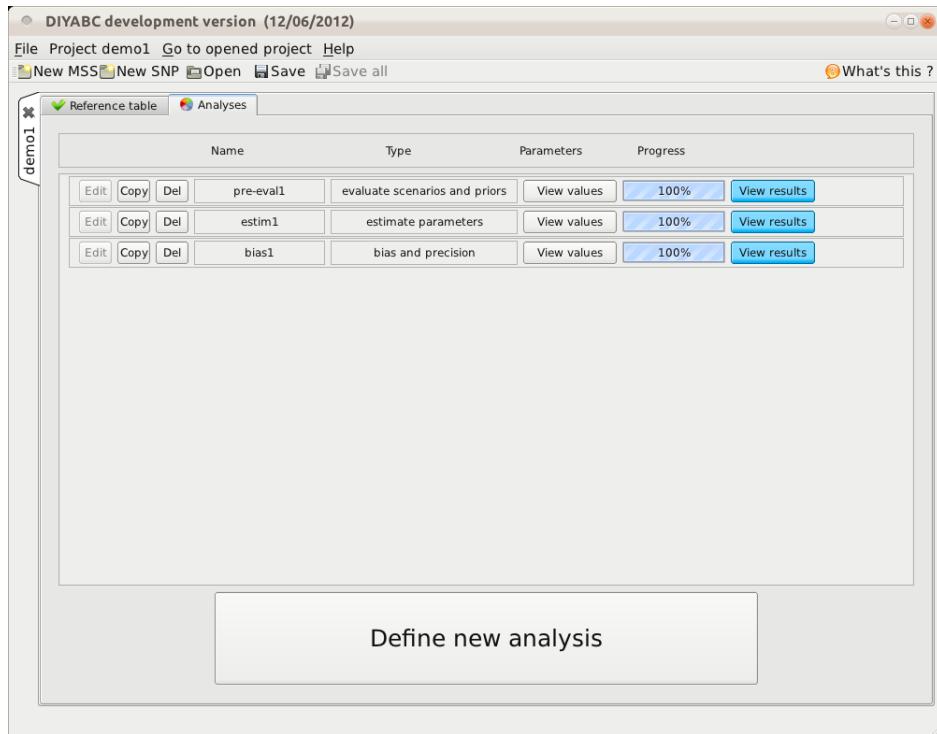


1

2

3 The analysis takes some time to run compared to the previous one, because it simulates hundreds
 4 datasets and on each one, a full ABC estimation is performed. Then after some time, the analysis is
 5 finished:

6

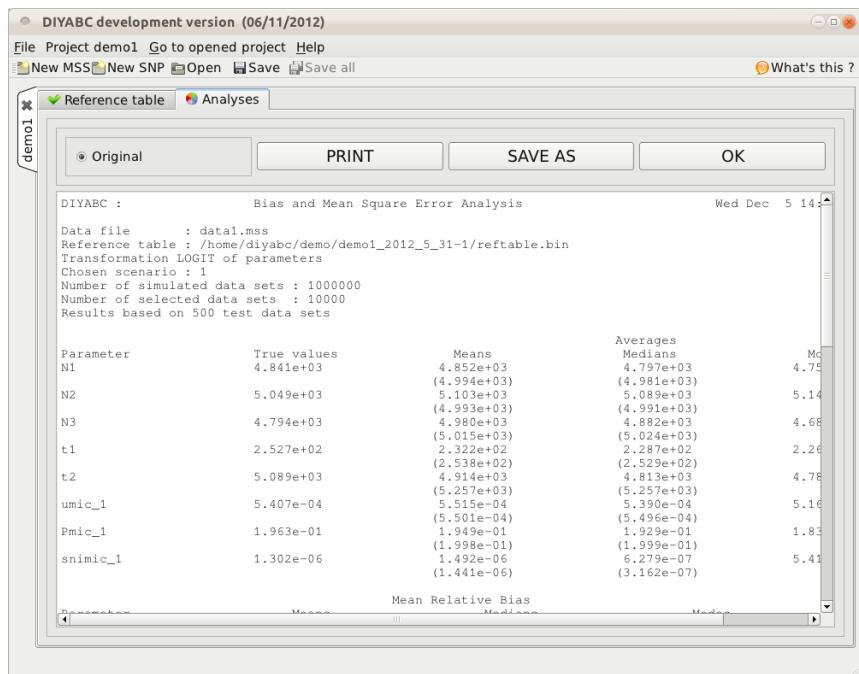


7

8

9 To view results, we click on the [View results](#) button.

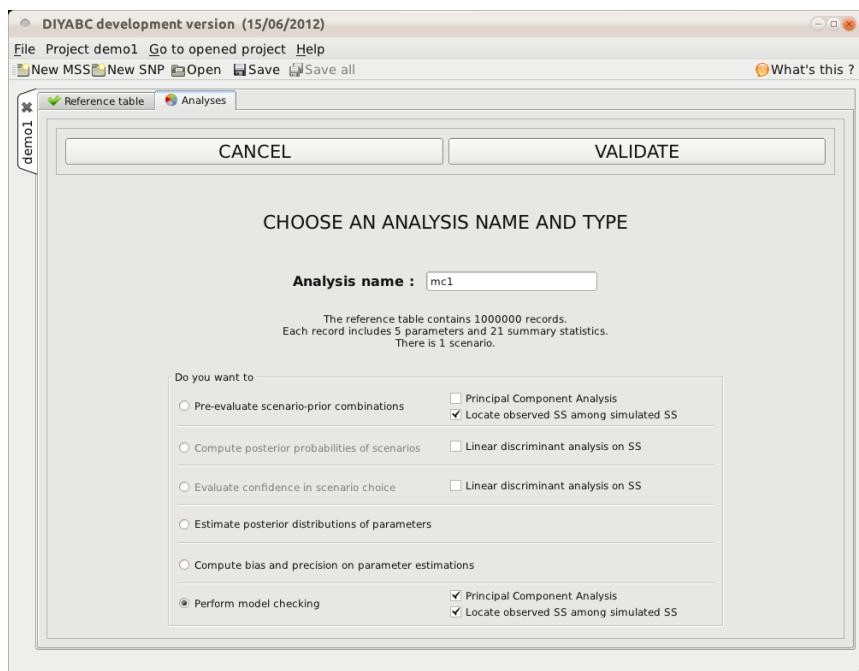
The results are visible in a scrolling window :



Note that there are two values given for each statistics. The upper value is that of the statistics computed from the *posterior* distribution of parameters, *i.e.* **using the genetic information provided by data**. The lower value, noted between parentheses, is that of the statistics computed from the *prior* distribution of parameters, *i.e.* **NOT using the genetic information provided by data but only that contained in prior distributions**.

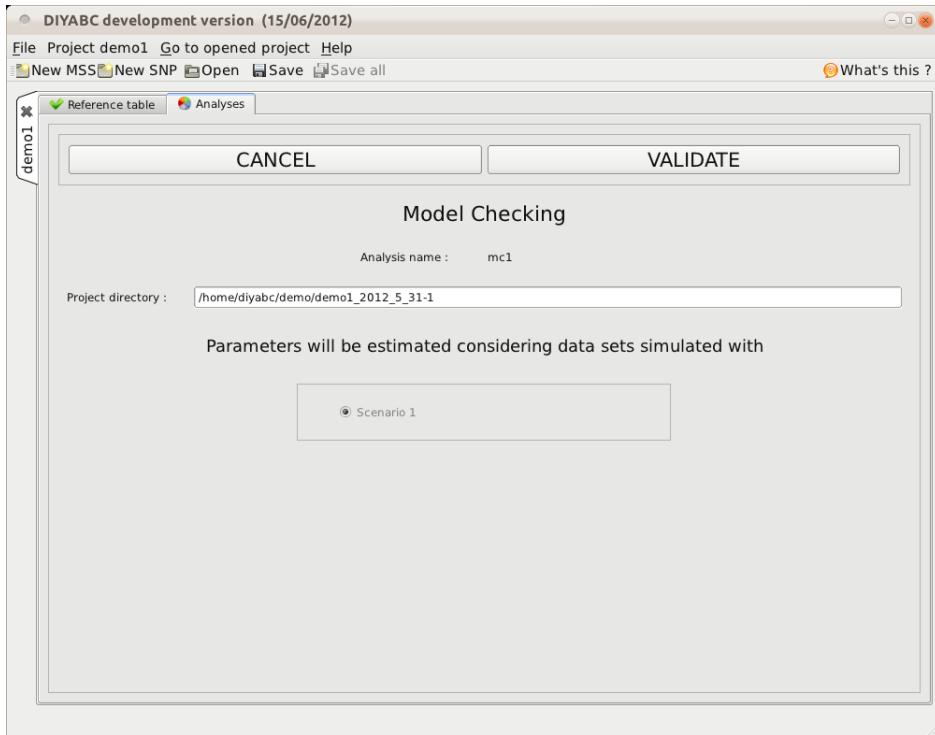
10 3.5.3 Model Checking

11 We now define another type of analysis called **Model Checking** which is used to evaluate how well the
12 scenario and priors of parameters fit the data summarized by summary statistics. This is the last option
13 on the following screen :

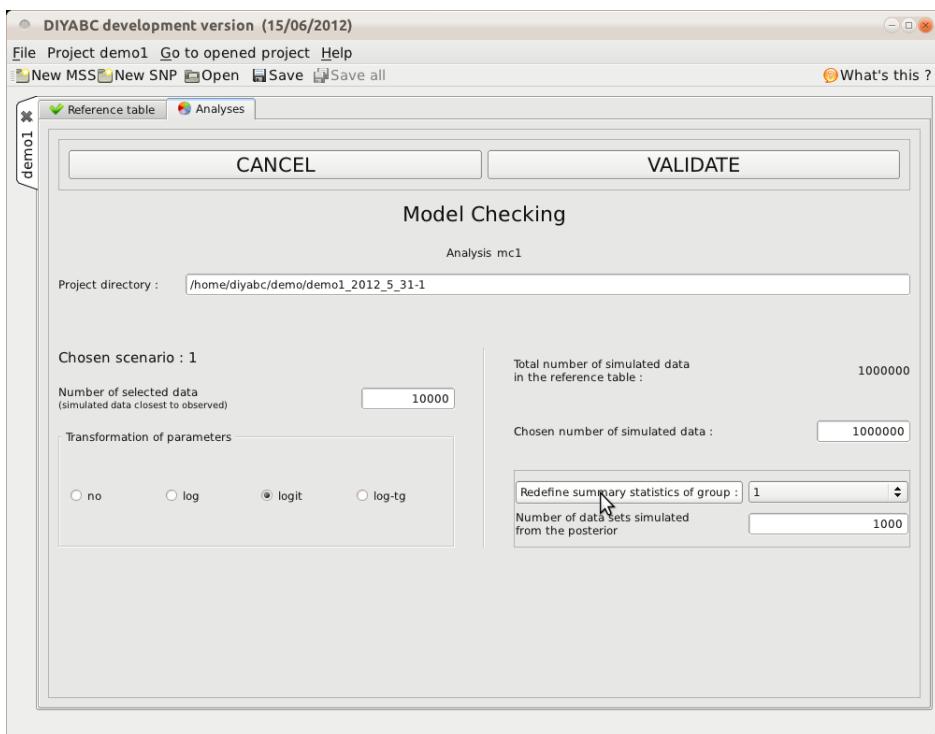


We call this analysis `mcl` and check the box to get a PCA performed. This PCA is computed in the same way compared to that of the first option (`Pre-evaluate scenario prior combinations`). However, new datasets simulated with parameters drawn from the posterior distributions of parameters are also represented on the different planes of the PCA (but not taken in the PCA computation).

We validate the above screen and get the usual next screen :



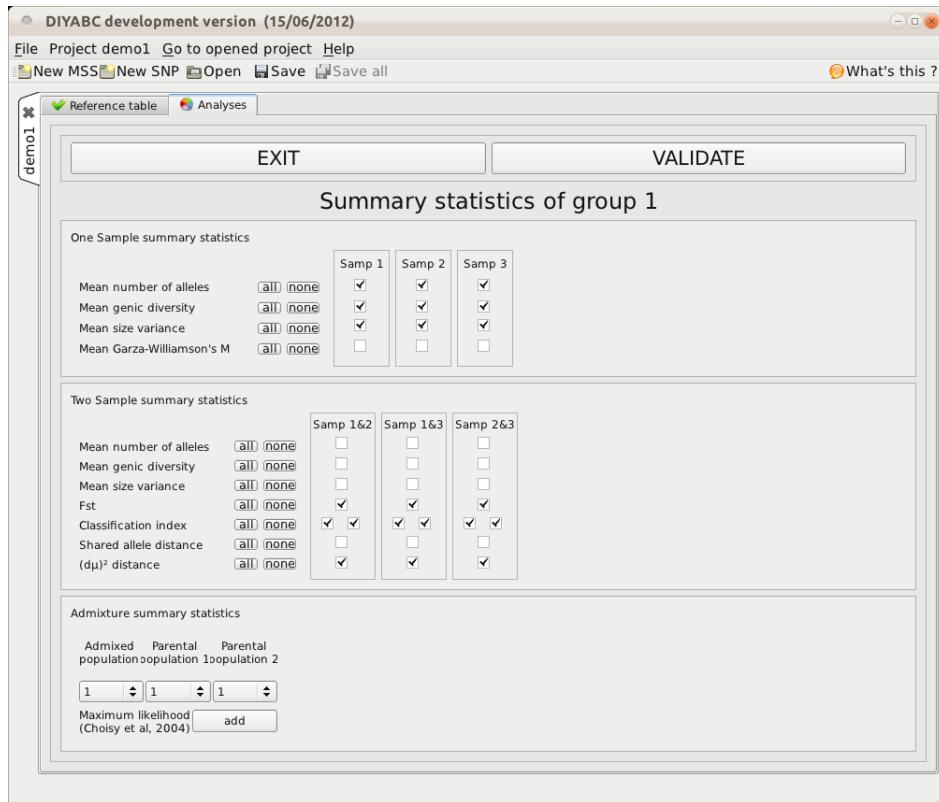
that we just validate to get the last screen :



In this screen which we have already seen, there is a new panel (bottom right) in which we can choose the number of datasets that we want to simulate from the posterior distributions of parameters. There is

- 1 also a button **Redefine summary statistics of group:** shown by the pointer. This button allows to change
 2 the set of summary statistics (for a given group of loci chosen through the drop list on the right). Clicking
 3 on this button opens up the usual following screen in which, by default, are checked the summary
 4 statistics in the reference table.

5

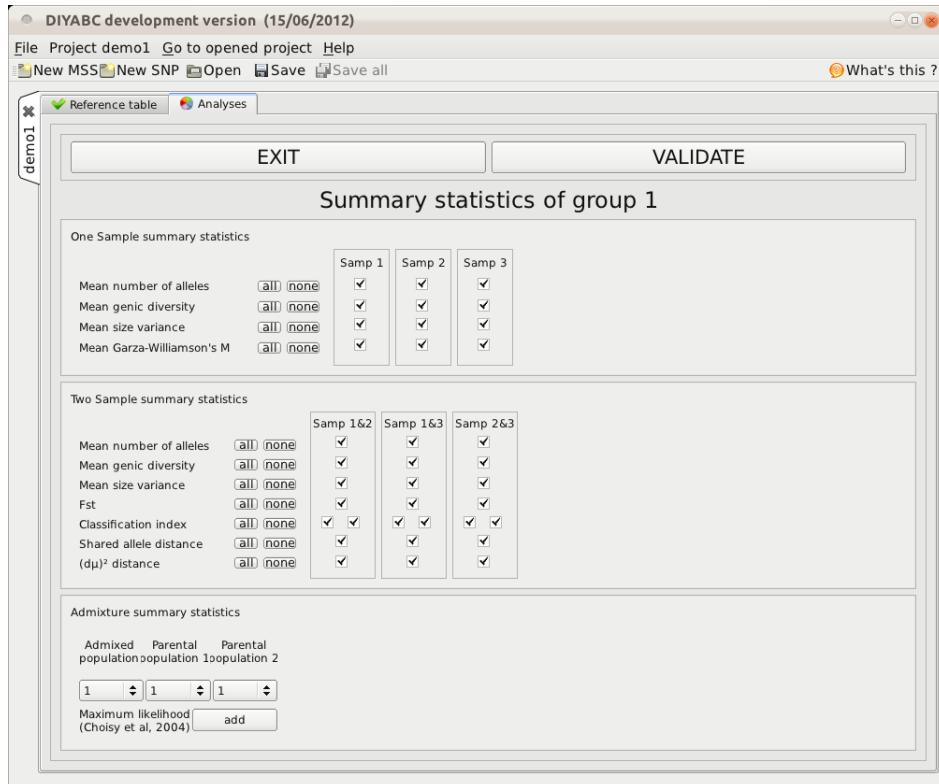


6

7

We decide to use all one-sample and two-sample summary stats :

8



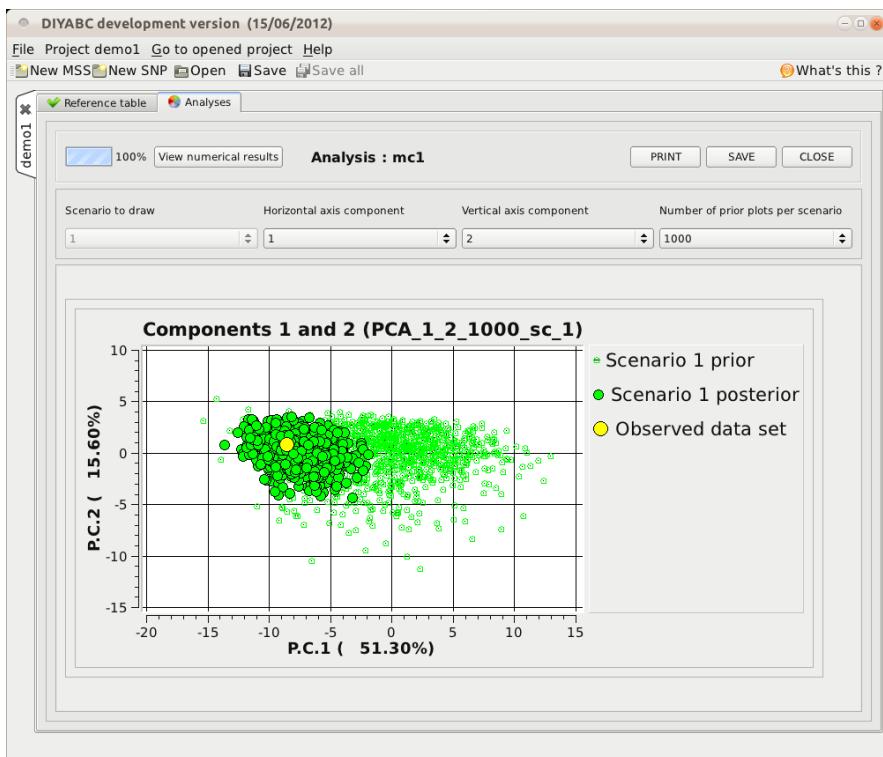
9

10

1

2 Note that when the set of summary statistics is changed (as here), it is necessary to also simulate a
 3 large number of datasets using priors to get a new reference table.

4 We validate twice and launch the analysis. When it is finished, we click on the **View results** button and
 5 get this screen:



7

8

9 Clicking on the **View numerical results** leads to the following screen which provides, for each individual
 10 summary statistics, the value in the observed dataset as well as the proportion of data sets (simulated
 11 from the posterior) that have a value lower than the observed data set.

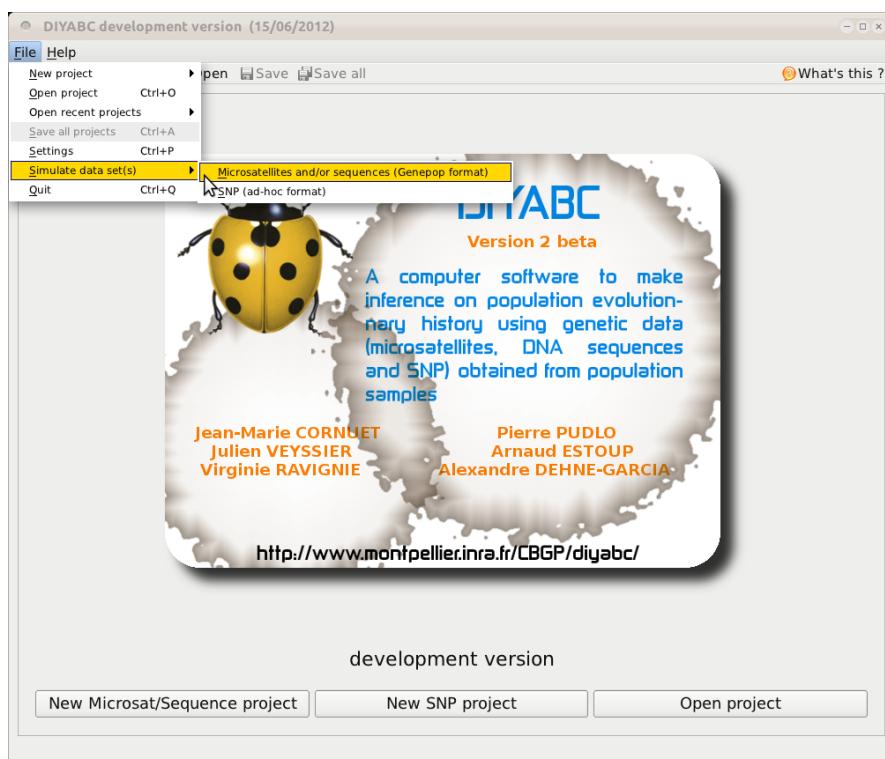
POSTERIOR CHECKING		
DIYABC :		Tue Jun 19 09:06:14 2012
Data file	observed	proportion
Reference table	value	(simulated<observed)
Chosen scenario	: 1	
Number of simulated data sets used to compute posterior	: 1000000	
Number of simulated data sets used in the local regression	: 10000	
Number of data sets simulated from the posterior	: 1000	
Transformation of parameters	: Logit	
summary statistics	observed	proportion
NAL_1_1	4.0933	0.3590
NAL_1_2	3.5000	0.2760
NAL_1_3	2.7500	0.3755
HET_1_1	0.4562	0.2715
HET_1_2	0.4782	0.5535
HET_1_3	0.3754	0.4290
VAR_1_1	1.3355	0.4480
VAR_1_2	1.0758	0.4035
VAR_1_3	0.6574	0.2665
MGW_1_1	1.0000	0.5800
MGW_1_2	1.0000	0.5640
MGW_1_3	1.0312	0.6845
N2P_1_1&2	4.9167	0.2050
N2P_1_1&3	4.7500	0.2855
N2P_1_2&3	3.8333	0.2405
H2P_1_1&2	0.5540	0.3560
H2P_1_1&3	0.5235	0.2995
H2P_1_2&3	0.4523	0.4545
V2P_1_1&2	1.3393	0.2505
V2P_1_1&3	1.1330	0.1810
V2P_1_2&3	0.8998	0.3075
FST_1_1&2	0.2703	0.5000
FST_1_1&3	0.3404	0.5645
FST_1_2&3	0.1064	0.4720
LIK_1_1&2	1.7198	0.3025

13

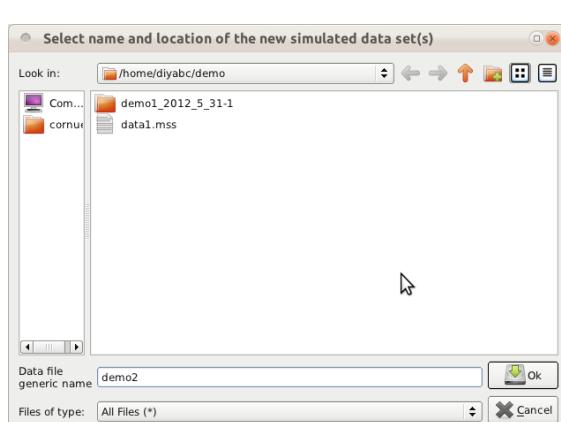
1
2 Notice that in this computation, values that are in the interval $[s_{obs} - 0.001, s_{obs} + 0.001]$ are counted
3 for one half those that are outside the interval. This explains why the fourth digit of the proportion can
4 be 0 or 5 while having simulated 1000 data sets.

5 3.6 Simulating data sets

6 The *DIYABC* program can also be used to simulate data sets, either microsatellite and/or DNA se-
7 quence data sets using our Genepop format, or SNP data sets using our specific format. This option is
8 reachable through the main File menu as shown below :

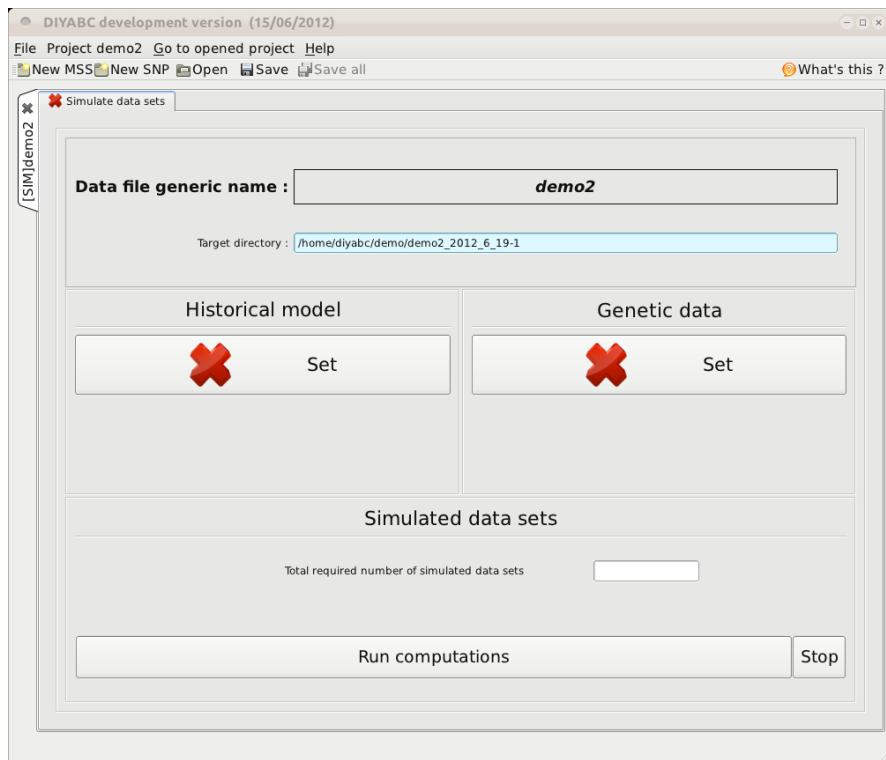


10
11 Clicking on e.g. the Microsatellites and/or sequences (Genepop format) opens up a dialog window in
12 which one can choose the directory into which will be located the project and the future data files :
13

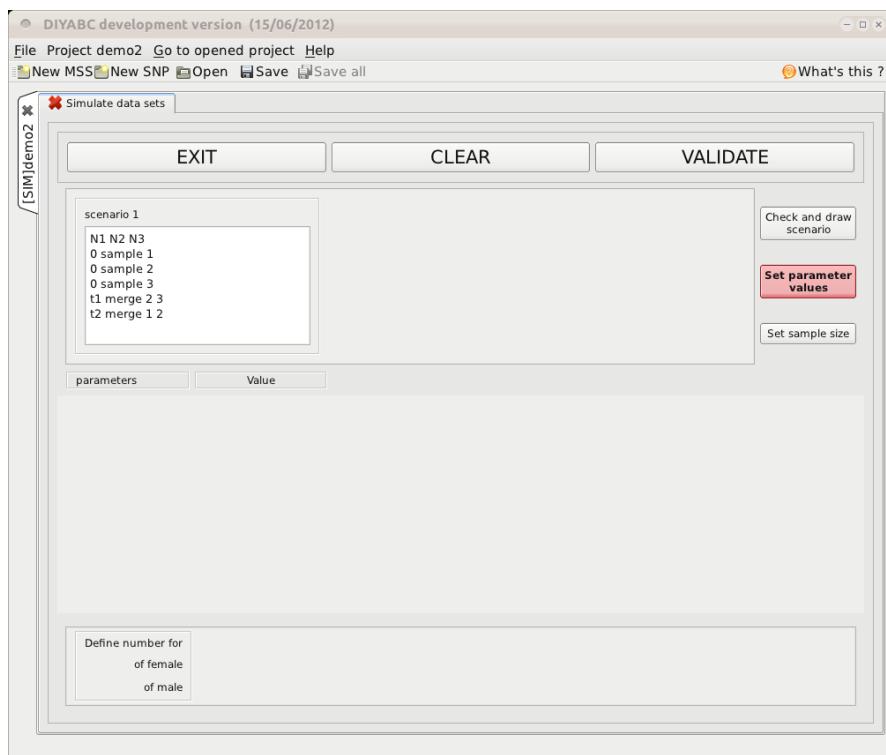


14
15 Above, we decided to call **demo2** this new directory and to locate it in the `home/diyabc/demo` directory.
16 Clicking on **OK** leads to usual screen:
17

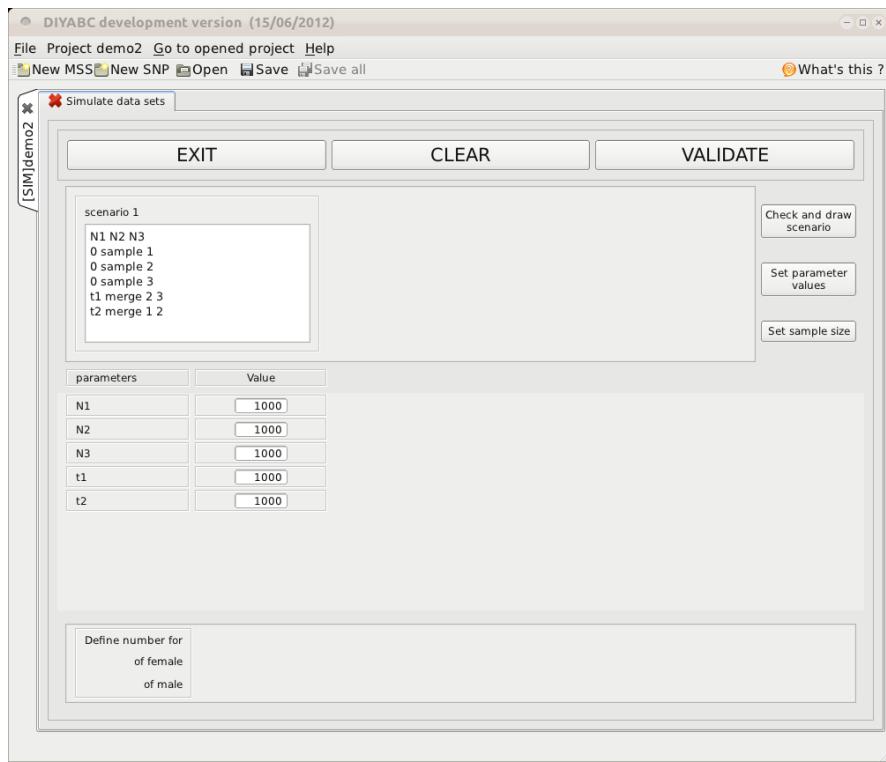
18



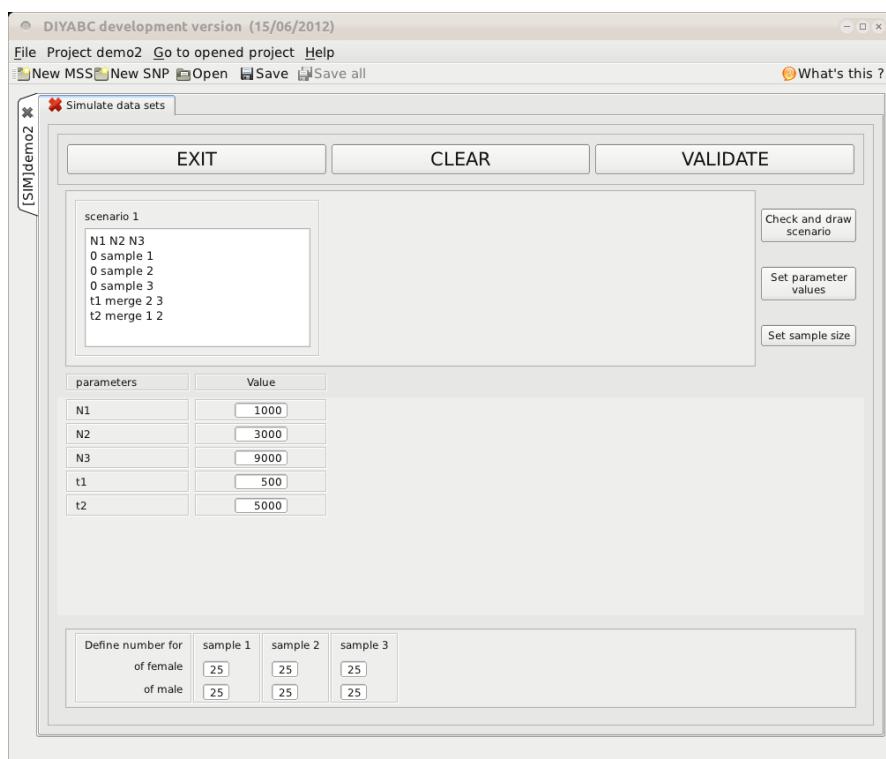
We first inform the historical model clicking on the **Set** button under **Historical model**. We edit the scenario box as below:



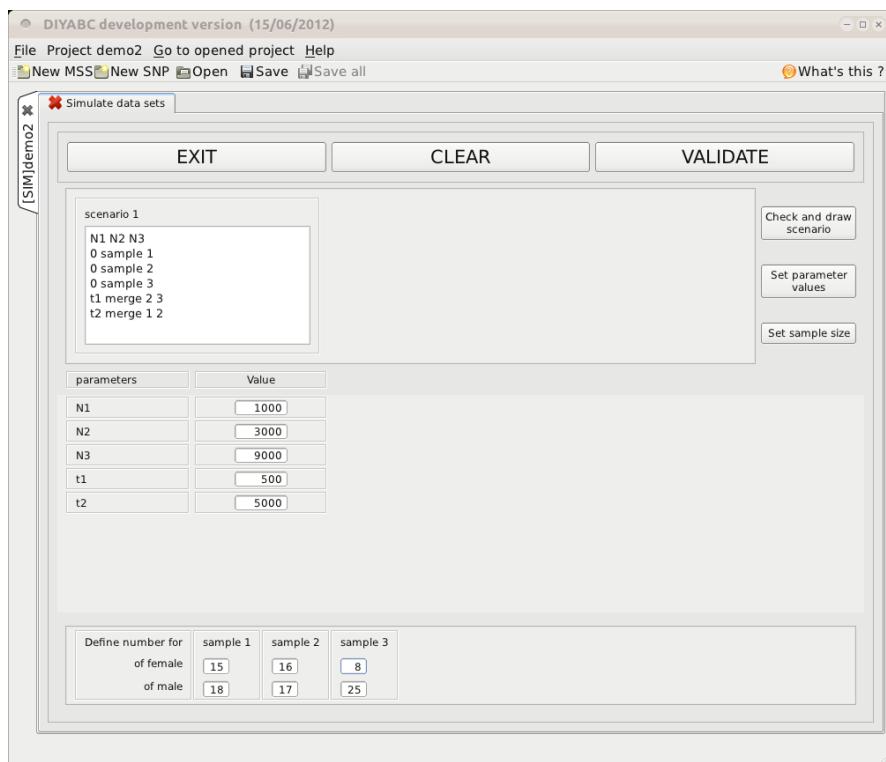
We click on the **Set parameter values** button. Arbitrary default values appear :



1
2 We change these values according to our needs and we click the **Set sample size** button, getting this
3 screen:



7 We input the needed sample sizes as below :
8

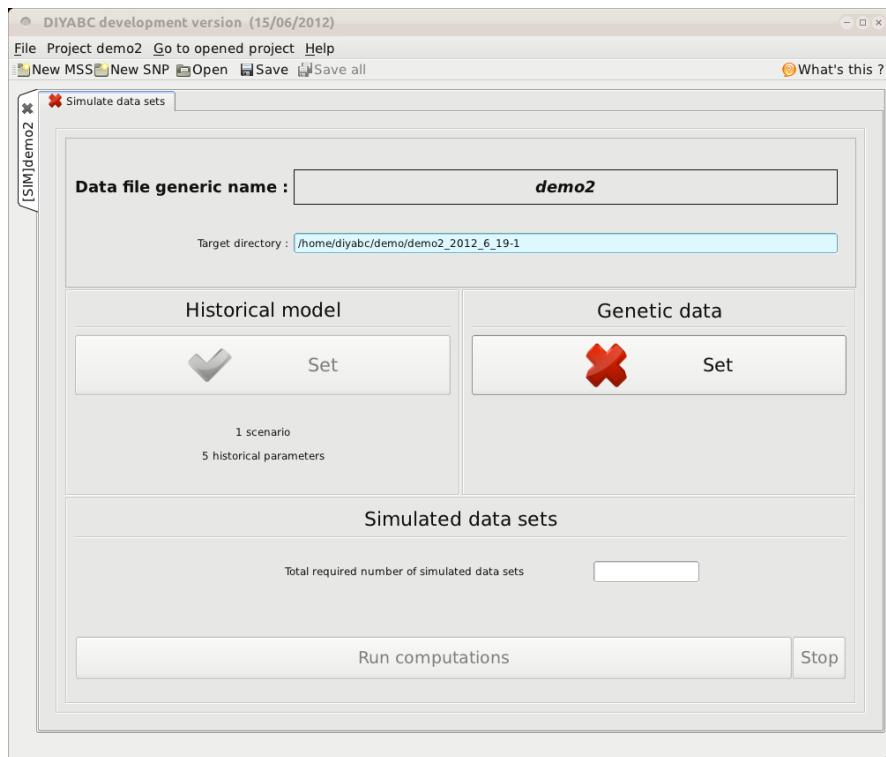


1

2

3 Clicking on the **VALIDATE** button, we get back to the previous screen showing that the Historical
4 model is now completed:

5

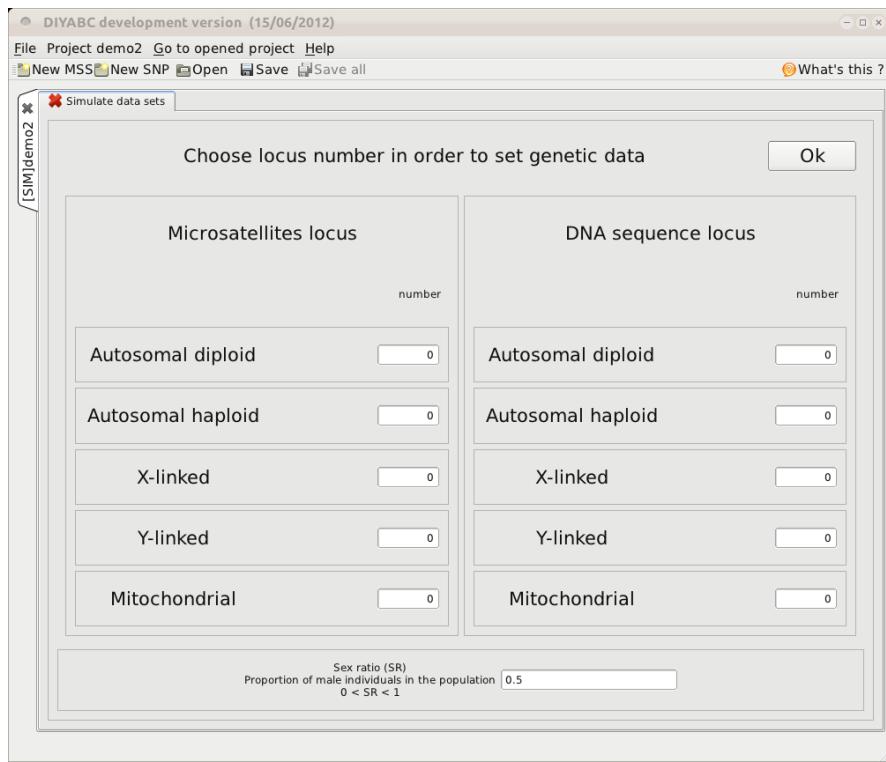


6

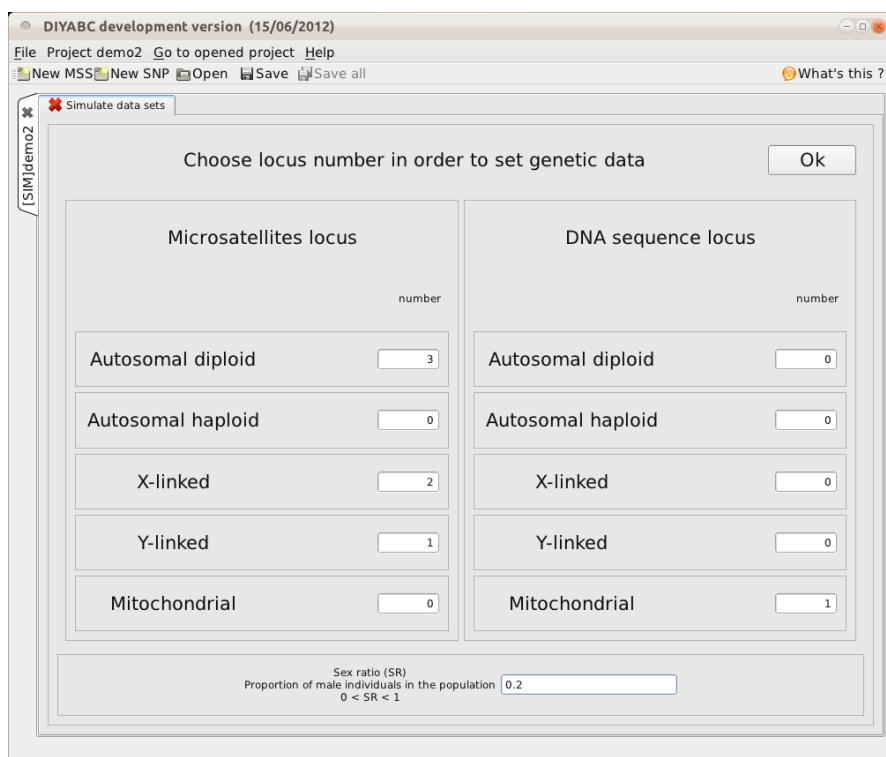
7

8 We have now to complete the Genetic data (click on the **Set** button under Genetic data). The fol-
9 lowing screen appears:

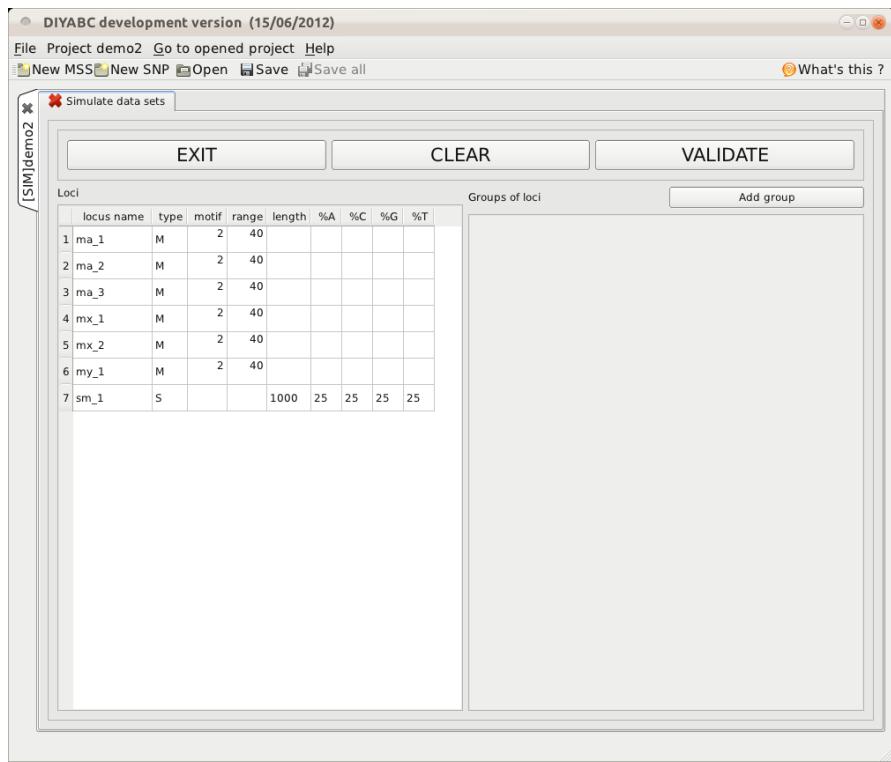
10



We want a data set including three autosomal, two X-linked and one Y-linked diploid microsatellite loci and one mitochondrial sequence. We also need a sex ratio of one male for four females :



We click on the **OK** button and get the following screen :

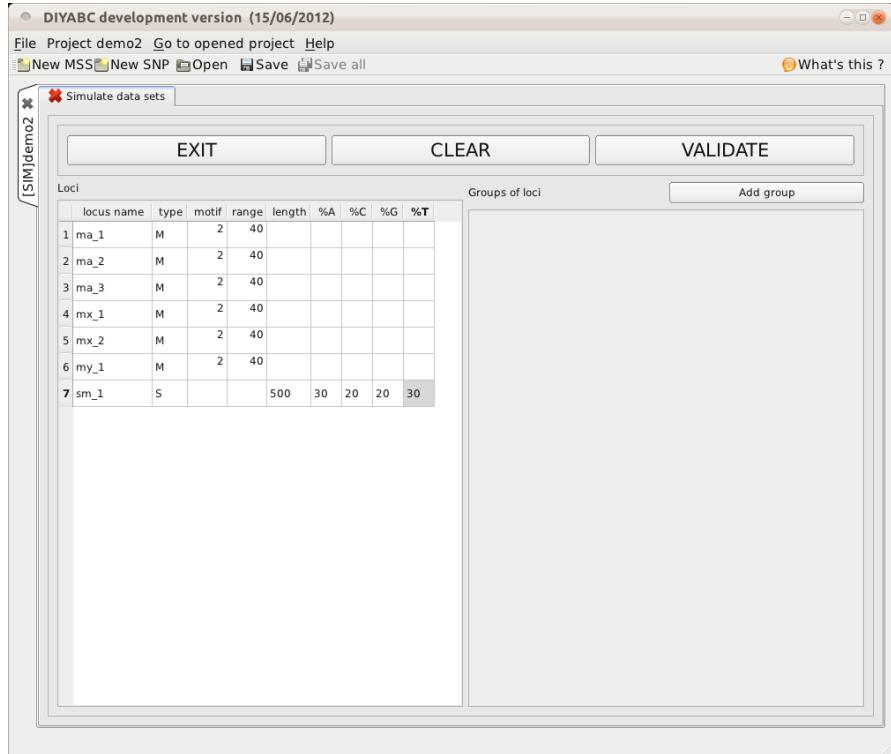


1

2

3 Our mitochondrial DNA sequence is only 500 nucleotides long and there is a slight excess of A+T
 4 (60%). We edit the corresponding cells :

5

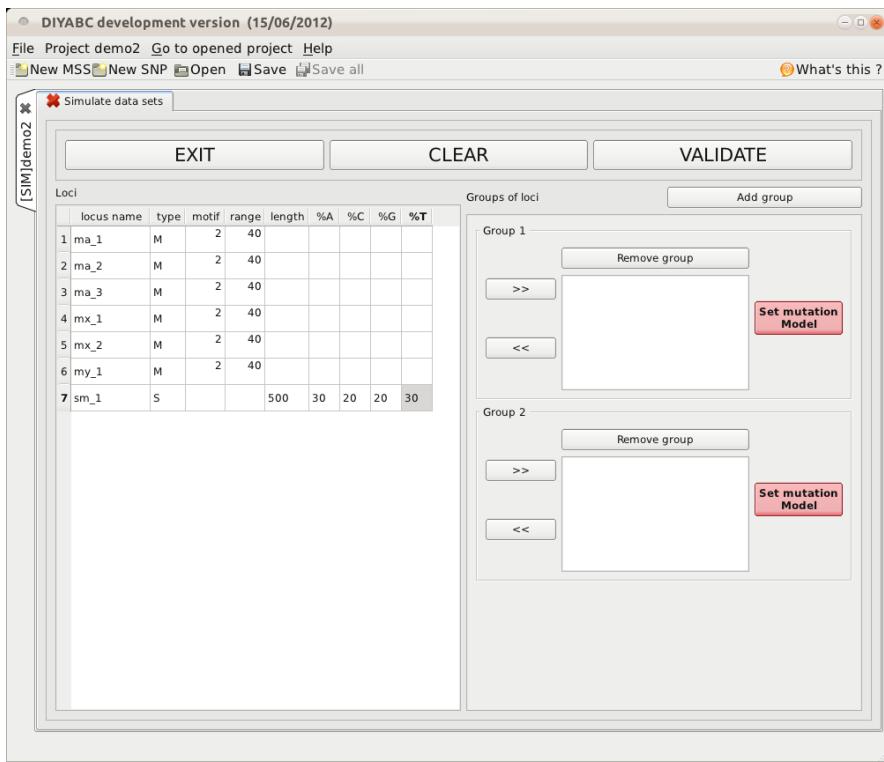


6

7

8 Since mutation models are different for microsatellites and DNA sequences, we define two groups by
 9 clicking twice on the **Add group** button :

10

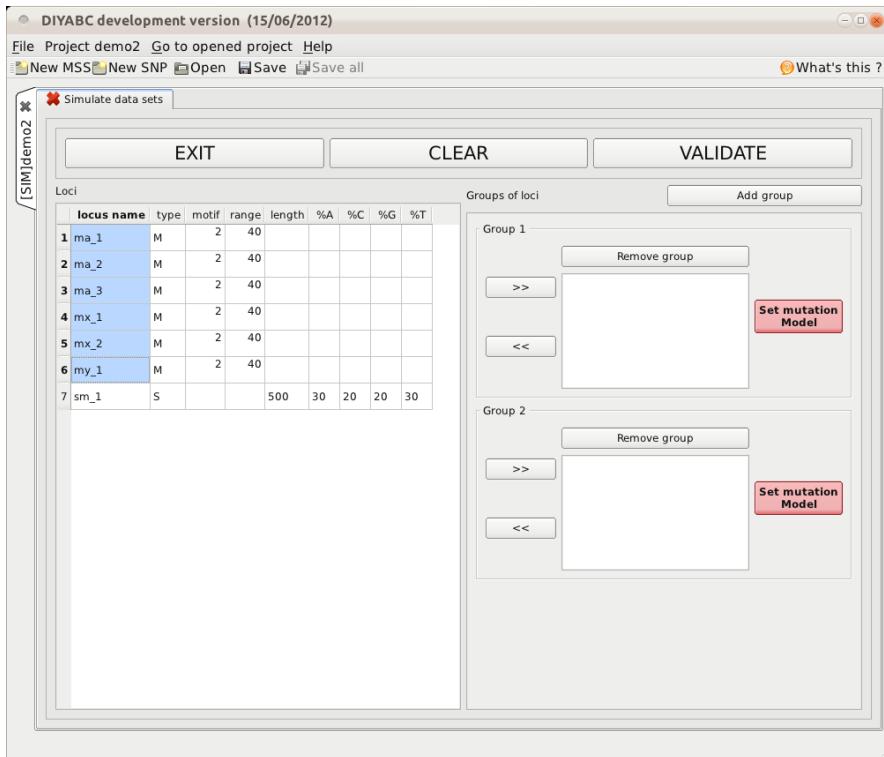


1

2

3 We select the 6 microsatellite loci by clicking on the first locus name cell and shift-clicking on the
4 sixth locus name cell :

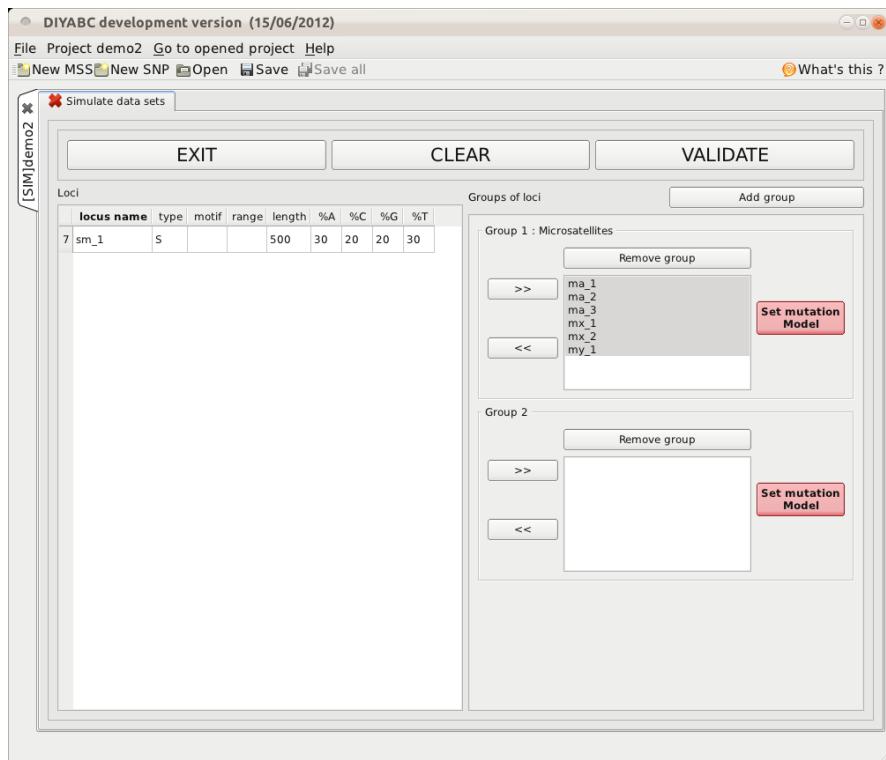
5



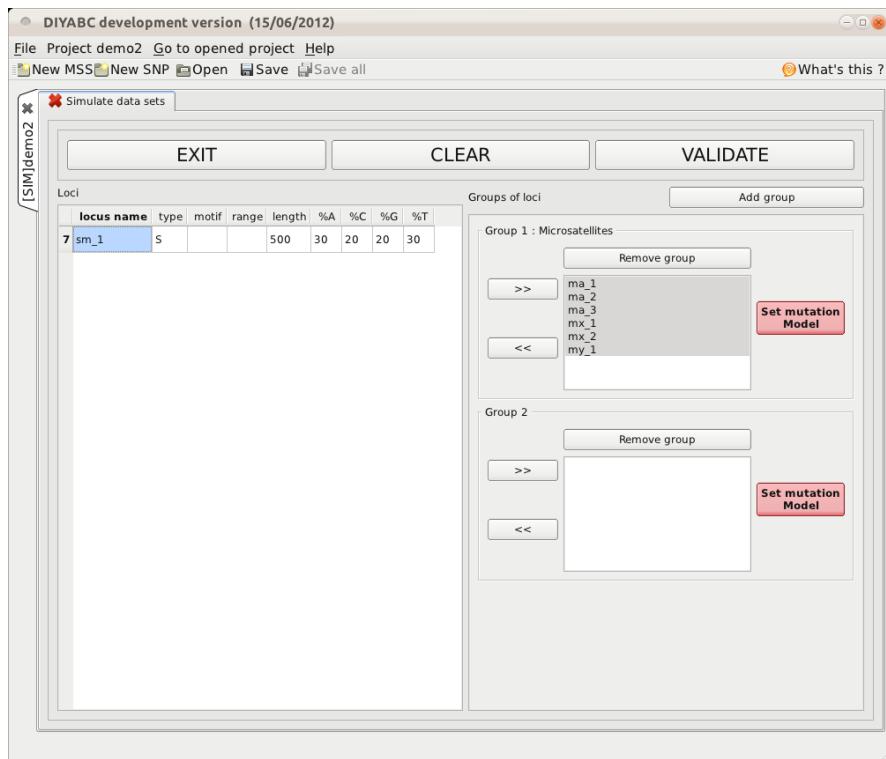
6

7

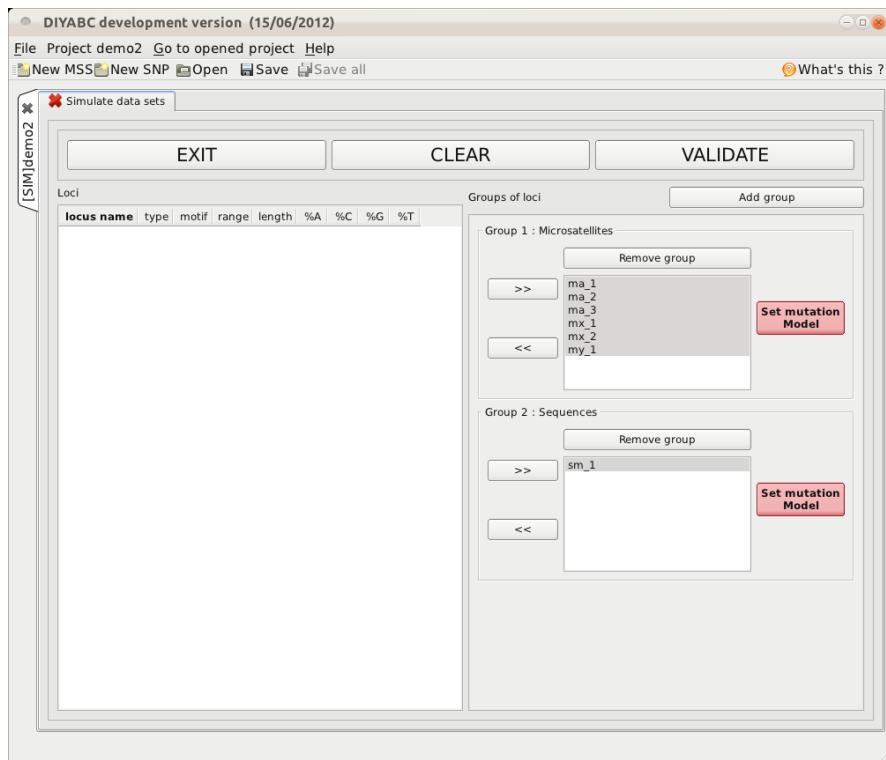
8 The six locus names are transferred into group 1 by clicking on the button :
9

1
2
3
4

Then the DNA sequence locus is selected :

5
6
7
8

and transferred into group 2 in the same way :

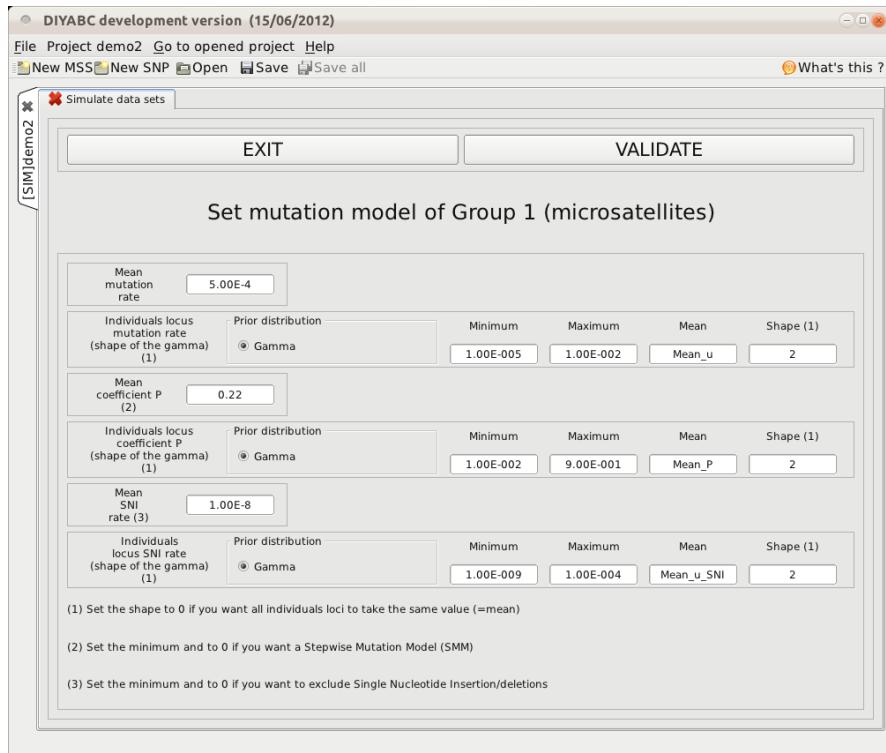


1

2

3 We need now to define the mutation model of each group. Let's click on the **Set Mutation Model**
4 button of group 1:

5



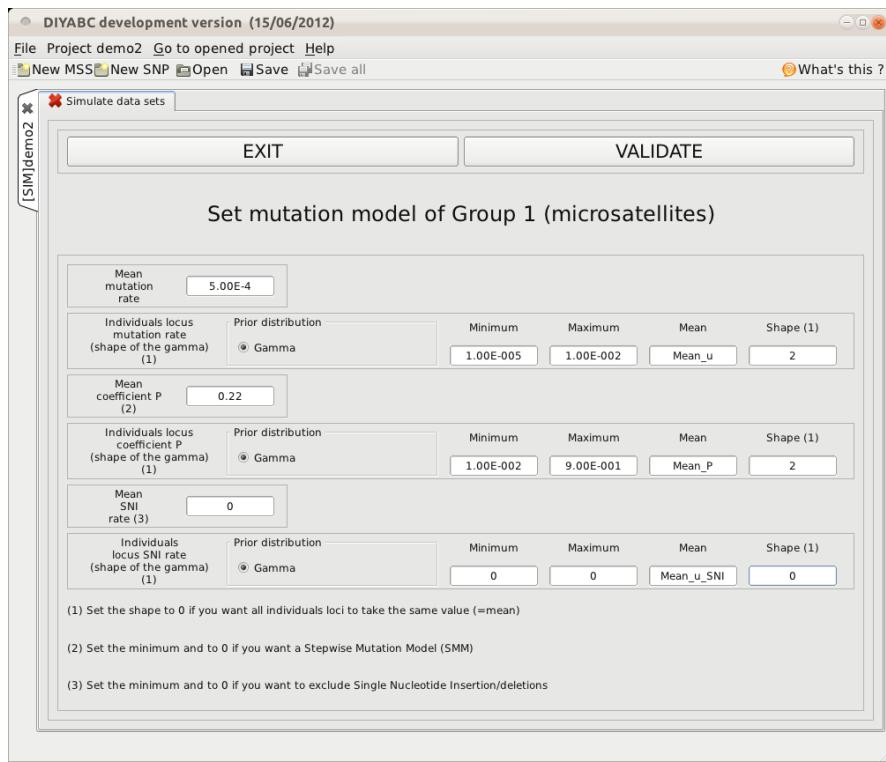
6

7

8 The usual default values appear. We want to exclude single nucleotide insertions/deletions (SNI mutations). So we set to 0 the Mean SNI rate and Minimum, Maximum and Shape of individual loci SNI rates :

9

10

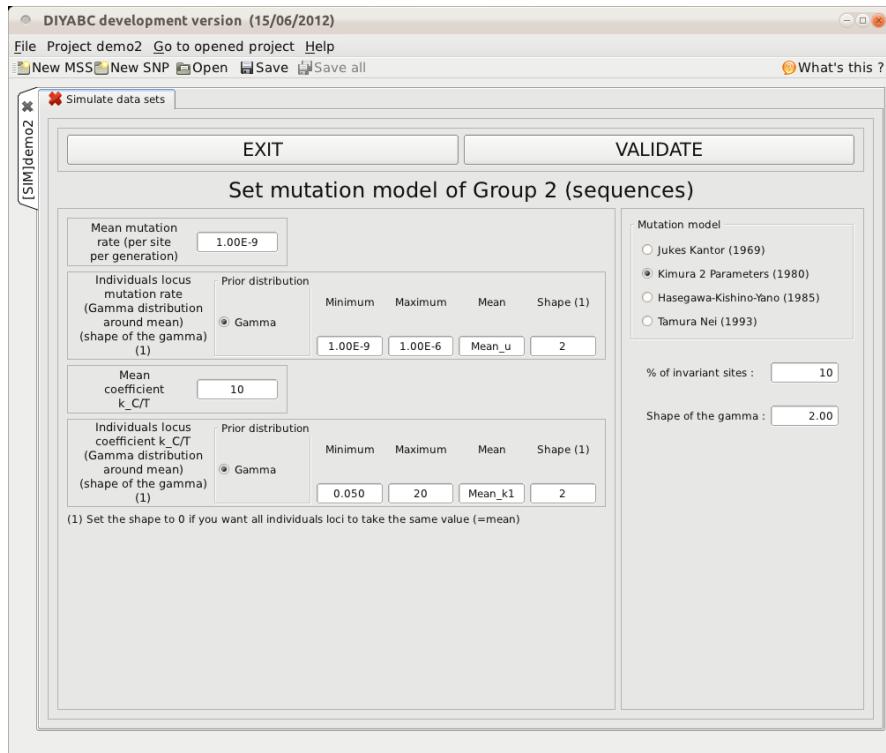


1

2

Once this done, we go back to the previous screen by clicking on the **VALIDATE** button. Then we set the mutation model of the mitochondrial DNA sequence. The default values are as follows :

3



4

5

6

7

8

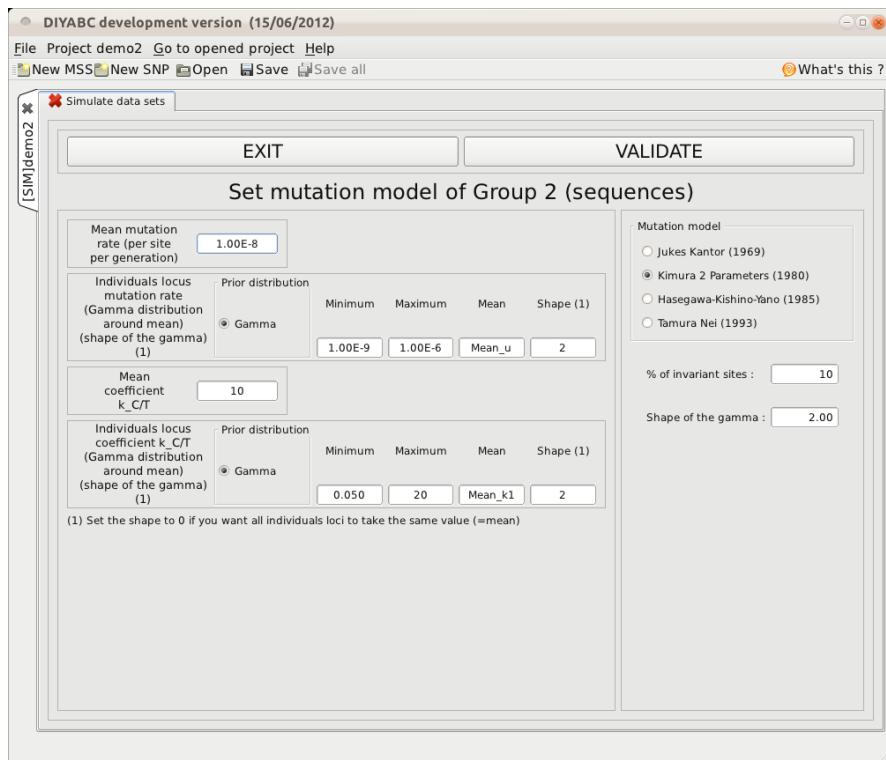
9

10

11

The default mean mutation rate is not suited to mitochondrial DNA which generally evolves at a faster rate than nuclear DNA. So we set its value to 10^{-8} . For all other parameters, we just keep the default values:

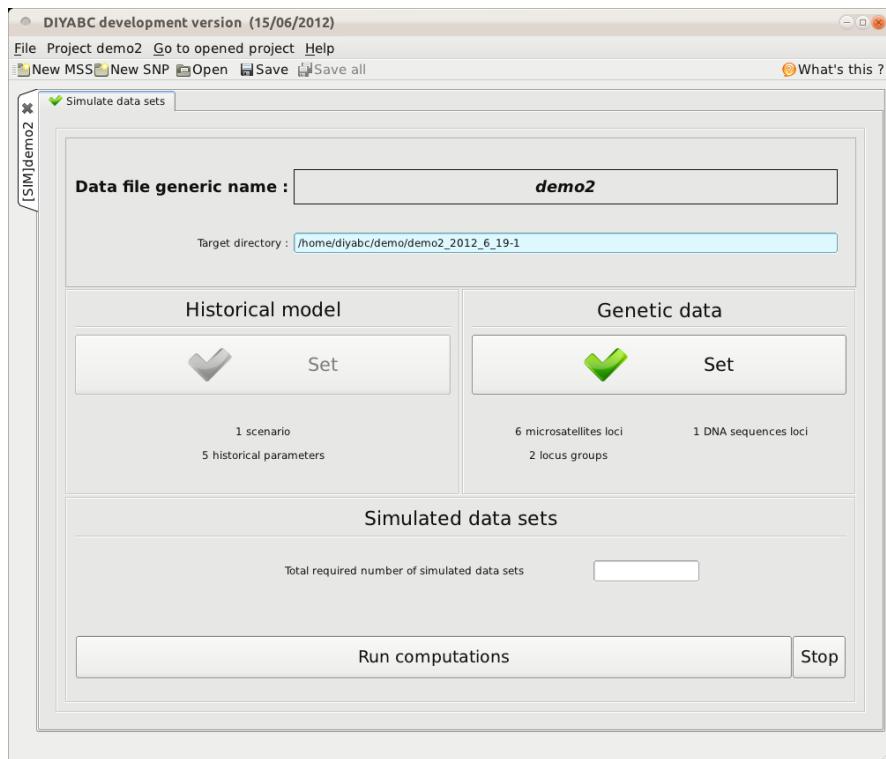
12



1

2

After validating twice, we get back to the main screen :

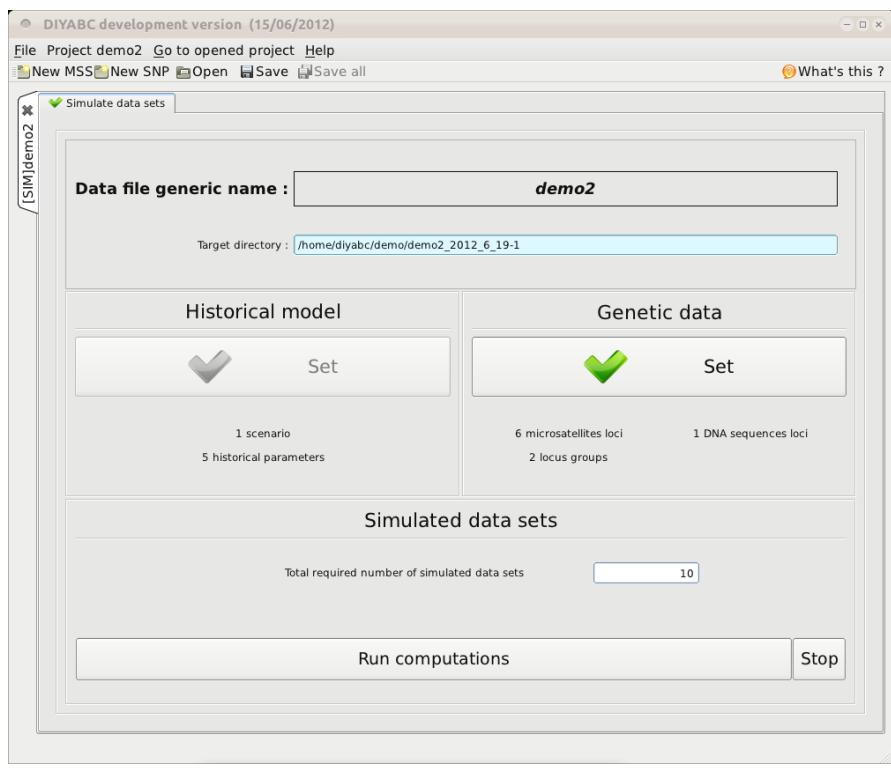


5

6

We require 10 simulated data sets :

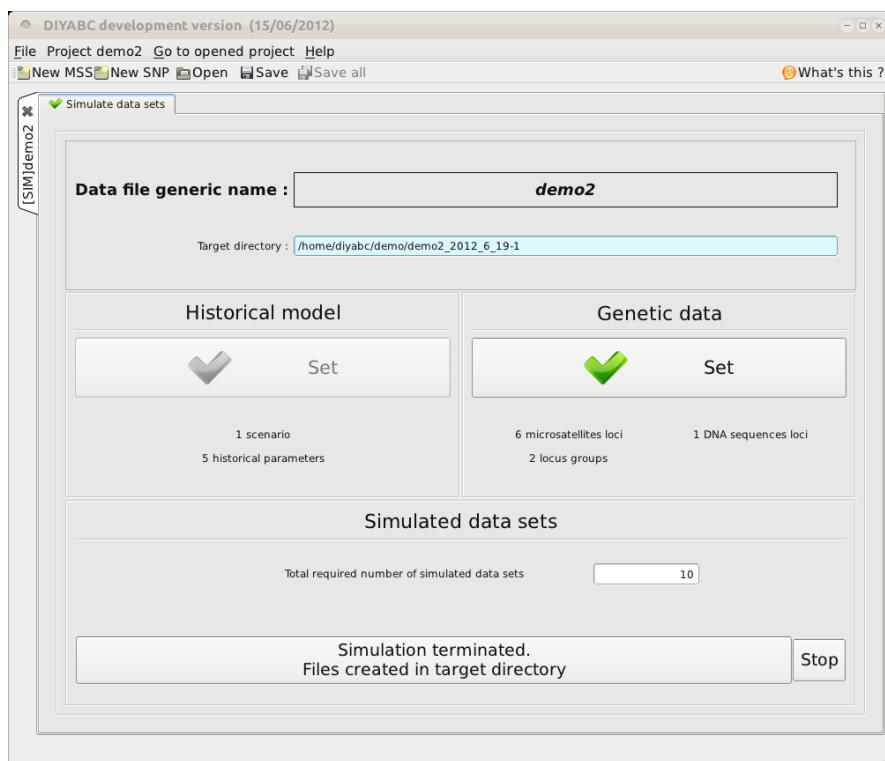
8



1

2

1 We then click on the **Run computation** button. In a matter of seconds, the computation ends up:



2
3 Using the file manager, we can check that ten new files (demo2_001.mss to demo2_010.mss) have been added to new directory :



4
5
6
7
8
9
10
11
12 Opening e.g. the second one with a text editor, we can have a partial view of the simulated genotypes of the first population sample :

file:///home/diyabc/demo/demo2_2012_6_19-1/demo2_002.mss – Kate

Fichier Édition Affichage Signets Sessions Outils Configuration Aide

Nouveau Ouvrir Enregistrer Enregistrer sous Fermer Annuler Refaire Reload All Save All Précédent Suivant

Explorateur de systèmes de fichiers Documents .02_2012_6_19-1 demo2_002.mss

```

1 Simulated genepop file <NM=0.25NF>
2 Locus M_A 1 <A>
3 Locus M_A 2 <A>
4 Locus M_A 3 <A>
5 Locus M_X_4 <X>
6 Locus M_X_5 <X>
7 Locus M_Y_6 <Y>
8 Locus S_M_7 <M>
9 POP
10 1-001 , 205205 207207 197199 207205 203205 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
11 1-002 , 209211 201209 197201 205205 205203 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
12 1-003 , 205207 207209 197199 205207 205203 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
13 1-004 , 209209 207209 197199 205207 199203 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
14 1-005 , 209209 201201 199201 207205 203203 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
15 1-006 , 209209 207203 197197 205205 203201 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
16 1-007 , 209205 201201 197199 205207 199199 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
17 1-008 , 205205 207207 199201 205207 203205 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
18 1-009 , 205209 203203 197197 205205 203203 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
19 1-010 , 209209 201207 197201 205205 203203 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
20 1-011 , 209209 201207 197197 205205 203203 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
21 1-012 , 209209 209201 201199 205205 203199 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
22 1-013 , 209209 209201 197197 205205 203199 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
23 1-014 , 205205 207207 197197 205205 203203 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
24 1-015 , 209209 207207 197197 205205 203203 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
25 1-016 , 209209 207207 203199 205 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
26 1-017 , 209209 201207 197199 205 199 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
27 1-018 , 209209 207207 195197 205 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
28 1-019 , 205209 201201 197197 205 199 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
29 1-020 , 205205 209209 199197 205 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
30 1-021 , 205205 201207 197201 205 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
31 1-022 , 209209 201201 197197 207 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
32 1-023 , 209209 209207 199197 205 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
33 1-024 , 205209 207207 199197 207 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
34 1-025 , 209209 207209 197201 205 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
35 1-026 , 205209 209201 197197 205 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
36 1-027 , 205209 207209 197197 205 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
37 1-028 , 209209 201207 197197 205 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
38 1-029 , 205209 201207 197197 205 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
39 1-030 , 209209 207207 199201 205 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
40 1-031 , 209209 201201 201197 205 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
41 1-032 , 207209 201203 201201 205 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
42 1-033 , 209209 201207 201197 201 203 195 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
43 POP
44 2-001 , 205207 199199 205205 201201 201207 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
45 2-002 , 207211 199195 205201 199201 201199 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
46 2-003 , 205207 201197 205205 201199 207203 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG
47 2-004 , 205207 199199 199199 199203 205207 000 <[CATATGTAAGTTTATCCACTGACTCAGGAACATGATGGAACTATCCATAACGTGTGTG

```

Line: 1 Col: 1 INS LINE UTF-8 demo2_002.mss

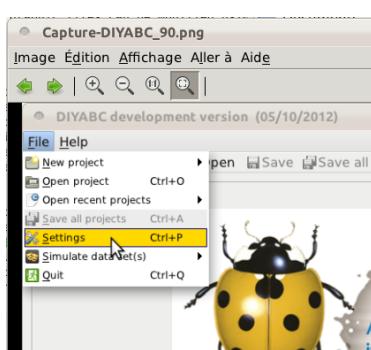
Chercher dans les fichiers Terminal

We can check that the sex ratio is correct : the number of males is one fourth the number of females.

The type of each locus given after the name is also correct. All microsatellite allelic values are odd, in agreement with the motif length (2) and the absence of single nucleotide insertion/deletion. More interestingly, it gives an example of how X- and Y-linked microsatellite loci must be written for each sex (here 15 females and 18 males) in our Genepop format.

3.7 The Settings option of the File menu

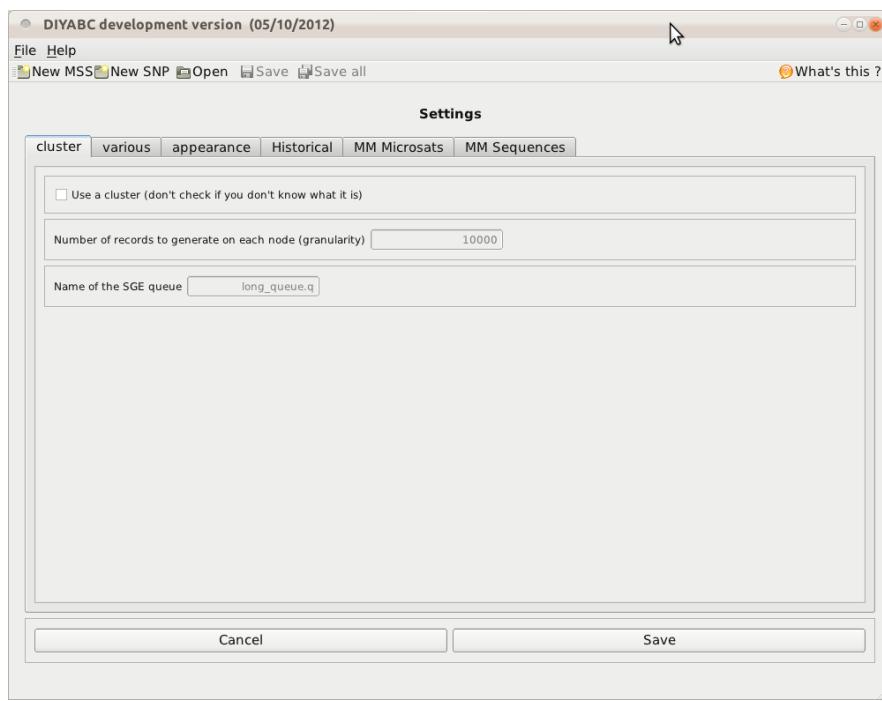
Let us now detail what is under the Settings option of the File menu shown below :



11

12

1 Clicking on the **Settings** option opens up the following multitab window :



3
4

5 3.7.1 Tab “cluster”

6 The first tab (on the left) is related to the use of a computer cluster to perform computations of the
7 reference table. If your computer is connected with a computer cluster on which the necessary files
8 (detailed in section 5) have been correctly installed, and if the computer cluster is run under SGE
9 queuing system, then you can use it for generating the reference table. You then need to :

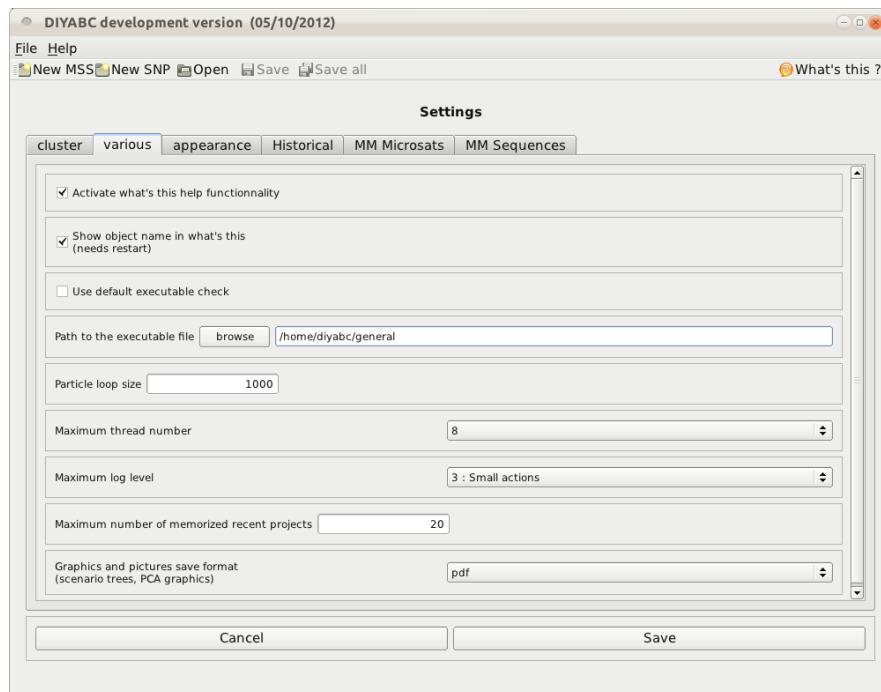
- 10 1. check the box **Use a cluster** (...)
- 11 2. indicate the number of data sets produced by each single job of the queue
- 12 3. give the name of the SGE queue

13 Clicking on the **Run computation** button will send the right script to the computer cluster, monitor
14 the progression of computations and retrieve the reference table once all computations are done.

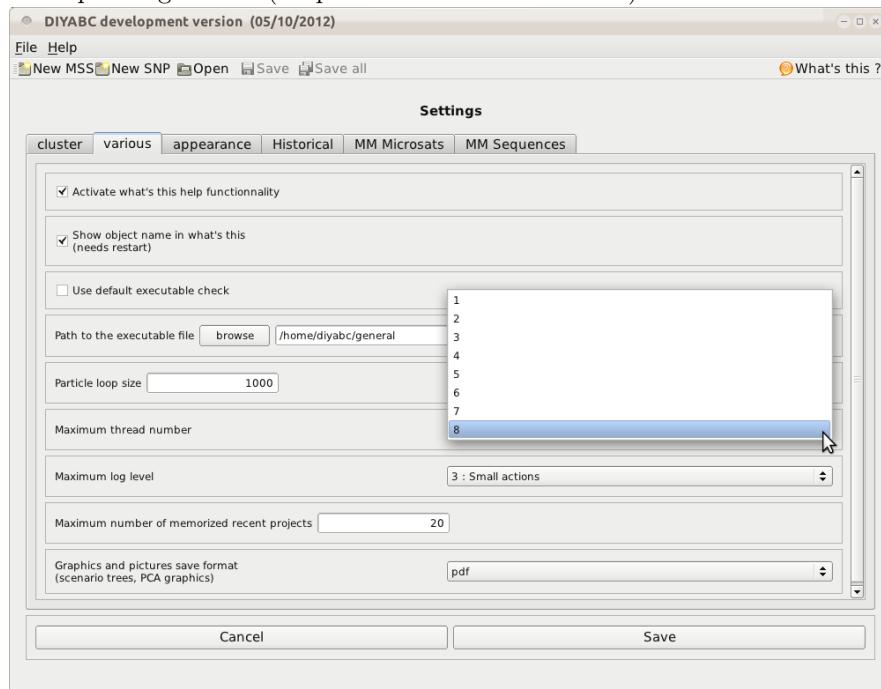
15 3.7.2 Tab “various”

16 The second tab “various” contains the following settings :

- 17 1. **What's this** is a help functionnality that allows the user to obtain a help message when pointing
18 towards a specific feature of the graphic interface such as a button or an edit field. This help
19 functionnality can be activated by checking the correspoding box.
- 20 2. Checking this box is mainly for debugging purpose or signalling a bug.
- 21 3. DIYABC is made of two programs : the graphic interface and a computation program. When the
22 user clicks on buttons such as **Run computation** or **Launch**, the graphic interface programs sends
23 a command that launches the computation program. To issue this command, the graphic interface
24 needs to know where the computation program executable is located. There is a default location
25 which depends on the operating system. Clicking on the box **Use default executable** check will direct
26 the graphic interface to use the executable located in this default directory.
- 27 4. You can also choose another location (e.g. if you want to use a distinct version of the executable)
28 by clicking on the **browse** button.

1
2

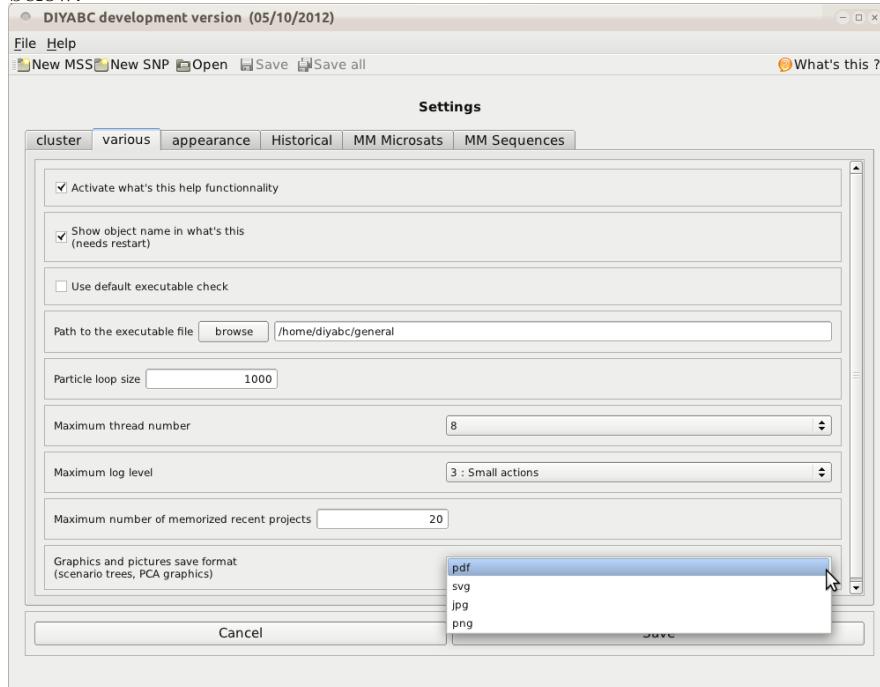
- 3 5. The next setting **Particle loop size** defines the number of data sets (n) that are simulated in a single block when building the reference table. The computation program proceeds as follows : it first simulate and compute summary statistics of n data sets. When this is done, it writes the results to the reference table file. The reason of doing like this is that computation can be multithreaded but not the file writing.
- 4 6. The graphic interface can detect the number of cores of the computer processor. By default, it sets the number of threads of the computation program to this core number. However, if the user wants to keep some cores for other purposes, the number of threads can be reduced by on the corresponding button (drop down menu shown below).

12
13

- 14 7. The next setting is for debugging purpose and/or signalling a bug.

1 8. The graphic interface memorizes recently opened projects. The edit field is used to set the maximum
 2 of memorized recent projects.

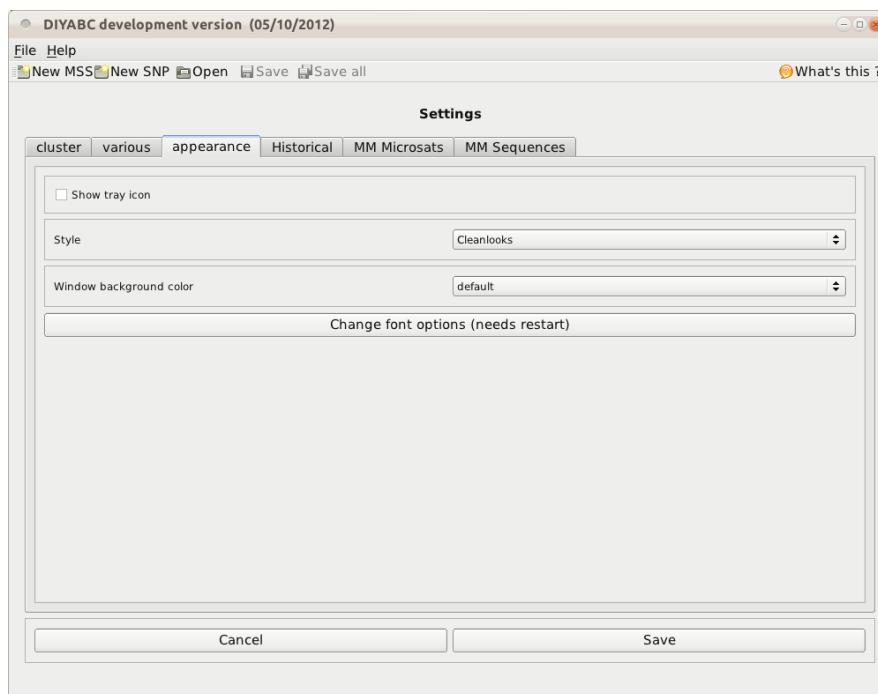
3 9. The last setting concerns the format of graphic files output by different analyses. Choice is shown
 4 below:



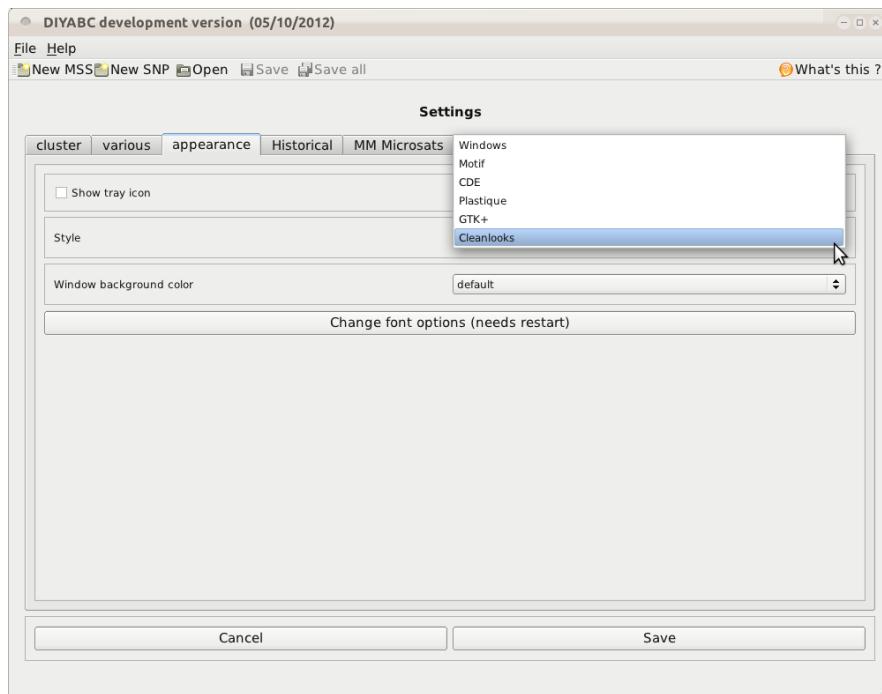
5
 6 Eventually, if changes have been made, they can be either saved or cancelled (two bottom buttons).

8 3.7.3 Tab “appearance”

9 Clicking on this tab results in the following screen :



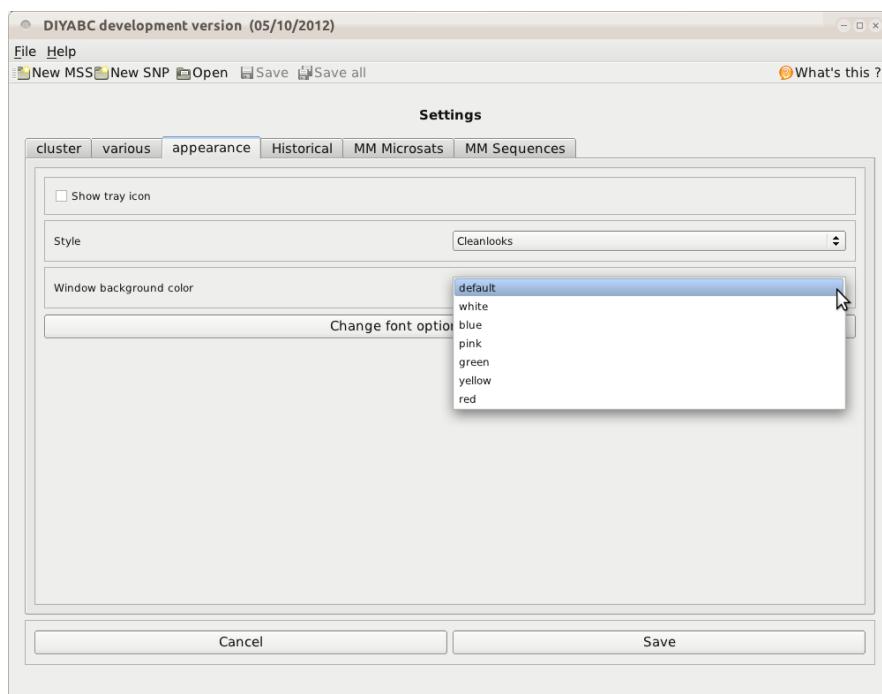
11
 12 The window style can be chosen among the following (click on the upper drop down menu) :



1

2

Likewise, the background color can be chosen among the following colours :

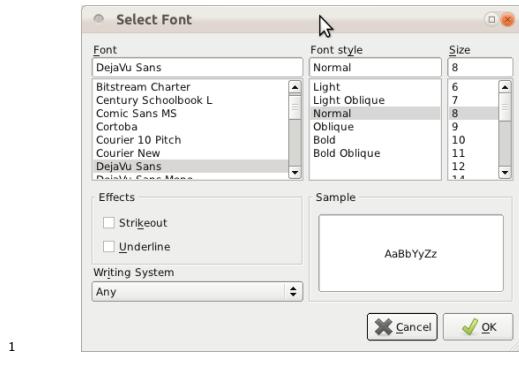


5

6

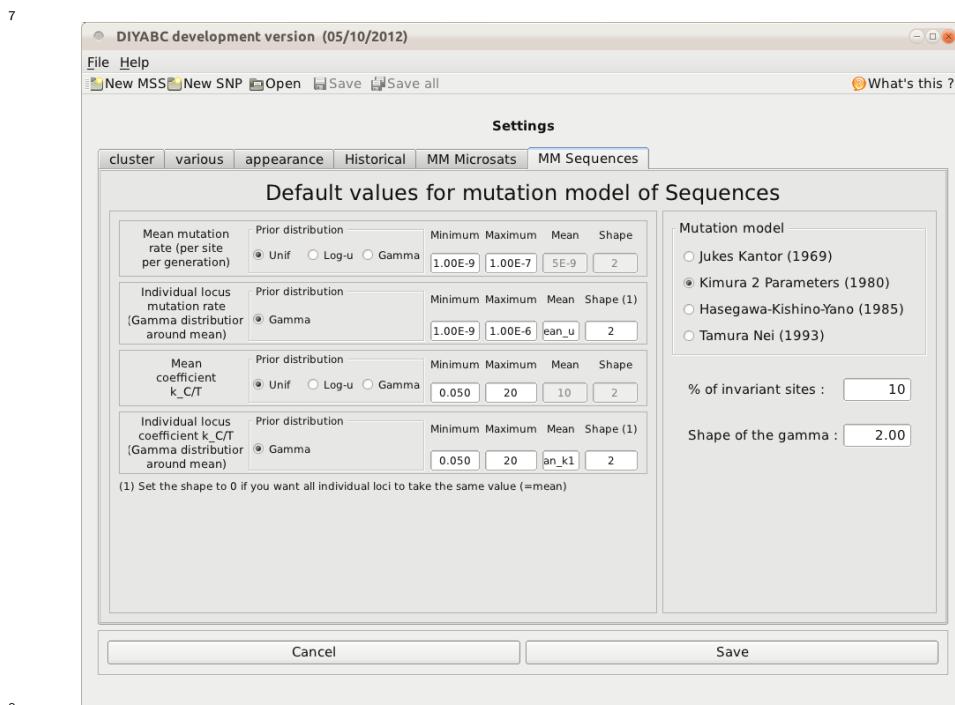
Eventually, one can change the font of texts appearing in the different windows by clicking on the corresponding button. A usual font menu then appears allowing the desired change :

9



3.7.4 Tabs “MM Microsats” and “MM Sequences”

4 These two tabs are used to modify the default values of mutation parameters (MM means Mutation
 5 Model), for microsatellites and DNA sequences respectively. As an example, here is the screen corre-
 6 sponding to the tab “MM Sequences” :



10 The initial default values have been obtained through literature compilation and are valid for a large
 11 number of species. However, some species may have values that differ substantially from most species.
 12 For instance, the mutation rate of some *Drosophila* species are much lower than the values encountered
 13 in many other species (Schug *et al.*, 1997; Vzquez *et al.*, 2000) and is outside the range indicated in the
 14 initial default values.

1 4. Implementation details

2 4.1 Software design

3 DIYABC v2 has been designed in a very different way compared to version 1. Version 1 was a single
 4 executable file were the GUI⁵ and computation codes were highly intricated and both written in the same
 5 language (*Delphi*). In version 2, the GUI and the computation codes have been completely separated.
 6 Actually, the GUI is a script written in *python* and all computations are included in a program written
 7 in *C++*. In opposition to *Delphi* which is restricted to a single OS (*Windows*), *python* and *C++* can
 8 be used with the main three OS (*Linux*, *Mac* and *Windows*), allowing version 2 to be operated under all
 9 three OS.

10 The GUI uses the *Qt* graphic library. The computation code is linked to the *openmp* library allowing a
 11 better use of multicore/multiprocessor computers.

12 The GUI can launch the computation program with the right parameters and keeps track of the progress
 13 of the latter through small log files. The GUI can launch as many computation programs as there are
 14 open projects, but no more than one computation program per project. A *lock* file located in the project
 15 directory is created when the computation program is launched by the GUI and removed when the
 16 computation program has normally terminated. When the computation program has exited anomalously,
 17 the GUI issues an error message trying to explain where the programm failed.

18 4.2 Files

19 The program uses and produces various files which we will describe now.

20 4.2.1 data files

21 Data files are text files that contain information about the samples : number and names of microsatellite
 22 markers, multilocus genotypes of individuals. The basic format is that of the Genepop software (Ray-
 23 mond and Rousset, 1995) and data files produced by DIYABC are under this format. **Microsatellite**
24 genotypes must be noted with 3 (haploid) or 6 (diploid) digits, these three digit numbers
25 being the length in nucleotides of the corresponding PCR products. In addition, we have added
 26 some features to this basic format in order to use sequence data. All these additions are explained in
 27 section 4.4.

28 Any extension is accepted for datafile names, including no extension at all. If the data file is simulated
 29 with DIYABC, the extension is **mss** for microsatellite/DNA sequence data and **snp** for SNP data. The
 30 next page shows examples of data sets saved.

31 4.2.2 reference table files

32 Reference table files are binary files which include two successive parts :

- 33 • The first part is a header which contains information necessary to read the second part, such as the
 34 number of scenarios, or the number of parameters of each scenario.
- 35 • The second part contains simulated data set records, each record containing the scenario number,
 36 the parameter and summary statistics values.

37 Each time a reference table is created or increased (each time the **Run computation** button is pressed),
 38 a text file is created in the project directory with the name **first_records_of_the_reference_table_X.txt**
 39 in which X is an integer number starting at 0 and increasing each time the **Run computation** button is
 40 pressed. This file provides a text version of the first n newly created records of the reference table (n
 41 being equal to the *Particle loop size*, see section 3.7.3).

43 4.2.3 output files

44 As already seen, DIYABC achieves different analyses : comparison of scenarios, estimation of posterior
 45 distribution of parameters, model checking, computation of bias and mean square errors and evaluation of
 46 confidence in scenario choice. Each analysis has its own output which can be printed and saved. Graphs

⁵Graphic User Interface

1 are saved under the chosen format and non-graphic output are saved in text files.

2
3 We now describe all the files produced by each type of analysis. These files are located in directories
4 (one directory per analysis) gathered in the **analysis** subdirectory of the project directory. Below is an
5 example of the TOYTEST2_2012_9_26-1 project directory substructure:

```
6   └── TOYTEST2_2012_9_26-1
7     └── analysis
8       ├── bias1-1_bias
9       ├── bias1-2_bias
10      ├── bias1-3_bias
11      ├── bias1-4_bias
12      ├── bias1-5_bias
13      ├── bias1_bias
14      ├── compscen_comparison
15      ├── conf1_confidence
16      ├── conf2_confidence
17      ├── estim_s2_estimation
18      ├── mc_scen2_modelChecking
19      ├── mc_scen2_newmsstat-1_modelChecking
20      ├── mc_scen2_newmstat-1_modelChecking
21      ├── mc_scen2_newmstat-2_modelChecking
22      ├── mc_scen2_newmstat-3_modelChecking
23      ├── mc_scen2_newmstat-4_modelChecking
24      ├── mc_scen2_newmstat-5_modelChecking
25      ├── mc_scen2_newmstat_modelChecking
26      ├── mc_scen2_newsstat_modelChecking
27      ├── new4_modcheck_stats_-1_modelChecking
28      ├── new4_modcheck_stats_-2_modelChecking
29      ├── new4_modcheck_stats_-3_modelChecking
30      ├── new4_modcheck_stats_-et+_modelChecking
31      └── NEW_NEW_TEST_MODCHECK_modelChecking
32
33   └── preval_pca
34     └── pictures
```

7
8 Note that each directory name starts with the name of analysis followed by the type of analysis, *e.g.*
9 **bias** for a bias/precision analysis or **comparison** for a comparison of scenarios. In addition, when a
10 picture has been saved, the corresponding file is located under a subdirectory named **pictures** (*e.g.* at
11 the bottom of the figure above).

12 **Pre-evaluate scenario prior combinations :** This analysis can produce two output files named **ACP.txt**
13 and **locate.txt**. The former is the output of the Principal Component Analysis and the latter
14 that of the analysis giving the proportion of simulated data sets which have a value below the
15 observed value for every summary statistics. This latter file is exactly what appears in the GUI.
16 The structure of the **ACP.txt** file is the following. The first line indicates the number of points
17 of the PCA, the number of PCA components (axes) and the inertia of each component, all values
18 are separated by a single space. The second line provides the components of the observed data. It
19 starts with a zero which corresponds to the scenario number in the following lines. Each subsequent
20 line provides the components of data simulated according to a given scenario which number is at
21 the beginning of the line. If one or more PCA figures have been saved, the corresponding files are
22 saved in the **pictures** subdirectory. They are named as **refTable_PCA_X.Y.N.pdf**, with X and Y
23 giving the axis numbers and N being the number of represented points.

24
25 **Compute posterior probabilities of scenarios :** This analysis produces three output text files : **compdirect.txt**,
26 **complogreg.txt** and **commdirlog.txt**. The latter is directly visualized in the GUI when clicking
27 the **view numerical results** button. The first two files are used by the GUI to elaborate the two
28 graphics (Direct approach and Logistic regression). Again, if graphics have been saved, the corre-
29 sponding file(s) is(are) in the **pictures** subdirectory of the analysis directory.

30 **Evaluate confidence in scenario choice :** This analysis produces a single output file, **confidence.txt**,
31 the content of which is visualized in the GUI.

32 **Estimate posterior distributions of parameter :** Nine files are written as output of this type of
33 analysis :

- three files `mmmq_original.txt`, `mmmq_composite.txt` and `mmmq_scaled.txt` contain the statistics (mean, median, mode and quantiles) for the original, composite and scaled parameters, respectively. They are visualized in the GUI when clicking the `view numerical results` button.
- three files `paramstatdens_original.txt`, `paramstatdens_composite.txt` and `paramstatdens_scaled.txt` are used by the GUI to produce the graphics showing prior/posterior distribution.
- three files `phistar_original.txt`, `phistar_composite.txt` and `phistar_scaled.txt` contains the ϕ^* values of the original, composite and scaled parameters, respectively.

As already mentionned, saved graphics are located in a `pictures` subdirectory.

Compute bias and precision of parameter estimations : Three files `bias_original.txt`, `bias_composite.txt` and `bias_scaled.txt` are produced by this type of analysis. All three files are visualized in the GUI.

Perform model-checking The output files of this type of analysis are the same as those of the *Pre-evaluate scenario prior combinations* analysis (see above). The only difference is in the names of the two text files which start with `mc` for `model checking`.

In addition, the GUI program writes several files in the project directory :

command.txt : this text file contains the history of commands issued by the GUI to be achieved by the computation program.

conf.analysis : this text file contains information about analyses.

conf.gen.tmp : this text file contains information about the loci, the genetic parameters and the summary statistics.

conf.hist.tmp : this text file contains information about the scenario and the historical parameters.

conf.th.tmp : This text file contains the title line of the reference table.

conf.tmp : This text file contains the name of the dataset and the number of parameters and summary statistics.

header.txt : This text file is a concatenation of the previous four files and is red by the computation program.

xxx.diyabcproject : This text file contains the path to the `xxx` project.

RNG_state_0000.bin : This binary file contains the current state of the random generator.

init_rng.out : This text file contains information about the initialization of the random generator.

The computation program writes the following files in the project directory :

reftable.log : This text file is produced when a reftable is increased. It provides the GUI with information about the progress of computations : achieved number of records, time left.

statobs.txt : This text file is written every time an analysis is performed. It contains the values of summary statistics for the observed data set.

The following files are output by the computation program everytime it has been launched by a specific command of the GUI (their use is only for debugging purposes and they are all in the project directory) :

general.out : when computing a reftable.

pre-ev.out : when performing a *Pre-evaluate scenario prior combinations* analysis.

compare.out : when performing a *Compare scenarios* analysis.

- confidence.out** : when performing a *Confidence in scenario choice* analysis.
 - estimate.out** : when performing a *ABC parameter estimation* analysis.
 - bias.out** : when performing a *bias-precision* analysis.
 - modelChecking.out** : when performing a *model checking* analysis.
- When performing a *Bias-precision* or a *Confidence in scenario choice* analysis, the computation program simulates what we call *pseudo-observed datasets*. The parameter and summary statistics values of these pseudo-observed datasets are written in a text file named **pseudo-observed_datasets_xxx.txt** in which **xxx** is the name given to the analysis.

4.3 Missing data

Missing or undetermined genotypes should be coded as 000 (haploid microsatellites, 000000 (diploid microsatellites), < [] > (haploid sequences) or < [][] > (diploid sequences) and 9 (SNP) in the data file. Missing data are taken into account in the following way. For each appearance of a missing genotype in the observed data set, the programs records the individual and the locus. When simulating data sets, the program replaces the simulated genotype (obtained through the coalescence process algorithm) by the missing data code at all corresponding locations. All summary statistics are thus computed with the same missing data as for the observed data set.

4.4 Data files

There are two different incompatible formats for data files, one for SNP loci and the other for microsatellite/DNA sequence data.

For the latter, the format already presented in version 1 of DIYABC is an extended Genepop format. The additional features are :

1. In the title line appears the sex ratio noted between < and > under the form < NM = rNF >, in which r is the ratio of the number of females per male (e.g. < NM = 2.5NF > means that the number of males is 2.5 times the number of females). Since the title is generally only copied, this addition should not interfere with other programs using Genepop datafiles. Also if there is no such sex ratio addition, DIYABC will consider by default that NM=NF.
2. After the locus name, there is an indication for the category of the locus which is < A > for autosomal diploid loci, < H > for autosomal haploid loci, < X > for X-linked (or haplo-diploid) loci, < Y > for Y-linked loci and < M > for mitochondrial loci. If no category is noted, DIYABC will consider the locus as autosomal diploid or autosomal haploid depending on the corresponding genotype of the first typed individual.
3. Genotypes of microsatellite loci are noted with six digit numbers (e.g. 190188) if diploid and by three digit numbers (e.g. 190) if haploid.
4. Sequence locus are noted between < and > . In addition each sequence allele is noted between brackets. For instance, a haploid sequence locus will be noted < [GTCTA] > and a diploid sequence locus < [GTCTA][GTCTT] >.
5. Missing microsatellite genotypes are noted 000 if haploid or 000000 if diploid.
6. Missing sequence genotypes are noted < [] > if haploid or < [][] > if diploid.

For SNP data, the format includes:

- a first line providing the sex-ratio as above
- a second line starting with the three keywords IND SEX POP, separated by at least one space, followed by as many letters as SNP loci, the letter giving the location of the locus as above (< A > for autosomal diploid loci, < H > for autosomal haploid loci, < X > for X-linked (or haplo-diploid) loci, < Y > for Y-linked loci and < M > for mitochondrial loci). Letters are separated by a single space.

- 1 • as many lines as there are genotyped individuals, with the code-name of the individual, a letter (*M*
 2 or *F*) indicating its sex, a code-name for its population and the values (0, 1 or 2) of the number of
 3 the reference allele at each SNP locus.

4 Below are three examples of data sets simulated by DIYABC.

- 5 In the first example, this data set includes two population samples, each of 12 diploid individuals (8
 6 females and 4 males in the first sample and 5 females and 7 males in the second sample). As deduced
 7 from the letter between < and > on the locus name lines (see page 25), these individuals have been
 8 genotyped at 3 microsatellite loci (1 autosomal <*A*>, 1 X-linked <*X*> and 1 Y-linked <*Y*>) and 3
 9 DNA sequence loci (1 autosomal, 1 X-linked and 1 mitochondrial <*M*>). The species sex-ratio, given
 10 in the title line, is of three males for one female (<*NM* = 3*NF*>) or in other words, the number of
 males equals three times the number of females.

```
18/03/2010 17:02:03 N=1000000 <NM=3.00NF>
locus M_A_1 <A>
locus M_X_1 <X>
locus M_Y_1 <Y>
locus S_A_1 <A>
locus S_X_1 <X>
locus S_M_1 <M>
POP
 1-001 , 172176 156172 000 <[TAGCTA][TAGCTA]> <[GATCG][GATCG]> <[ATTACTTA]>
 1-002 , 180204 156184 000 <[TAGCTA][TAGCTA]> <[GATCG][GATCG]> <[ATTACTTG]>
 1-003 , 150162 164168 000 <[TAGCTA][TAGCTA]> <[GATCG][GATCG]> <[ATTACTTG]>
 1-004 , 210218 154168 000 <[TAGCTA][TAGCTA]> <[GATCG][GATCG]> <[ATTACTTA]>
 1-005 , 188218 160152 000 <[TAGCTA][TAGCTA]> <[GATCG][GATCG]> <[ATTACTTG]>
 1-006 , 180190 172216 000 <[TAGCTA][TAGCTA]> <[GATCG][GATCG]> <[ATTACTTG]>
 1-007 , 174168 180210 000 <[TAGCTA][TAGCTA]> <[GATCG][GATCG]> <[ATTACTTA]>
 1-008 , 168208 160202 000 <[TAGCTA][TAGCTA]> <[GATCG][GATCG]> <[ATTACTTG]>
 1-009 , 176162 180 216 <[TAGCTA][TAGCTA]> <[GATCG]> <[ATTACTTA]>
 1-010 , 150194 220 196 <[TAGCTA][TAGCTA]> <[GATCG]> <[ATTACTTA]>
 1-011 , 158176 182 226 <[TAGCTA][TAGCTA]> <[GATCG]> <[ATTACTTG]>
 1-012 , 154156 166 218 <[TAGCTA][TAGCTA]> <[GATCG]> <[ATTACTTA]>
POP
 2-001 , 168164 184216 000 <[TAGCTA][TAGCTA]> <[GATCG][GATCG]> <[ATTACTTG]>
 2-002 , 164160 152208 000 <[TAGCTA][TAGCTA]> <[GATCG][GATCG]> <[ATTACTTG]>
 2-003 , 166222 180170 000 <[TAGCTA][TAGCTA]> <[GATCG][GATCG]> <[ATTACTTG]>
 2-004 , 212150 228166 000 <[TAGCTA][TAGCTA]> <[GATCG][GATCG]> <[ATTACTTG]>
 2-005 , 210152 228196 000 <[TAGCTA][TAGCTA]> <[GATCG][GATCG]> <[ATTACTTA]>
 2-006 , 156212 178 192 <[TAGCTA][TAGCTA]> <[GATCG]> <[ATTACTTA]>
 2-007 , 200226 174 202 <[TAGCTA][TAGCTA]> <[GATCG]> <[ATTACTTG]>
 2-008 , 196200 168 178 <[TAGCTA][TAGCTA]> <[GATCG]> <[ATTACTTG]>
 2-009 , 174174 172 192 <[TAGCTA][TAGCTA]> <[GATCG]> <[ATTACTTG]>
 2-010 , 178194 212 182 <[TAGCTA][TAGCTA]> <[GATCG]> <[ATTACTTA]>
 2-011 , 204160 180 204 <[TAGCTA][TAGCTA]> <[GATCG]> <[ATTACTTA]>
 2-012 , 190226 160 226 <[TAGCTA][TAGCTA]> <[GATCG]> <[ATTACTTG]>
```

- 11 In the second example, the species is haploid. Individuals have been genotyped at three autosomal
 12 microsatellite loci and one mitochondrial DNA sequence locus. The species being haploid (deduced from
 13 the presence of autosomal haploid loci), no indication of the sex-ratio appears in the title line.

```
18/03/2010 17:44:54 N=100000
locus M_H_1 <H>
locus M_H_2 <H>
locus M_H_3 <H>
locus S_M_1 <M>
POP
 1-001 , 164 184 210 <[TCCTTCCGGTTGTGCGACCACTTCGTACGTT]>
 1-002 , 164 184 214 <[TCCTTCCGGTTGTGCGACCACTTCGTACGTT]>
 1-003 , 150 186 214 <[TCCTTCCGGTTGTGCGACCACTTCGTACGTT]>
 1-004 , 160 188 220 <[TCCTTCCGGTTGTGCGACCACTTCGTACGTT]>
 1-005 , 152 176 214 <[TCCTTCCGGTTGTGCGACCACTTCGTACGTT]>
POP
 2-001 , 154 188 210 <[TCCTTCCGGTTGTGCGACCACTTCGTACGTT]>
 2-002 , 154 196 206 <[TCCTTCCGGTTGTGCGACCACTTCGTACGTT]>
 2-003 , 166 194 202 <[TCCTTCCGGTTGTGCGACCACTTCGTACGTT]>
 2-004 , 150 194 200 <[TCCTTCCGGTTGTGCGACCACTTCGTACGTT]>
POP
 3-001 , 202 222 202 <[TCCTTCCGGTTGTGCGACCACTTCGTACGTT]>
 3-002 , 226 206 198 <[TCCTTCCGGTTGTGCGACCACTTCGTACGTT]>
 3-003 , 216 206 208 <[TCCTTCCGGTTGTGCGACCACTTCGTACGTT]>
```

In the third example, the species is diploid and has been genotyped at a large number of SNP autosomal loci. The first line provides the species sex-ratio. The second line indicates what's in the different columns : individual name in column 1, individual sex in column 2, population name in column 3 and one column per SNP locus. Columns are separated by one or more spaces. SNP are coded 0, 1 or 2 according to the number of reference alleles. Only the top left part of the data file is represented below :

<NM=0.428571NF>																													
IND	SEX	POP	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A	A						
P1_1	F	P1	0	0	2	0	0	1	0	0	0	1	0	0	0	0	0	2	1	1	0	1	0	2					
P1_2	F	P1	0	0	2	0	2	0	0	0	0	1	0	0	0	2	0	0	0	0	0	2	0	1	0	2			
P1_3	F	P1	0	0	2	0	2	0	0	0	0	2	0	0	0	1	0	0	0	0	1	0	0	1	0	0	2		
P1_4	F	P1	0	0	2	0	0	1	0	0	0	2	0	0	0	1	0	0	0	1	0	0	0	1	1	0	1	2	
P1_5	F	P1	0	0	2	0	1	0	0	0	0	1	1	0	0	1	0	1	0	0	0	2	0	1	1	0	1	2	
P1_6	F	P1	0	0	2	1	0	1	0	0	0	1	0	0	0	1	0	0	0	0	2	0	0	0	0	0	1	1	
P1_7	F	P1	0	0	2	0	0	0	0	0	1	1	0	0	0	1	0	0	1	0	0	0	0	1	0	0	0	2	
P1_8	F	P1	0	0	1	0	0	0	0	0	0	1	0	0	0	2	0	0	0	0	1	0	0	0	0	0	0	2	
P1_9	F	P1	0	0	1	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	1	0	1	1	2	0	2	
P1_10	F	P1	0	0	0	0	0	0	0	1	1	2	1	0	0	2	1	0	0	0	0	0	0	1	1	0	0	2	
P1_11	F	P1	0	0	2	0	0	2	0	0	0	1	0	0	0	2	0	0	0	0	2	0	0	1	0	0	2		
P1_12	F	P1	0	0	1	0	1	1	0	0	1	0	0	0	0	1	0	1	0	0	0	1	0	0	2	1	0	2	
P1_13	F	P1	0	0	1	0	0	2	0	0	0	1	0	0	0	2	0	0	0	0	1	0	0	0	1	0	0	2	
P1_14	F	P1	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	2	0	1	0	1	2	
P1_15	F	P1	0	0	1	0	0	2	0	0	0	2	0	0	0	2	0	0	0	0	2	0	0	0	0	1	0	2	
P1_16	F	P1	0	0	2	0	0	2	0	0	0	1	0	0	0	1	0	0	0	0	1	0	1	0	0	0	2		
P1_17	F	P1	0	0	2	0	1	2	0	0	0	1	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	2	
P1_18	F	P1	0	0	1	0	2	2	0	0	0	1	0	0	0	1	0	0	0	0	1	0	1	1	1	0	0	2	
P1_19	F	P1	0	0	1	0	0	2	0	0	0	1	0	0	0	2	0	0	0	0	2	0	0	0	0	0	0	2	
P1_20	F	P1	0	0	2	0	0	1	0	0	0	1	0	0	0	1	0	0	0	0	1	0	1	0	1	0	0	2	
P1_21	F	P1	0	0	0	0	1	0	0	0	1	0	0	0	1	0	1	0	0	0	2	0	1	0	0	0	1	2	
P1_22	F	P1	0	0	1	0	0	1	0	0	1	0	0	0	0	0	0	0	0	0	2	0	2	1	1	0	0	2	
P1_23	F	P1	0	0	1	0	0	1	0	0	0	2	0	0	0	1	0	0	0	0	2	0	1	2	0	0	2		
P1_24	F	P1	0	0	2	0	0	0	0	0	2	0	0	0	2	0	0	0	0	0	2	0	2	1	2	0	0	2	
P1_25	F	P1	0	0	2	0	0	2	0	0	0	1	0	0	0	1	0	0	0	0	2	0	1	2	0	0	1	2	
P1_26	M	P1	0	0	1	0	1	1	0	1	1	0	0	0	2	0	0	0	0	1	1	0	1	1	1	0	0	2	
P1_27	M	P1	0	0	1	0	0	1	0	0	1	1	0	0	0	1	0	0	0	0	1	0	1	1	2	0	0	1	
P1_28	M	P1	0	0	2	0	0	0	0	0	0	2	0	0	0	0	1	0	0	0	0	1	0	0	0	0	0	2	
P1_29	M	P1	0	0	1	0	0	0	0	0	0	1	0	0	0	2	0	0	0	0	0	1	0	0	0	0	0	0	2
P1_30	M	P1	0	0	1	0	1	2	0	0	0	2	1	0	0	0	0	0	0	0	0	1	0	1	0	0	0	1	2
P1_31	M	P1	0	0	1	0	0	1	1	0	0	0	0	0	1	0	0	0	0	0	2	0	1	1	1	0	1	2	

1 5. Cluster version

2 The process of simulating data sets is generally a time consuming part of the ABC approach. Typically,
 3 one to several millions data sets are needed to build up a reference table and this process can last several
 4 hours to several days. For those that can have access to a computer grid cluster, two additional programs
 5 are available. They both run under Linux. Source files are written in Delphi pascal and can be compiled
 6 with Free Pascal Compiler (<http://www.freepascal.org/download.var>) using the command line :

7 >fpc -Sd -m Delphi <filename>

8 diyabc_sim :

9 This program simulates a given number of data sets according to the information given in a
 10 **reftableHeader** file. The simulated data sets are output in a **reftable** file. This program re-
 11 quires 3 parameters : the name of the reftableHeader file, a string that will be affixed to the name
 12 of the reftable file (to distinguish it from other reftable files) and the number of desired simulated
 13 data sets. This programs has two input files : the reftableHeader file and the data file. The data file
 14 must match the one written at the beginning of the reftableHeader file. The data file is just needed
 15 to evaluate sample sizes and number of loci as well as determining missing data. As an example,
 16 supposing a reftableHeader file named **mydataanalysis.reftableHeader**, a command such as :

17 >diyabc_sim mydataanalysis 01 5000

18 will produce a reference table file named **mydataanalysis01.reftable** containing 5,000 simulated
 19 data sets according to the observed data set and analysis summarized in **mydataanalysis.reftableHeader**.

20 diyabc_cat :

21 This program pools all reftable files of a directory into a single reftable file. This programs requires
 22 2 parameters : the name of the reftableHeader file used to generate all the reftable files to be pooled
 23 and the name of the output reftable. **diyabc_cat** reads the reftableHeader file (which has to be
 24 in the same directory), scans all available reftable files present in the directory and select all those
 25 with a header that matches the reftableHeader file. The selected files are then copied to a single
 26 reftable file. Following with the example above, a command such as :

27 >diyabc_cat mydataanalysis mda

28 will pool all reftable files with header matching **mydataanalysis.reftableHeader** in an output file
 29 named **mda.reftable**.

30 The cluster is just used to produce a reftable file ready to be analyzed by DIYABC. All treatments
 31 are thus performed (under Windows) with the DIYABC programs. A typical session will include :

32

- 33 1. Using DIYABC, input data file name, provide scenario(s), provide priors for historico-demographic
 34 and mutationnal parameters, provide motif lengths and ranges, select summary statistics, simulate
 35 a few datasets (in one wants to check that the ouput is OK, but this last step is not compulsory).

36 2. Transfer the reftableHeader file and the data file on the cluster and using **diyabc_sim** and **diyabc_cat**
 37 as explained above, create a large reftable file.

38 3. Transfer the reftable file back to the Windows directory where the analysis began and rerun DIYABC
 39 to perform comparison of scenarios, estimation of parameters or computation of bias/precision as
 40 needed.

41 For those who have access to a cluster with Sun's Grid Engine software, they can use the following script
 42 to create multiple reference tables (see next page). However, this script will not exempt the user to
 43 manually run **diyabc_cat** to concatenate all reference tables in a single one. The script is run with the
 44 following four parameters :

45 >ScriptABC file.dat file.reftableHeader n_jobs n_iter_per_job

46

Figure 1: Bash script that can be run with Sun Grid Engine software

```

#!/bin/bash
#-----
#####TESTS#####
#-----
if [ "$#" -ne "4" ]
then
    echo "Wrong number of arguments"
    echo "./ScriptABC file.dat file.reftableHeader n_jobs"
    exit
fi
if [ ! -e diyabc_sim ]
then
    echo "Executable file diyabc_sim missing"
    exit
fi
if [ ! -e $1 ]
then
    echo "Data file " $1 " missing or bad written"
    exit
fi
if [ ! -e $2 ]
then
    echo "Reference table header file " $2 " missing or bad
written"
    exit
fi
if [ ! "$(echo $3 | grep '^[:digit:]*$')" ]
then
    echo $3 " is not a number"
    exit
fi
if [ ! "$(echo $4 | grep '^[:digit:]*$')" ]
then
    echo $4 " is not a number"
    exit
fi
#-----
#####VARIABLES#####
#-----
njob=$3
niter_intra_job=$4
output_dir_name=output2
aux_script=./clusaux
rep=`pwd`#
#-----
#####AUX SCRIPT#####
#-----
mkdir $output_dir_name
echo 'mkdir /tmp/travail$$' > $aux_script
echo 'cp ./diyabc_sim /tmp/travail$$' >> $aux_script
echo 'cp '$1' /tmp/travail$$' >> $aux_script
echo 'cp '$2' /tmp/travail$$' >> $aux_script
echo 'cd /tmp/travail$$' >> $aux_script
echo './diyabc_sim ${2%.[^.]*}' '$1 $2' >> $aux_script
echo 'cp *.reftable '$rep'/'$output_dir_name'/.>> $aux_script
echo 'cd /tmp' >> $aux_script
echo 'rm -rf travail$$' >> $aux_script
#-----
#####SEND TO CLUSTER#####
#-----
for ((j=1; j<=$njob; j++))
do
qsub -cwd $aux_script $j $niter_intra_job
done
#-----
#####CLEANING#####
#-----
rm -rf $aux_script

```

5.1 A note about the random number generator used in DIYABC

- 2 By nature, a random number generator (RNG) is a sequential algorithm, as described in Figure 2 below.
 3 Indeed, we shall describe a RNG by its updating function f changing deterministically the internal state.
 4 Each time the user requires a new realization of the uniform distribution over $[0; 1)$, the algorithm derives
 5 a value u_k from the current internal state i_k and then updates this state with f . Hence a first and
 6 important issue for parallel Monte Carlo computations is to design independent RNGs that might run in
 7 parallel while minimizing the communications between processors. It is quite standard to use as many
 8 RNGs as computing cores in the computer or in the cluster of computers.

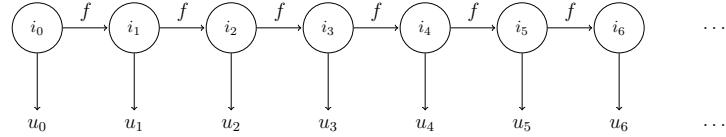


Figure 2: Random Number Generator. A RNG is an algorithm that produces a sequence of floating numbers, say u_0, u_1, \dots , that resembles a sequence of independent random numbers, uniformly distributed over $[0; 1)$. It uses a sequence of internal states, say i_0, i_1, \dots , which are computed by recurrence, namely, $i_{k+1} = f(i_k)$. The first internal state i_0 is often named the seed.

The second version of DIYABC uses the Dynamic Creator (DCMT) of Matsumoto and Nishimura (2000) to look for a set of independent Mersenne-Twister generators. Actually, the updating function f of a Mersenne-Twister generator is parametrized by a few integer numbers. The output of the DCMT is a set of N updating functions, say $\{f^{(1)}, \dots, f^{(N)}\}$, producing independent streams. That is, the n -th RNG is a sequence of internal states $i_0^{(n)}, i_1^{(n)}, i_2^{(n)}, \dots$ satisfying $i_{k+1}^{(n)} = f^{(n)}(i_k^{(n)})$ that gives rise to a sequence of independent, uniformly distributed numbers $u_0^{(n)}, u_1^{(n)}, u_2^{(n)}, \dots$. We found that the DCMT was simple to use and gave good results. There is no limitation on the number N of RNGs it produces. Once initialized, the different RNGs do not require any communication between them and each of them runs as quickly as a single Mersenne-Twister generator. But an important limitation is that it is impossible to add a new RNG to the set produced by the DCMT. Practically, this means that we have to know *a priori* a bound on the number of jobs working together in parallel. See XXX Alex XXX.

¹ Bibliography

- ² Beaumont, M. A., W. Zhang and D. J. Balding, 2002. Approximate Bayesian Computation in Population
³ Genetics. *Genetics* **162**, 2025-2035.
- ⁴ Beaumont, M.A., 2008. Joint determination of topology, divergence time, and immigration in population
⁵ trees. In Simulation, Genetics, and Human Prehistory, eds. S. Matsumura, P. Forster, C. Renfrew.
⁶ McDonald Institute Press, University of Cambridge (*in press*).
- ⁷ Begg, C.B. and R. Gray, 1984. Calculation of polychotomous logistic regression parameters using indi-
⁸ vidualized regressions. *Biometrika*, **71**, 11-18.
- ⁹ Belkhir K., Borsa P., Chikhi L., Raufaste N. and F. Bonhomme, 1996-2004 GENETIX 4.05, logiciel sous
¹⁰ Windows TM pour la gntique des populations. Laboratoire Gnome, Populations, Interactions, CNRS
¹¹ UMR 5171, Universit de Montpellier II, Montpellier (France).
- ¹² Bertorelle, G. and L. Excoffier, 1998. Inferring admixture proportion from molecular data. *Mol. Biol.*
¹³ *Evol.* **15**, 1298-1311.
- ¹⁴ Choisy, M., P. Franck and J.M. Cornuet, 2004. Estimating admixture proportions with microsatellites :
¹⁵ comparison of methods based on simulated data. *Mol. Ecol.* **13**, 955-968.
- ¹⁶ Chakraborty R and L Jin, 1993. A unified approach to study hypervariable polymorphisms: statistical
¹⁷ considerations of determining relatedness and population distances. EXS. **67**, 153175.
- ¹⁸ Cornuet, J. M., M. A. Beaumont, A. Estoup and M. Solignac, 2006. Inference on microsatellite mutation
¹⁹ processes in the invasive mite, *Varroa destructor*, using reversible jump Markov chain Monte Carlo.
²⁰ *Theoret. Pop. Biol.* **69**, 129-144.
- ²¹ Cornuet J.M., V. Ravigné and A. Estoup, 2010. Inference on population history and model checking using
²² DNA sequence and microsatellite data with the software DIYABC (v1.0). *submitted*.
- ²³ Cornuet J.M., F. Santos, M.A. Beaumont, C.P. Robert, J.M. Marin, D.J. Balding, T. Guillemaud and
²⁴ A. Estoup, 2008. Inferring population history with DIYABC: a user-friendly approach to Approximate
²⁵ Bayesian Computations. *Bioinformatics*, **24** (23), 2713-2719.
- ²⁶ Estoup, A., M. Solignac, M. Harry and J.M. Cornuet, 1993. Characterization of $(GT)_n$ and $(CT)_n$
²⁷ microsatellites in two insect species: *Apis mellifera* and *Bombus terrestris*. *Nucl. Ac. Res.*, **21**, 1427-
²⁸ 1431.
- ²⁹ Estoup, A., I. J. Wilson, C. Sullivan, J. M. Cornuet and C. Moritz, 2001 Inferring population history
³⁰ from microsatellite et enzyme data in serially introduced cane toads, *Bufo marinus*. *Genetics*, **159**,
³¹ 1671-1687.
- ³² Estoup, A., P. Jarne and J.M. Cornuet, 2002. Homoplasy and mutation model at microsatellite loci and
³³ their consequences for population genetics analysis. *Mol. Ecol.*, **11**, 1591-1604.
- ³⁴ Estoup, A. and S. M. Clegg, 2003. Bayesian inferences on the recent islet colonization history by the bird
³⁵ *Zosterops lateralis lateralis*. *Mol. Ecol.* **12**: 657-674.
- ³⁶ Estoup, A., M.A. Beaumont, F. Sennedot, C. Moritz and J.M. Cornuet, 2004. Genetic analysis of complex
³⁷ demographic scenarios : spatially expanding populations of the cane toad, *Bufo marinus*. *Evolution*, **58**,
³⁸ 2021-2036.

- 1 Excoffier, L., A. Estoup and J.M. Cornuet, 2005. Bayesian analysis of an admixture model with mutations
2 and arbitrarily linked markers. *Genetics* **169**, 1727-1738.
- 3 FAGUNDES, N.J.R., N. RAY, M.A. BEAUMONT, S. NEUENSCHWANDER, F. SALZANO, S.L. BONATTO
4 AND L. EXCOFFIER, 2007. Statistical evaluation of alternative models of human evolution. *Proc. Natl.
5 Acad. Sc.*, **104** : 17614-17619.
- 6 Fu, Y.X. and Chakraborty, R., 1998. Simultaneous estimation of all the parameters of a stepwise mutation
7 model. *Genetics*, **150**, 487-497.
- 8 Garza JC and E Williamson, 2001. Detection of reduction in population size using data from microsatellite
9 DNA. *Mol. Ecol.* **10**, 305-318.
- 10 Gelman, A., J.B. Carlin, H.S. Stern and D.B. Rubin, 1995. *Bayesian Data Analysis*. Chapman et Hall,
11 London, 526p.
- 12 Goldstein DB, Linares AR, Cavalli-Sforza LL, and Feldman MW, 1995. An evaluation of genetic distances
13 for use with microsatellite loci. *Genetics* **139**, 463-471.
- 14 Goudet, J. ,1995. FSTAT (Version 1.2): A computer program to calculate F- statistics. *J. Hered.* **86**,
15 485-486.
- 16 Griffiths, R.C. and S. Tavaré, 1994. Simulating probability distributions in the coalescent. *Theor. Pop.
17 Biol.* **46**, 131-159.
- 18 Guillemaud T., M.A. Beaumont, M. Ciosi, J.M. Cornuet and A. Estoup, 2010. Inferring introduction
19 routes of invasive species using approximate Bayesian computation on microsatellite data. *Heredity*,
20 **104**, 88-99.
- 21 Haag-Liautard C., N. Coffey, D. Houle, M. Lynch, B. Charlesworth and P.D. Keightley, 2008. Direct
22 estimation of the mitochondrial DNA mutation rate in *Drosophila melanogaster*. *Plos Biol.* **6**, e204.
- 23 Hamilton, G., M. Stoneking and L. Excoffier, 2005. Molecular analysis reveals tighter social regulation
24 of immigration in patrilocal populations than in matrilocal populations. *Proc. Natl. Acad. Sci. USA*,
25 **102**, 7476-7480.
- 26 Hasegawa, M., Kishino, H and Yano, T., 1985. Dating the human-ape splitting by a molecular clock of
27 mitochondrial DNA. *Journal of Molecular Evolution* 22:160-174.
- 28 Hudson,R. R., M. Slatkin and W.P. Maddison, 1992. Estimation of levels of gene flow fom DNA sequence
29 data. *Genetics*, 132, 583-589.
- 30 Ihaka R. and R. Gentleman, 1996. *R*: a language for data analysis and graphics. *J. Comput. Graph. Stat.*,
31 **5**, 299-314
- 32 Ingvarsson P.K., 2008. Multilocus patterns of nucleotide polymorphism and the demographic history of
33 *Populus tremula*. *Genetics*, 180: 329-340.
- 34 Jukes, TH and Cantor, CR., 1969. Evolution of protein molecules. Pp. 21-123 in H. N. Munro, ed.
35 *Mammalian protein metabolism*. Academic Press, New York.
- 36 Kimura, M., 1980. A simple method for estimating evolutionary rate of base substitution through com-
37 parative studies of nucleotide sequences. *Journal of Molecular Evolution* 16:111-120.
- 38 Lombaert E., T. Guillemaud, J.M. Cornuet, T. Malausa, B. Facon and A. Estoup, 2010. Bridgehead effect in the worldwide invasion of the biocontrol harlequin ladybird. *PLoS ONE*,
39 <http://dx.plos.org/10.1371/journal.pone.0009743>.
- 41 Matsumoto M and T Nishimura, 2000. Dynamic Creation of Pseudorandom Number Generators. *Monte
42 Carlo and Quasi-Monte Carlo Methods 1998*, Springer, pp 56-69.
- 43 Miller N, A. Estoup, S. Toepfer, D Bourguet, L. Lapchin, S. Derridj, K.S. Kim, P Reynaud, F. Furlan and
44 T. Guillemaud, 2005. Multiple Transatlantic Introductions of the Western Corn Rootworm. *Science*,
45 **310**, p. 992

- ¹ Nei M., 1972. Genetic distance between populations. *Am. Nat.* **106**:283-292
- ² Nei M., 1987. *Molecular Evolutionary Genetics*. Columbia University Press, New York, 512 pp.
- ³ Ohta, T. and Kimura, M., 1973. A model of mutation appropriate to estimate the number of electrophoretically detectable alleles in a finite population.
- ⁴
- ⁵ Pascual, M., M.P. Chapuis, F. Mestres, J. Balanyá, R.B. Huey, G.W. Gilchrist, L. Serra and A. Estoup, 2007. Introduction history of *Drosophila subobscura* in the New World : a microsatellite based survey using ABC methods. *Mol. Ecol.*, **16**, 3069-3083.
- ⁶
- ⁷
- ⁸ Pritchard, J., M. Seielstad, A. Perez-Lezaun and M. Feldman, 1999. Population growth of human Y chromosomes: a study of Y chromosome microsatellites. *Mol. Biol. Evol.* **16**, 1791-1798.
- ⁹
- ¹⁰ Rannala, B., and J. L. Mountain, 1997. Detecting immigration by using multilocus genotypes. *Proc. Nat. Acad. Sci. USA* **94**, 9197-9201.
- ¹¹
- ¹² Raymond M., and F. Rousset, 1995. Genepop (version 1.2), population genetics software for exact tests and ecumenicism. *J. Hered.*, **86**, 248-249
- ¹³
- ¹⁴ Schug M.D., T.F. Mackay and C.F. Aquadro, 1997. Low mutation rates of microsatellite loci in *Drosophila melanogaster*. *Nat Genet.* **15**, 99-102.
- ¹⁵
- ¹⁶ Stephens, M. and P. Donnelly, 2000, Inference in molecular population genetics (with discussion). *J. R. Stat. Soc. B* **62**, 605-655.
- ¹⁷
- ¹⁸ Tajima, F., 1989. Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics* **123**: 585-595
- ¹⁹
- ²⁰ Tamura, K., and M. Nei., 1993. Estimation of the number of nucleotide substitutions in the control region of mitochondrial DNA in humans and chimpanzees. *Molecular Biology and Evolution* **10**:512-526.
- ²¹
- ²² Vzquez F.J., T. Prez, J. Albornoz and A. Domnguez, 2000. Estimation of microsatellite mutation rates in *Drosophila melanogaster*. *Genet Res.*, **76**, 323-6.
- ²³
- ²⁴ Weir BS and CC Cockerham , 1984. Estimating F-statistics for the analysis of population structure. *Evolution* **38**: 1358-1370.
- ²⁵
- ²⁶ Whittaker, J.C., Harbord, R.M., Boxall, N., Mackay, I., Dawson, G. and Sibly, R.M. 2003. Likelihood-based estimation of microsatellite mutation rates, *Genetics*, **164**, 781-787.
- ²⁷