

December 9, 2015

# BayPass version 2.1

## User Manual

---

BAYPASS code © [INRA](#)

This document © Mathieu Gautier 2015

# Contents

<b>1</b>	<b>Overview</b>	<b>4</b>
<b>2</b>	<b>Before you start</b>	<b>4</b>
2.1	How to get BAYPASS?	4
2.2	How to compile BAYPASS?	5
2.3	Input file format	6
2.3.1	The genotyping data file [always required]	6
2.3.2	The pool (haploid) size file [only required for Pool-Seq data]	8
2.3.3	The covariate data file [required for the covariate modes]	8
2.3.4	The covariance matrix file [optional, required for the AUX covariate mode]	9
<b>3</b>	<b>Running BayPass</b>	<b>9</b>
3.1	Overview of the different models available in BAYPASS	9
3.1.1	The core model	10
3.1.2	The standard covariate model and extensions	11
3.1.3	The auxiliary covariate model	12
3.2	Detailed overview of all the options	15
3.3	Format of the output files	21
<b>4</b>	<b>Miscellaneous R functions</b>	<b>24</b>
4.1	The R function <i>simulate.baypass()</i>	25
4.1.1	Description	25
4.1.2	Usage	25
4.1.3	Arguments (in alphabetic order)	25
4.1.4	Values	27
4.1.5	Examples	29
4.2	The R function <i>fmd.dist()</i>	29
4.2.1	Description	29
4.2.2	Usage	29
4.2.3	Arguments	29
4.2.4	Values	30
4.2.5	Example	30
4.3	The R function <i>geno2YN()</i>	30
4.3.1	Description	30
4.3.2	Usage	30
4.3.3	Arguments	30
4.3.4	Values	30
4.3.5	Example	31

<b>5</b>	<b>Worked Examples</b>	<b>31</b>
5.1	Cattle allele count data . . . . .	31
5.1.1	Analysis under the core model mode: MCMC is run under the core model . . . . .	31
5.1.2	Analysis under the IS covariate mode: MCMC is run under the core model . . . . .	33
5.1.3	Analysis under the MCMC covariate mode: MCMC is run under the STD model . . . . .	34
5.1.4	Analysis under the AUX covariate mode: MCMC is run under the AUX model . . . . .	35
5.2	Littorina Pool-Seq read count data . . . . .	36
<b>6</b>	<b>Credits</b>	<b>36</b>
<b>7</b>	<b>Copyright</b>	<b>37</b>
<b>8</b>	<b>Contact</b>	<b>37</b>
	<b>Bibliography</b>	<b>38</b>
	<b>Appendix A Comparisons of the computational efficiency of the different version of BayPass</b>	<b>40</b>

# 1 Overview

The package BAYPASS is a population genomics software which is primarily aimed at identifying genetic markers subjected to selection and/or associated to population-specific covariates (e.g., environmental variables, quantitative or categorical phenotypic characteristics). The underlying models explicitly account for (and may estimate) the covariance structure among the population allele frequencies that originates from the shared history of the populations under study. Note that, apart from standard population genetics studies, BAYPASS is generic enough to be also suited to the analyses of data from other kinds of experiments in which the allele frequency covariance structure is simpler (e.g., experimental evolution). The genetic data typically consists of allele (when derived from individual genotype calls) or read (when derived from Pool-Seq experiments) counts at several markers (for now, BAYPASS is restricted to the analysis of bi-allelic markers) in several populations. Note that BAYPASS can handle missing data (no count available in one or several populations) which might be helpful in some contexts.

The core BAYPASS model is based on the BAYENV model which was proposed by Coop *et al.* (2010) and Günther and Coop (2013). However, as detailed in Gautier (2015), in addition to a complete and independant reprogramming of the core Markov Chain Monte Carlo (MCMC) algorithm, BAYPASS allows to monitor most of the parameters and the priors of the original models and to introduce various extensions (via e.g., the optional addition of hyper-parameters, the modeling of spatial dependency among consecutive markers).

BAYPASS is written in Fortran90. The source code and compilation instructions for various platforms (OS X, Windows, Linux) are available. The executable reads data file(s) supplied by the user, and a number of options can be passed through the command line. Some R functions are also provided in the package to facilitate interpretation of the resulting outputs.

This document provides information about how to format the data file, how to specify the user-defined parameters, and how to interpret the results.

## 2 Before you start

### 2.1 How to get BayPass?

Download the archive from <http://www1.montpellier.inra.fr/CBGP/software/bypass/>, and extract it from a terminal:

```
tar -zxvf bypass_2.1.tar.gz
```

## 2.2 How to compile BayPass?

The source files are to be found in the `src` subdirectory. BAYPASS is coded in Fortran90 and can therefore be compiled for any system supporting a Fortran90 compiler using the provided `Makefile`. This `Makefile` is designed to work with either the free compiler `gfortran`<sup>1</sup> or the commercial `ifort` intel® Fortran compiler<sup>2</sup>. As a consequence, using another Fortran90 compiler requires modifying the `Makefile` accordingly. Note also that BAYPASS uses `OpenMP` (<http://openmp.org/wp/>) to implement multithreading, which allows parallel calculation on computer systems that have multiple CPUs or CPUs with multiple cores. Users thus have to make sure that the corresponding libraries are installed (which is usually the case, at Linux OS or following compiler installation previously described<sup>1</sup>). The following instructions run within the `src` subdirectory allows to compile the code and to produce a binary:

- using the `gfortran` free compiler (the command should automatically produce an executable called `g_baypass`):  
`make clean all FC=gfortran`
- using the `ifort` intel® Fortran compiler (the command should automatically produce an executable called `i_baypass`):  
`make clean all FC=ifort`

From my (limited) experience, I would recommend using the `ifort` intel® Fortran compiler instead of `gfortran`. Indeed, although `gfortran` options could probably be better adapted to further improve the speed of the executable (any feedback about this is welcome), using the options considered in the `Makefile`, the `ifort` executable is more than two times faster than the `gfortran` one as illustrated in the comparison below<sup>3</sup>. An example comparison of the computational performances of the different version and binaries of the program is illustrated in Table 1.

---

<sup>1</sup> If not already installed in your system, binaries are available at <https://gcc.gnu.org/wiki/GFortranBinaries> and easy to install for most Windows, Mac and Linux OS versions (many thanks to Andrew Beckerman for pointing this webpage to me!)

<sup>2</sup> Intel proposes a free personal non-commercial unsupported version that is essentially identical to the commercial version: <https://software.intel.com/en-us/qualify-for-free-software/academicresearcher>

<sup>3</sup> I however noticed that the `ifort` executable might sometimes crash in an unpredictable way when reading the command line arguments (for a reason that remains very obscure to me and that I am trying to figure out!)

Under Linux (or MacOS), before the first use, make sure to give appropriate execution rights to the program. For instance you may run:

```
chmod +x baypass
```

In the following, it is assumed that the program was made executable and accessible in your path. For instance, under Linux, this may be achieved by copying the executable in a directory declared in the path (e.g., `/usr/local/bin`) or by adding the program to the `$PATH` system variable (using the `export` command)

## 2.3 Input file format

Depending on the type of analyses, different data files might be required by the program. The following examples of the different input files are available in the `examples` directory:

- `geno.bta14`: this file contains allele count data for 18 French cattle breeds at 1,394 SNPs mapping to the BTA14 bovine chromosome (see Gautier (2015) for details).
- `bta.pc1`: this file contains the SMS (Synthetic Morpholy Score) for the 18 French cattle breeds (see Gautier (2015) for details).
- `omega.bta`: this file contains the matrix  $\Omega$  for the 18 French cattle breeds ( $\hat{\Omega}_{BTA}^{bpas}$ ) as estimated under the core model from the whole genome SNP data (see Gautier (2015) for details).
- `lsa.geno`: this file contains read count data (Pool-Seq data) for 12 populations from the *Littorina saxatilis* marine snail (Westram *et al.*, 2014) at 2,500 SNPs randomly chosen among the ones analysed in Gautier (2015) (but including the ca. 150 outlier SNPs identified).
- `lsa.poolsize`: this file contains the haploid pool sizes of the 12 *Littorina saxatilis* populations.
- `lsa.ecotype`: this file contains the code for the ecotype of the 12 *Littorina saxatilis* populations ( $-1$  for the "crab" habitat and  $1$  for the "wave" habitat).

### 2.3.1 The genotyping data file [always required]

The genotyping data files contain allele or read count (for PoolSeq experiment) data for each of the  $nsnp$  markers assayed in each of the  $npops$  populations sampled. The genotyping data file is simply organised as a matrix with

*n*snp rows and  $2 * n$ pop columns. The row field separator is a space. More precisely, each row corresponds to one marker and the number of columns is twice the number of populations because each pair of numbers corresponds to each allele (or read counts for PoolSeq experiment) counts in one population<sup>4</sup>.

As a schematic example, the genotyping data input file for allele count data should read as follows:

```
--- file begins here ---  
81 19 86 14 2 98 8 92 32 68 23 77  
89 11 81 19 9 91 1 99 27 73 27 73  
89 11 91 9 0 0 15 85 77 23 80 20
```

[...97 more lines...]

```
--- file ends here ---
```

In this example, there are 6 populations and 100 SNP markers. At the first SNP, in the first population, there are 81 copies of the first allele, and 19 copies of the second allele. In the second population, there are 86 copies of the first allele, and 14 copies of the second allele, etc. Note that both alleles in the third SNP in the third population have 0 copie. This marker will be treated as a missing data in the corresponding population. The file named `geno.bta14` in the `example` directory provides a more realistic example.

Similarly, as a schematic example, the genotyping data input file for allele count data should read as follows:

```
--- file begins here ---  
71 8 115 0 61 36 51 39 10 91 69 58  
82 0 91 0 84 14 24 57 28 80 18 80  
93 28 112 30 0 0 0 113 33 68 0 106
```

[...97 more lines...]

```
--- file ends here ---
```

In this example, there are also 6 populations and 100 SNP markers. At the first SNP, in the first population, there are 71 reads of the first allele, and 8 reads of the second allele. In the second population, there are 115 reads of the first allele, and 0 read of the second allele, etc. Note that both alleles in the third SNP in the third population have 0 copie. This marker will be treated as a missing data in the corresponding population. The file named `lsa.geno` in the `example` directory provides a more realistic example.

---

<sup>4</sup>For now, BAYPASS is restricted to bi-allelic marker

For Pool-Seq data to be analyzed properly (i.e. not like allele count data), it is necessary to provide a file with the (haploid) size of each pool (see 2.3.2).

### 2.3.2 The pool (haploid) size file [only required for Pool-Seq data]

For Pool-Seq experiment, the haploid size (twice the number of pooled individuals for diploid species) of each population should be provided. As a schematic example, the pool (haploid) size file should read as follows:

```
--- file begins here ---  
60 75 100 90 80 50  
--- file ends here ---
```

In this example, there are 6 populations with respective haploid sample sizes of 60 (first population), 75 (second population), 100 (third population), 90 (fourth population), 80 (fifth population) and 50 (sixth population). The order of the populations in the pool size file must be the same as in the allele count (and the covariate) data file(s). The file named `lsa.poolsize` in the `example` directory provides a more realistic example.

### 2.3.3 The covariate data file [required for the covariate modes]

The values of the covariates (e.g., environmental data, phenotypic traits, etc.) for the different populations should be provided in a file with the format exemplified in the following:

```
--- file begins here ---  
150 1500 800 300 200 2500  
181.5 172.6 152.3 191.8 154.2 166.8  
1 1 0 0 1 1  
0.1 0.8 -1.15 1.6 0.02 -0.5  
--- file ends here ---
```

In this example, there are 6 populations (columns) and 4 covariates (row). The first covariate might be viewed as a typical environmental covariate, like altitude in meters (the first population is living at ca. 150m above the sea level, the second at ca. 1,500m, and so on), the second as a quantitative traits like average population sizes in cm (individuals from the first population are 181.5 cm height on average, individuals from the second population 172.6 cm, and so on), the third covariate as a typical binary trait like presence (1, for the first, second, fifth and sixth populations) or absence (0, for the third and fourth populations) and the last might be viewed as a synthetic variable



like the first principal components of a PCA. The order of the populations (columns) in the covariate data file must be the same as in the allele count (and the pool size) data file(s).

The files named `bta.pc1` and `lsa.ecotype` in the `example` directory provide alternative real-life examples.

Note that in most cases, it is (strongly) recommended to scale each covariate (so that  $\hat{\mu} = 0$  and  $\hat{\sigma}^2 = 1$  for each covariable). The `scalecov` option allows to perform this step automatically prior to analysis, if needed.

### 2.3.4 The covariance matrix file [optional, required for the AUX covariate mode]

For some applications (see below), it might be interesting (e.g., to parallelize some analyses) or required (when using the AUX covariate mode) to provide the population covariance matrix  $\Omega$ . As a schematic example, the covariance matrix file reads as follows:

```
--- file begins here ---
0.098053 0.019595 0.032433 -0.029601 -0.024190 -0.029247
0.019595 0.160147 0.018942 -0.027348 -0.039733 -0.039010
0.032433 0.018942 0.149962 -0.054973 -0.058700 -0.057288
-0.029601 0.027348 0.054973 0.187511 0.221914 0.165862
-0.024190 0.039733 0.058700 0.221914 0.562666 0.260231
-0.029247 0.039010 0.057288 0.165862 0.260231 0.219761
--- file ends here ---
```

In this example, there are 6 populations. Hence, the population covariance matrix is a  $6 \times 6$  squared symmetric matrix. The order of the populations (columns and rows) in the matrix  $\Omega$  should be the same as the columns in the allele count (and the pool size and the covariate) data file(s). Note that this file is produced in the appropriate format by the program when running BAYPASS under the core model (see 3.3).

The file named `omega.bta` provides a real-life example.

## 3 Running BayPass

### 3.1 Overview of the different models available in BayPass

Directed Acyclic Graphs (DAG) of the different family of models are represented in Figure 1 (see Gautier (2015) for details). Briefly, three types of

(closely related) models might be investigated using BAYPASS, considering either Allele count data (left panel in Figure 1) or Read count data (right panel in Figure 1) as obtained from Pool-Seq experiments.

### 3.1.1 The core model

The core model depicted in Figure 1A might be viewed as a generalisation of the model proposed by Nicholson *et al.* (2002) and was first proposed by Coop *et al.* (2010). This model is the one considered by BAYPASS when no covariate data file is provided and is actually nested in the others models.

The main parameter of interest is the (scaled) covariance matrix of population allele frequencies  $\mathbf{\Omega}$  resulting from their (possibly unknown and complex) shared history. Conversely, one might rely on this matrix for demographic inference. For instance,  $\mathbf{\Omega}$  might easily be converted (e.g., using the `cov2cor()` function in R stats package) into a correlation matrix  $\mathbf{\Sigma}$  further interpreted as a similarity matrix. From this latter matrix, one may define a distance (dissimilarity) matrix (e.g.,  $d_{ij} = 1 - |\rho_{ij}|$  where  $d_{ij}$  is the distance between populations  $i$  and  $j$  and  $\rho_{ij}$  is the element  $ij$  of  $\mathbf{\Sigma}$ ) to perform hierarchical clustering<sup>5</sup> and summarise the history of the population as a bifurcating phylogenetic tree (without gene flow). A more complex demographic inference based on an interpretation the matrix  $\mathbf{\Omega}$  (although estimated in a less accurate way) in terms of tree with migration has been recently proposed by Pickrell and Pritchard (2012).

The core model allows to perform genome scan for differentiation (covariate-free) using the XtX statistics as introduced by Günther and Coop (2013) which is computed by default in BAYPASS. The main advantage of this approach stems is to explicitly account for the covariance structure in population allele frequencies (via estimating  $\mathbf{\Omega}$ ) resulting from the demographic history of the populations. To identify outlier loci (based on the XtX statistics), the R function `simulate.baypass()` provided in the `utils` directory of the package (see 4) allows to simulate data under the inference model (e.g., using posterior estimates of  $\mathbf{\Omega}$  and any other hyperparameters) which might further be analysed to calibrate the neutral XtX distribution (Gautier, 2015).

---

<sup>5</sup>For an interesting discussion and examples in R, see <http://research.stowers-institute.org/mcm/efg/R/Visualization/cor-cluster/index.htm>

In the current implementation of BAYPASS, the prior distribution for  $\mathbf{\Omega}$  is an Inverse-Wishart:  $\mathbf{\Omega} \sim W_J^{-1}(\rho \mathbf{I}_J, \rho)$  (where  $J$  is the number of populations). By default  $\rho = 1$  (rather than  $\rho = J$  as in BAYENV) which was found as the most reliable value (Gautier, 2015). Similarly, the hyperparameters  $a_\pi$  and  $b_\pi$  of the prior  $\beta$  distribution for the overall (across population) SNP allele frequencies are estimated by default. However, they might be fixed to  $a_\pi = b_\pi = 1$  (as in e.g., BAYENV) using `fixpibetapar` option or to any other values using further the `betaprior` option (3.2).

### 3.1.2 The standard covariate model and extensions

The standard covariate model is represented in Figure 1B and is the one considered by default in BAYPASS when a covariate data file is provided using `efile` option (3.2). This model allows to evaluate to which extent the population covariable(s)  $k$  is (linearly) associated to each marker  $i$  (which are assumed independent given  $\mathbf{\Omega}$ ) by the introduction of the regression coefficients  $\beta_{ik}$  (for convenience the indices  $k$  for covariables are dropped in Figure 1B).

In the current implementation of BAYPASS, the prior distribution for the  $\beta_{ik}$ 's is Uniform:  $\beta_{ik} \sim \text{Unif}(\beta_{\min}, \beta_{\max})$ . By default,  $\beta_{\min} = -0.3$  and  $\beta_{\max} = 0.3$  but these values might be changed by the user with the `minbeta` and `maxbeta` options respectively (3.2). Note that in BAYENV (Coop *et al.*, 2010),  $\beta_{\min} = -0.1$  and  $\beta_{\max} = 0.1$ .

The estimation of the  $\beta_{ik}$  regression coefficients for each SNP  $i$  may be performed using two different approaches (Gautier, 2015):

- Using an importance sampling estimator (IS) which is the default option and also allows the computation of Bayes Factor to compare on an individual SNP (and covariable) basis the two alternative models, namely the model with association ( $\beta_{ik} \neq 0$ ) against the null model ( $\beta_{ik} = 0$ ). Bayes Factor ( $\text{BF}_{\text{IS}}$ ) and  $\beta_{ik}$  IS algorithm are inspired from Coop *et al.* (2010) and are described in details elsewhere (Gautier, 2015). Note that the IS estimation procedure is based on a numerical integration that requires the definition of a grid covering the whole support of the  $\beta_{ik}$  prior distribution. In BAYPASS, the grid consists in  $n_\beta$  (by default  $n_\beta = 201$ ) equidistant points from  $\beta_{\min}$  to  $\beta_{\max}$  (including the boundaries) leading to a lag between two successive values equal to  $\frac{\beta_{\max} - \beta_{\min}}{n_\beta - 1}$  (i.e., 0.003 with default values). Other values for  $n_\beta$  might be supplied by the user with the `nbetagrid` option (3.2).
- Using an MCMC algorithm (activated via the `covmcmc` option). In this case, the user should provide the matrix  $\mathbf{\Omega}$  (e.g., using posterior

estimates available from a previous analysis) and it is recommended to consider only one covariable at a time (particularly if some covariables are correlated).

### 3.1.3 The auxiliary covariate model

The auxiliary covariate model, represented in Figure 1C and activated with the `auxmodel` option, is an extension of the previous model (Figure 1B) involving the introduction of a Bayesian (binary) auxiliary variable  $\delta_{ik}$  for each regression coefficient  $\beta_{ik}$  (Gautier, 2015). In a similar population genetics context, this modelling was also proposed by Riebler *et al.* (2008) to identify markers subjected to selection in genome-wide scan of adaptive differentiation based on a  $\mathcal{F}$ -model.

Here, the auxiliary variable actually indicates whether a specific SNP  $i$  can be regarded as associated to a given covariable  $k$  ( $\delta_{ik} = 1$ ) or not ( $\delta_{ik} = 0$ ). By looking at the posterior distribution of each auxiliary variable, it is then straightforward to derive a Bayes Factor ( $\text{BF}_{\text{mc}}$ ) to compare both models while dealing with multiple testing issues (Gautier, 2015). In addition, the introduction of a Bayesian auxiliary variable makes it easier to account for spatial dependency among markers. In BAYPASS, the general form of the  $\delta_{ik}$  prior distribution is indeed that of an 1D Ising model with a parametrization analogous to the one proposed in a similar context by Duforet-Frebourg *et al.* (2014):  $\pi(\delta_{\mathbf{k}}) \propto P^{s_1} (1 - P)^{s_0} e^{\eta b_{\text{is}}}$ , where  $\delta_{\mathbf{k}}$  is the vector of the nsnp auxiliary variables for covariable  $k$ ,  $s_1$  and  $s_0$  are the number of SNPs associated (i.e. with  $\delta_{ik} = 1$ ) and not associated (i.e. with  $\delta_{ik} = 0$ ) to the covariable<sup>6</sup>, and  $\eta$  corresponds to the number of pairs of consecutive markers (neighbors) that are in the same state at the auxiliary variable<sup>7</sup> (i.e.,  $\delta_{i,k} = \delta_{i+1,k}$ ). The parameter  $P$  broadly corresponds to the prior proportion of SNP associated to the covariable. In the BAYPASS auxiliary model,  $P$  is assumed a priori beta distributed:  $P \sim \beta(a_P, b_P)$ . By default,  $a_P = 0.02$  and  $b_P = 1.98$  (this values might be changed by the user with the `auxPbetaprior` option) which amounts to consider that only a small fraction of the SNPs ( $\frac{a_P}{a_P+b_P} = 1\%$ ) are a priori expected to be associated to the covariable while allowing some uncertainty on this key parameter (e.g., the prior probability of  $P > 10\%$  being equal to 0.028 with these parameters). The parameter  $b_{\text{is}}$ , called the inverse temperature in the Ising (and Potts) model literature, determines the level of spatial homogeneity of the auxiliary variables between neighbors. In

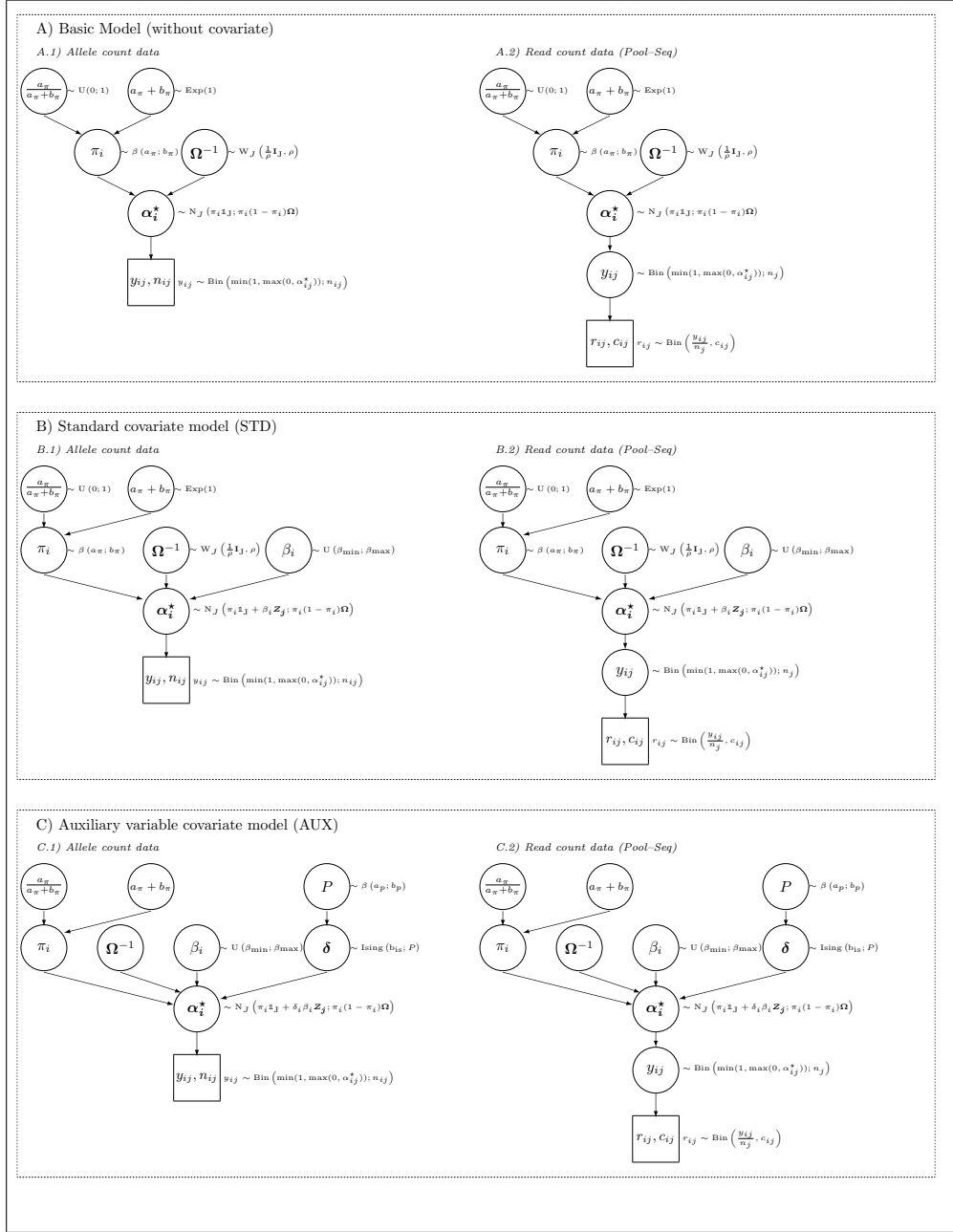
---


$$^6 s_1 = \sum_{i=1}^{\text{nsnp}} \delta_{ik} = 1 \text{ and } s_0 = \sum_{i=1}^{\text{nsnp}} \delta_{ik} = 0$$

$$^7 \eta = \sum_{i \sim j} \mathbf{1}_{\delta_{ik} = \delta_{jk}}$$

BAYPASS,  $b_{is} = 0$  by default implying that auxiliary variables are independent (no spatial dependency). Note that  $b_{is} = 0$  amounts to assume the  $\delta_{ik}$  follows a Bernoulli distribution with parameter  $P$ . Conversely,  $b_{is} > 0$  leads to assume that the  $\delta_{ik}$  with similar values tend to cluster according to the underlying SNP positions (the higher the  $b_{is}$ , the higher the level of spatial homogeneity). In biological terms, SNP associated to a given covariable might be expected to cluster due to Linkage Disequilibrium with the underlying (possibly not genotyped) causal variant(s). In practice,  $b_{is} = 1$  is commonly used and a value of  $b_{is} \leq 1$  is recommended.

In technical terms, the overall parametrisation of the Ising prior assumes no external field and no weight (as in the so-called compound Ising model) between the neighboring auxiliary variables. In the current implementation of the BAYPASS auxiliary covariate model (when  $b_{is} > 0$ ), the information about the distances between SNPs is therefore not accounted for. Only the relative position of markers are considered. For the applications where this modeling might be relevant (whole genome scan), this corresponds to assuming a relative homogeneity in marker spacing as measured by genetic (rather than physical) distances (which might be unavailable, in practice).



**Figure 1:** Directed Acyclic Graphs of the different hierarchical Bayesian models available in BAYPASS (see 3.1). For each model, optional (hyper-)parameters are displayed in orange.

## 3.2 Detailed overview of all the options

BAYPASS is a command-line executable. The ASCII hyphen-minus ("-") is used to specify options. As specified below, some options take integer or float values and some options do not. Here is an example call of the program:

```
baypass -npop 12 -gfile data.geno -efile env.data -outprefix ana1
```

The full list of the options accepted by BAYPASS are printed out using the command: `baypass -help` as follows:

```
Version 2.1

Usage: BayPass [options]

Options:
I) General Options:
  -help          Display the help page
  -npop          INT      Number of populations (always required)
  -gfile         CHAR     Name of the Genotyping Data File (always required)
  -efile         CHAR     Name of the covariate file: activate covariate mode (def="")
  -scalecov     CHAR     Scale covariates (def="")
  -poolsizefile CHAR     Name of the Pool Size file => activate PoolSeq mode (def="")
  -outprefix    CHAR     Prefix used for the output files (def="")

II) Modeling Options:
  -omegafile    CHAR     Name of the omega matrix file => inactivate estim. of omega (def="")
  -rho          INT      Rho parameter of the Wishart prior on omega (def=1)
  -setpibetapar Inactivate estimation of the Pi beta priors parameters
  -betaprior    FLOAT2   Pi Beta prior parameters (if -setpibetapar) (def=1.0 1.0)
  -minbeta     FLOAT    Lower beta coef. for the grid (def=-0.3)
  -maxbeta     FLOAT    Upper beta coef. for the grid (def= 0.3)

I.1) IS covariate mode (default covariate mode):
  -nbetagrid    INT      Number of grid points (IS covariate mode) (def=201)

I.2) MCMC covariate mode:
  -covmcmc      Activate mcmc covariate mode (deactivate estim. of omega)
  -auxmodel     Activate Auxiliary variable mode to estimate BF
  -isingbeta    FLOAT    Beta (so-called inverse temperature) of the Ising model (def=0.0)
  -auxPbetaprior FLOAT2  auxiliary P Beta prior parameters (def=0.02 1.98)

III) MCMC Options:
  -nthreads     INT      Number of threads (def=1)
  -nval         INT      Number of post-burnin and thinned samples to generate (def=1000)
  -thin         INT      Size of the thinning (record one every thin post-burnin sample) (def=25)
  -burnin       INT      Burn-in length (def=5000)
  -npilot       INT      Number of pilot runs (to adjust proposal distributions) (def=20)
  -pilotlength  INT      Pilot run length (def=1000)
  -accinf       FLOAT    Lower target acceptance rate bound (def=0.25)
  -accsup       FLOAT    Upper target acceptance rate bound (def=0.40)
  -adjrate      FLOAT    Adjustment factor (def=1.25)
  -d0pi         FLOAT    Initial delta for the pi all. freq. proposal (def=0.5)
  -upalphaalt   Alternative update of the pij
  -d0pij        FLOAT    Initial delta for the pij all. freq. proposal (alt. update) (def=0.05)
  -d0yij        INT      Initial delta for the yij all. count (PoolSeq mode) (def=1)
  -seed         INT      Random Number Generator seed (def=5001)
```

In this menu, each option is followed by the kind of argument (if any) required (e.g., INT for integer, FLOAT for real, FLOAT2 for a pair of real number space separated), a brief description of its function, and the default value.

In the following, we detailed all the options of BAYPASS:

`-help`

This option prints out the help menu (see above). Note that this option is dominating all the other options, i.e. if `-help` is used in conjunction

with any other option of the program, the help menu is displayed. No argument is required for this option.

#### `-npop`

This option (mandatory) gives the number of population considered in the data set (half the number of column in the genotype data file). The required argument must be an integer (INT). (e.g., `-npop 12` if 12 populations are studied).

#### `-gfile`

This option (mandatory) gives the name of the genotyping input file. See [2.3.1](#) for a description of the corresponding input file format. The required argument must be character chain (name of the input file) without space (e.g., `-gfile data.geno` if the input file is named "data.geno").

#### `-efile`

This option gives the name of the covariate input file. See [2.3.3](#) for a description of the corresponding input file format. The required argument must be character chain (name of the input file) without space (e.g., `-gfile data.env` if the input file is named "data.env").

#### `-scalecov`

This option allows to perform scaling of each covariable in the covariate input file (See [2.3.3](#)). No argument is required for this option. If activated, an output file named "covariate.std" containing the scaled covariables is produced.

#### `-poolsizefile`

This option gives the name of the input file containing the haploid sample size of each population. See [2.3.2](#) for a description of the corresponding input file format. The required argument must be character chain (name of the input file) with no space (e.g., `-poolsizefile data.poolsize` if the input file is named "data.poolsize"). Note that this option automatically activates the Pool-Seq mode, i.e., the PoolSeq version of the different models are considered (as represented in Figures [1A2](#), [B2](#) and [C2](#)).

#### `-outprefix`



This option allows to add a prefix to all the output files. The required argument must be a character chain without space. For instance, if using `-outprefix ana1`, the name of all the output files will begin by "ana1\_". By default, no prefix is added.

#### `-omegafile`

This option gives the name of the input file for the population covariance matrix ( $\mathbf{\Omega}$  in 3.1.1 and Figure 1). See 2.3.4 for a description of the corresponding input file format. The required argument must be character chain (name of the input file) with no space (e.g., `-omegafile matrix.dat` if the input file is named `matrix.dat`). This option inactivates the estimation of  $\mathbf{\Omega}$  and is mandatory in the covariate models involving estimation of the regression coefficients via MCMC i.e. the standard model (see 3.1.2) with the `-covmcmc` option and the auxiliary variable model (see 3.1.3).

#### `-rho`

This option allows to specify the value of  $\rho$  for the Inverse-Wishart prior of  $\mathbf{\Omega}$  (see Figure 1 and 3.1.1). The required argument must be a positive integer. By default, `-rho 1` (i.e.,  $\rho = 1$ ).

#### `-setpibetapar`

This option allows to inactivate estimation of the two (hyper-)parameters  $a_\pi$  and  $b_\pi$  of the prior  $\beta$  distribution for the overall (across population) SNP allele frequencies (see Figure 1 and 3.1.1) (and set them to the values specified with the `-betapiprior` option). No argument is required for this option.

#### `-betapiprior`

This option allows to specify the values of the two (hyper-)parameters  $a_\pi$  and  $b_\pi$  (respectively) of the prior  $\beta$  distribution for the overall (across population) SNP allele frequencies (see Figure 1 and 3.1.1). The required argument must be two positive real numbers. By default `-betapiprior 1.0 1.0` (i.e.,  $a_\pi = b_\pi = 1$ ).

#### `-minbeta`

This option allows to specify the lower bound of the Uniform prior distribution on the regression coefficients (see Figure 1 and 3.1.2). The

required argument must be a real number (lower than `maxbeta` defined below). By default `-minbeta -0.3` (i.e.,  $\beta_{\min} = -0.3$ ).

`-maxbeta`

This option allows to specify the upper bound of the Uniform prior distribution on the regression coefficients (see Figure 1 and 3.1.2). The required argument must be a real number (greater than `minbeta` defined above). By default `-maxbeta 0.3` (i.e.,  $\beta_{\max} = 0.3$ ).

`-nthreads`

This option gives the number of threads to be used for parallel computations. By default, `-nthreads 1` (i.e., a single core is used hence no parallelization).

`-nval`

This option gives the number of post–burn–in (and thinned MCMC) samples recorded from the posterior distributions of the parameters of interest. The required argument must be a positive integer. By default, `-nval 1000` (i.e., 1,000 post burn-in and thinned samples are generated). Note that with default values, the total number of iterations of the MCMC sampler run after the burn-in period is equal to 25,000 (since by default, the thinning rate is equal to 25 , see `thin` option)

`-thin`

This option gives the size of the thinning (i.e., the number of iterations between any two records from the MCMC). The required argument must be a positive integer. By default, `-thin 25` (i.e., the size of the thinning is 25).

`-burnin`

This option gives the length of the burn-in period (i.e., the number of iterations before the first record from the MCMC). The required argument must be a positive integer. By default, `-burnin 5000` (i.e., 5,000 iterations are run during the burn-in period).

`-npilot`

This option gives the number of pilot runs (i.e., the number of runs used to adjust the parameters of the MCMC proposal distributions of parameters updated through a Metropolis-Hastings algorithm). The targeted acceptance rates are defined with the `accinf` and `accsup` options (by default, these are set to 0.25 and 0.40 respectively). The required argument must be a positive integer. By default, `-npilot 20` (i.e., 20 pilot runs are performed).

`-pilotlength`

This option gives the number of iterations of each pilot run (see `npilot` option above). The required argument must be a positive integer. By default, `-pilotlength 1000` (i.e., each pilot run consist of 1,000 iterations).

`-accinf`

This option gives the lower bound of the targeted acceptance rates to adjust the parameters of the MCMC proposal distributions of parameters updated through a Metropolis-Hastings algorithm, during the pilot runs. For instance, in the case of a uniform proposal distribution of the form  $\text{Unif}(x - \delta, x + \delta)$  (where  $x$  represents the current value of the parameter of interest and  $\delta$  specifies the size of the support), if acceptance rates are below this lower bound after a pilot run, then  $\delta$  is increased (e.g., multiplied by a factor defined with the `adjrate` parameter and set to 1.25 by default). The required argument must be a positive real number ( $< 1$  and lower than `accsup` defined below). By default, `-accinf 0.25` (i.e., acceptance rates should be at least equal to 25%).

`-accsup`

This option gives the upper bound of the targeted acceptance rates to adjust the parameters of the MCMC proposal distributions of parameters updated through a Metropolis-Hastings algorithm, during the pilot runs. For instance, in the case of a uniform proposal distribution of the form  $\text{Unif}(x - \delta, x + \delta)$  (where  $x$  represents the current value of the parameter of interest and  $\delta$  specifies the size of the support), if acceptance rates are above this upper bound after a pilot run, then  $\delta$  is decreased (e.g., divided by a factor defined with the `adjrate` parameter and set to 1.25 by default). The required argument must be a positive real number ( $\leq 1$  and greater than `accinf` defined above). By default, `-accinf 0.40` (i.e., acceptance rates should be less than 40%).

### -adjrate

This option gives the factor used to adjust the parameters of the MCMC proposal distributions of parameters updated through a Metropolis-Hastings algorithm, during the pilot runs. For instance, in the case of a uniform proposal distribution of the form  $\text{Unif}(x - \delta, x + \delta)$  (where  $x$  represents the current value of the parameter of interest and  $\delta$  specifies the size of the support), if acceptance rates are below (respectively above) the lower (respectively upper) bound of the targeted regions (as defined above with the `accinf` and `accsup` options) after a pilot run, then  $\delta$  is multiplied (respectively, divided) by this factor. The required argument must be a real number  $> 1$ . By default, `-adjrate 1.25`.

### -d0pi

This option gives the initial value of the  $\delta_\pi$ , which is half the window width from which proposal values of the overall SNP allele frequencies  $\pi_i$  (see Figure 1) are drawn uniformly around the current value  $p_{ij}$  in the Metropolis-Hastings update. The value of  $\delta_\pi$  is eventually adjusted for each locus during the pilot runs (see options `-npilot`, `-pilotlength`, `-accinf`, `-accsup` and `-adjrate`). The required argument must be a positive real number. By default, `-d0pi 0.5` (i.e.,  $\delta_p = 0.5$ ).

### -upalphaalt

This option activates an alternative Metropolis-Hastings algorithm for the population SNP allele frequencies  $\alpha_{ij}$  (see Figure 1). By default, the proposal is the same as the one described by Coop *et al.* (2010) (Appendix A). Briefly, denoting  $\alpha_i$  as the vector of allele frequencies for SNP  $i$  in each population, the vector  $\alpha_i^{*cdt}$  evaluated in a given Metropolis-Hastings update is sampled from the following multivariate Gaussian distribution:  $\alpha_i^{*cdt} \sim MNV(\alpha_i^*, \Gamma \sigma_i^2)$  where  $\Gamma$  is obtained by a Choleski decomposition of the matrix  $\Omega$  (i.e.,  $\Omega = {}^t\Gamma\Gamma$ ). The alternative proposal activated with the `-upalphaalt` option is defined on a SNP by population basis and is a uniform distribution centred on the current values of the parameters (i.e.,  $\alpha_{ij}^{*cdt} \sim \text{Unif}(\alpha_{ij}^* - \lambda_{ij}^\alpha, \alpha_{ij}^* + \lambda_{ij}^\alpha)$ ). The algorithm is slower than the default one but may perform better, in particular when sample sizes are heterogeneous across samples. No argument is required for this option.

### -d0pij

This option gives the initial value of the  $\delta_\alpha$  used in the proposal distribution of the population SNP allele frequencies  $\alpha_{ij}$  in the Metropolis–Hastings updates. Following the notations used above (see `-upalphaalt` option),  $\delta_\alpha = \sigma_i^2$  for the default algorithm and  $\delta_\alpha = \lambda_{ij}^\alpha$  for the alternative algorithm. The value of  $\delta_\alpha$  is eventually adjusted for each locus (and population, in the case of the alternative algorithm) during the pilot runs (see options `-npilot`, `-pilotlength`, `-accinf`, `-accsup` and `-adjrate`). The required argument must be a positive real number. By default, `-d0pij 0.05` (i.e.,  $\sigma_i^2 = 0.05$  for the default algorithm and  $\lambda_{ij}^\alpha = 0.05$  for the alternative algorithm).

`-d0yij`

This option gives, in the Pool–Seq mode, the initial value of the  $\delta_y$  used in the proposal distribution of the population SNP allele count  $Y_{ij}$  in the Metropolis–Hastings updates. The value of  $\delta_y$  is eventually adjusted for each locus and each population during the pilot runs (see options `-npilot`, `-pilotlength`, `-accinf`, `-accsup` and `-adjrate`). The required argument must be a positive integer number lower than the haploid pool sizes. By default, `-d0yij 1` (i.e.,  $\delta_y = 1$ ).

`-seed`

This option gives the initial seed of the (pseudo-) Random Number Generator. The required argument must be a positive integer number. By default, `-seed 5001`.

### 3.3 Format of the output files

While running, BAYPASS printed on the console several information regarding the execution of the analysis (these might be redirected in a log file using the `> log.file` unix syntax). At the end of the analysis BAYPASS produces several output files which might varied according to the options considered (see 3.2). The name of these different output files might be preceded by a prefix as defined with the `outprefix` options (see 3.2).

In the following, we detailed all the output files that may be generated by BAYPASS:

- `(outprefix_)summary_pij.out` (default mode; e.g., for allele count data) or `(outprefix_)summary_yij_pij.out` (Pool–Seq mode; e.g., for read count data)

These files contain for each locus (MRK column) within each population (POP column), the mean (M\_P column) and the standard deviation (SD\_P column) of the posterior distribution of the  $\alpha_{ij}^*$  parameter (see Figure 1) that is closely related to the frequency of the reference allele  $\alpha_{ij} = 1 \wedge (0 \vee \alpha_{ij}^*)$  except that its support is on the real line (hence possible values  $< 0$  or  $> 1$ ). It also contains the posterior mean (M\_Pstd column) and the standard deviation (SD\_Pstd column) of the standardized allele frequency  $\alpha^{\text{std}}$  ( $\alpha^{\text{std}} = \mathbf{\Gamma}^{-1}\alpha^*$ ). The final adjusted values of the  $\delta_\alpha$ 's, the parameters of the MCMC proposal distributions for the SNP allele frequencies (see `-upalphaa1t` option) are also reported in these files (DELTA\_P column) together with the post-burn-in final acceptance rates (ACC\_P column). Note that in the default Metropolis-Hastings algorithm, the  $\alpha_{ij}^*$  are updated for each SNP as a vector of allele frequencies across all populations. Hence, the  $\delta_\alpha$  and the acceptance rates have same values across all the populations for a given SNP. In the Pool-Seq mode (i.e., in the `(outprefix_)summary_yij_pij.out` file), the columns M\_Y, SD\_Y, DELTA\_Y and ACC\_Y similarly report the posterior mean, the posterior standard-deviation, the  $\delta_Y$  of the corresponding proposal distributions and the post-burn-in final acceptance rates for allele counts of each SNP within each population.

- `(outprefix_)summary_pi_xtx.out`

This file contains for each locus (MRK column), the mean (M\_P column) and the standard deviation (SD\_P column) of the posterior distribution of the (across populations) frequency  $\pi_i$  of the SNP reference allele (see Figure 1). The final adjusted values of the  $\delta_\pi$ 's, the parameters of the MCMC proposal distributions are also reported in these files (DELTA\_P column) together with the post-burn-in final acceptance rates (ACC\_P column). In addition, this files contains for each SNP, the posterior mean (M\_XtX column) and standard deviation (SD\_XtX column) of the XtX statistics introduced by Günther and Coop (2013) to identify outlier loci in genome-scan of adaptive differentiation (see 3.1.1).

- `(outprefix_)summary_lda_omega.out` and `(outprefix_)mat_omega.out`

The `(outprefix_)summary_lda_omega.out` file contains the posterior means and posterior standard deviations of each element of the  $n_{pop} \times n_{pop}$  scaled population allele frequencies covariance matrix  $\mathbf{\Omega}$  (M\_omega\_ij and SD\_omega\_ij columns respectively) as described in Figure 1 (see also 3.1.1), and its inverse  $\mathbf{\Lambda} = \mathbf{\Omega}^{-1}$  (M\_lambda\_ij and SD\_lambda\_ij columns respectively).

The `(outprefix_)mat_omega.out` file contains the posterior means of the elements of  $\mathbf{\Omega}$  in a matrix format. Note that this file is in the format required by the `-omegafile` option of BAYPASS.

- `(outprefix_)summary_beta_params.out` (generated with the `-estpibetapar` option)

This file contains the posterior mean (Mean column) and standard deviation (SD column) of the two parameters ( $a_\pi$  and  $b_\pi$ ) of the Beta prior distribution assumed for the (across populations) frequencies of the SNP reference allele (see Figure 1).

- `(outprefix_)summary_betai_reg.out`

This file is only produced in the BAYENV-like mode, i.e. the standard covariate mode (see Figure 1B and 3.1.2) where the estimation of the Bayes Factor (column BF(dB)) in dB units (i.e.,  $10 \times \log_{10}(\text{BF})$ ) measuring the support of the association of each SNP with each population covariable and the corresponding regression coefficients  $\beta_i$  (column Beta\_is) are done via an Importance Sampling algorithm (Coop *et al.*, 2010). In addition, the file also contains the empirical Bayesian P-value (eBPis) in the  $\log_{10}$  scale i.e.  $\text{eBPis} = -\log_{10}(1 - 2 | 0.5 - \Phi(\widehat{\mu}_\beta / \widehat{\sigma}_\beta) |)$  (where  $\Phi(x)$  represents the cumulative distribution function for the standard normal distribution) and thus allowing to evaluate the support in favour of a non-null regression coefficient (e.g.,  $\text{eBPis} > 3$ ). This file contains for each covariable and each SNP, the posterior mean and standard deviation of the Pearson correlation coefficient (columns M\_Pearson and SD\_Pearson respectively) between the scaled allele frequencies  $\widetilde{\boldsymbol{\alpha}}_i^* = \left\{ \frac{\alpha_{ij}^* - \pi_i}{\sqrt{\pi(1-\pi_i)}} \right\}_{(1..J)}$  and the given covariable after rotation of both vectors by  $\mathbf{\Gamma}^{-1}$  (see Günther and Coop, 2013) where  $\mathbf{\Gamma}$  is obtained by a Choleski decomposition of the matrix  $\mathbf{\Omega}$  (i.e.,  $\mathbf{\Omega} = {}^t\mathbf{\Gamma}\mathbf{\Gamma}$ ).

- `(outprefix_)summary_betai.out` (generated with the `-covmcmc` option)

This file is produced in place of the `(outprefix_)summary_betai_reg.out` described above when the `-covmcmc` option is activated (see 3.1.2).

Under the standard model (default), the file contains for each SNP, the posterior mean  $\widehat{\mu}_\beta$  (M\_Beta column) and standard deviation  $\widehat{\sigma}_\beta$  (SD\_Beta column) of the regression coefficient  $\beta_i$  together with the adjusted  $\delta_\beta$  parameter (DeltaB column) of the proposal distribution

and the post-burn-in acceptance rate (**AccRateB** column). In addition, the file also contains an approximated Bayesian P-value in the log10 scale (**eBPmc**) measured as  $eBPmc = -\log_{10}(1 - 2 | 0.5 - \Phi(\widehat{\mu}_\beta/\widehat{\sigma}_\beta) |)$  (where  $\Phi(x)$  represents the cumulative distribution function for the standard normal distribution) and thus allowing to evaluate the support in favour of a non-null regression coefficient (e.g.,  $eBPmc > 3$ ).

Under the model with auxiliary variables (**-auxmodel** option, see 3.1.3), the file contains for each SNP, the posterior mean (**M\_Beta** column) and standard deviation (**SD\_Beta** column) of the regression coefficient  $\beta_i$ , the posterior mean of the auxiliary variable (**M\_Delta** column) and the estimate of the Bayes Factor (column **BF(dB)**) in dB units (i.e.,  $10 \times \log_{10}(\text{BF})$ ) for comparison of models with ( $\beta_i \neq 0$ ) and without ( $\beta_i = 0$ ) correlation with the given covariable.

- *(outprefix\_)summary\_Pdelta.out* (covariate model with auxiliary variable, i.e. **-auxmodel** option, see 3.1.3)

This file contains the posterior mean (**M\_P** column) and standard deviation (**SD\_P** column) of the parameter  $P$  (see Figure 1C and 3.1.3) corresponding to the overall proportion of SNP associated to the corresponding covariable.

- *(outprefix\_)covariate.std* (generated by the **scalecov** option)

This file contains the scaled covariables.

- *(outprefix\_)DIC.out*

This files contains the average deviance (**bar(D)** column), the effective number of parameters of the models (**pD** column) and the Deviance Information Criterion (**DIC** column) as defined in Spiegelhalter *et al.* (2002) and that might be relevant for model comparison purposes. In addition, the logarithm of the pseudo-marginal likelihood of the model is also provided (**LPML** column).

## 4 Miscellaneous R functions

The **baypass\_utils.R** file in the **utils** directory contains three R functions (R Core Team, 2015) (**simulate.baypass()**; **fmd.dist()** and **geno2YN()**) that may be helpful to interpret some of the results obtained with **BAYPASS**. To use this functions, one may simply need to source the corresponding



files and ensure that the packages `mvtnorm` (Genz *et al.*, 2015) and `geigen` (Hasselmann, 2015) have been installed. In addition, although not required by these functions, the packages `corrplot` (Wei, 2013) and `ape` (Paradis *et al.*, 2004) might prove useful for the visualisation of the  $\Omega$  matrix (see 5).

## 4.1 The R function *simulate.baypass()*

### 4.1.1 Description

The R function `simulate.baypass()` allows to simulate either allele or read count data under the core inference model (Figure 1A) and possibly under the STD covariate model (Figure 1B). It produces several objects and output files in a format directly appropriate for analyses with BAYPASS and BAYENV2<sup>8</sup>. In practice, this function is useful to generate POD for calibration of the XtX differentiation measure (or any other measures). More broadly, because the  $\Omega$  matrix captures the demographic history of the populations, this function might also be viewed as an efficient simulator of population genetics data.

### 4.1.2 Usage

```
simulate.baypass(omega.mat, nsnp=1000, beta.coef=NA, beta.pi=c(1,1),
                 pop.trait=0, sample.size=100, pi.maf=0.05, suffix="sim",
                 remove.fixed.loci=FALSE, coverage=NA)
```

### 4.1.3 Arguments (in alphabetic order)

- `beta.pi` (def=c(1,1))

A two elements vector giving the parameters  $a_\pi$  and  $b_\pi$  respectively, for the Beta distribution of the  $\pi_i$  ("ancestral") allele frequencies.

- `beta.coef` (def=NA; required for simulation under the STD covariate model)

A vector giving the values of the regression coefficients ( $\beta_i$  in Figure 1) for the simulated associated SNPs (the number of the simulated associated SNPs is equal to the dimension of the vector).

- `coverage` (def=NA; required to activate simulation of read count data)

---

<sup>8</sup> For analyses with BAYENV2, make sure fixed loci have been removed, i.e., `remove.fixed.loci=TRUE`

Either a single value or a matrix giving the total read counts. In the latter case, the vector of total read counts for each simulated SNP are sampled with replacement from the row of the matrix. It is thus mandatory that the number of columns of the matrix equals the number of populations, but no restriction are set for the number of rows. For instance, if the matrix has only one row, all the SNPs will have the same read counts within a given population.

- `omega.mat` (always required)

A positive definite and symmetric matrix of rank `npop` corresponding to the covariance matrix of population allele frequencies ( $\Omega$  in Figure 1)

- `nsnp` (def=1000)

A single number giving the number of (neutral) SNPs to simulate.

- `pi.maf` (def=0.05)

A single value giving the MAF threshold on the simulated  $\pi_i$  ("ancestral") allele frequencies. In the simulation procedure, the  $p_{i_i}$ 's are sampled from the Beta distribution with parameters as specified in the `beta.pi` argument. For a given SNP  $i$ , if  $p_{i_i} < \text{pi.maf}$  (resp.  $p_{i_i} > 1 - \text{pi.maf}$ ) then  $p_{i_i}$  is set equal to `pi.maf` (resp. `1 - pi.maf`). Setting `pi.maf = 0` inactivates MAF filtering.

- `pop.trait` (def=0; required for simulation under the STD covariate model)

A vector of length `npop` giving each population-specific covariable values (ordering of the population is assumed to be the same as in the `omega.mat` matrix). By default all values are set to 0 (meaning the associated SNPs behave neutrally irrespective of their values at the regression coefficients).

- `remove.fixed.loci` (def=FALSE)

A logical indicating whether or not fixed loci (in the observed simulated data) should be discarded.

- `sample.size` (def=100)

If simulating allele count data, either a single value or a matrix giving the total allele counts (twice the number of genotyped individuals in

diploid species). In the latter case, the vector of total allele counts for each simulated SNP are sampled with replacement from the row of the matrix. It is thus mandatory that the number of columns of the matrix equals the number of populations, but no restriction are set for the number of rows. For instance, if the matrix has only one row, all the SNPs will have the same allele counts within a given population.

If simulating read count data, either a single value or a vector of length `npop` giving the pool haploid sample sizes for each population.

- `suffix` (def="sim")

A character string giving the suffix of the output files generated by the function.

#### 4.1.4 Values

The function produces an object which is a list with the following components:

- `omega.sim`

A matrix corresponding to the one used for simulation and declared with the `omega.mat` argument

- `alpha.sim`

A matrix with dimension `nsnp` rows and `npop` columns giving the simulated (unbounded) allele frequencies for each simulated SNP within each population (i.e.,  $\alpha_{ij}^*$  in Figure 1).

- `pi.sim`

A vector of length `nsnp` giving the simulated  $\pi_i$  "ancestral" allele frequencies for each simulated SNP.

- `N.sim`

A matrix with dimension `nsnp` rows and `npop` columns giving the total allele counts for each simulated SNP within each population.

- `Y.sim`

A matrix with dimension `nsnp` rows and `npop` columns giving the allele counts for the reference allele for each simulated SNP within each population.

- `N.pool` (read count data only)

A matrix with dimension `nsnp` rows and `npop` columns giving the total read counts for each simulated SNP within each population.

- `Y.pool` (read count data only)

A matrix with dimension `nsnp` rows and `npop` columns giving the read counts for the reference allele for each simulated SNP within each population.

- `betacoeff.sim` (simulation under the STD covariate model only)

A vector of length `nsnp` giving the regression coefficients of each SNP used to simulate the data.

In addition, the following output files are printed out (the extension `.suffix` is as defined in the `suffix` argument):

- `G.suffix`

The allele count data file in BAYPASS format (see 2.3).

- `Gpool.suffix` (when simulating read count data)

The read count data file in BAYPASS format (see 2.3).

- `bayenv_freq.suffix`

The allele count data file in BAYENV2 format (see 8).

- `bayenv_freq_pool.suffix` (when simulating read count data)

The read count data file in BAYENV2 format (see 8).

- `alpha.suffix`

The simulated (unbounded) allele frequencies (`nsnp` rows and `npop` columns) for each simulated SNP within each population (i.e.,  $\alpha_{ij}^*$  in Figure 1)

- `pi.suffix`

The simulated  $\pi_i$  "ancestral" allele frequencies for each simulated SNP.

- `betacoef.suffix` (when simulating under the STD covariate model)

The regression coefficients of each SNP used to simulate the data.

- `pheno.suffix` (when simulating under the STD covariate model)

The covariate data file in BAYPASS format (see 2.3).

- `poolsize.suffix` (when simulating read count data)

The haploid pool size data file in BAYPASS format (see 2.3).

#### 4.1.5 Examples

```
#source the baypass R functions (check PATH)
source("utils/baypass_utils.R")
#load the bovine covariance matrix
om.bta <- as.matrix(read.table("examples/omega.bta"))
#simulate allele count data for 1000 SNPs
simu.res <- simulate.baypass(omega.mat=om.bta)

#simulate allele count data for 1000 neutral SNPs and
#100 associated SNPs with varying regression coefficients
simu.res <- simulate.baypass(omega.mat=om.bta,beta.coef=runif(100,-0.2,0.2),
                             pop.trait=rnorm(18))

#simulate read count data for 1000 SNPs
simu.res <- simulate.baypass(omega.mat=om.bta,coverage=50)
```

## 4.2 The R function *fmd.dist()*

### 4.2.1 Description

This function computes the metric proposed by (Förstner and Moonen, 2003) to evaluate the distance between two covariance matrices (FMD distance).

### 4.2.2 Usage

```
fmd.dist(mat1,mat2)
```

### 4.2.3 Arguments

- `mat1` and `mat2`

Two positive-definite (symmetric) matrices

#### 4.2.4 Values

The function returns a numeric corresponding to the FMD distance between the two matrices.

#### 4.2.5 Example

```
#source the baypass R functions (check PATH)
source("utils/baypass_utils.R")
#load the bovine covariance matrix
om.bta <- as.matrix(read.table("examples/omega.bta"))
#create a dummy diagonal covariance matrix
#this might be obtained from a star-shaped phylogeny with
#branch length (Fst) equal to 0.1
star.bta<-diag(0.1,nrow(om.bta))

#compute the fmd.dist between the two matrices
fmd.dist(om.bta,star.bta)
```

### 4.3 The R function *geno2YN()*

#### 4.3.1 Description

This function reads the allele (or read) count data file in the BAYPASS format and extract both the counts for the reference allele and total counts.

#### 4.3.2 Usage

```
geno2YN(genofile)
```

#### 4.3.3 Arguments

- **genofile**

A character string giving the name of the allele (or read) count data file in the BAYPASS format.

#### 4.3.4 Values

The function returns two matrices:

- **genofile**

A character string giving the name of the allele (or read) count data file in the BAYPASS format.

The function produces an object which is a list containing the two following matrices:

- YY

A matrix with nsnp rows and npop columns containing allele or read counts for the reference allele.

- NN

A matrix with nsnp rows and npop columns containing the total allele or read counts.

### 4.3.5 Example

```
#source the baypass R functions (check PATH)
source("utils/baypass_utils.R")
#load the bovine BTA 14 data
counts.obj <- geno2YN("examples/geno.bta14")
```

## 5 Worked Examples

For illustration purposes, different types of analyses of the example files (see 2.3) are detailed step by step. Users might try to reproduce the corresponding examples using the files included in the `example` directory.

### 5.1 Cattle allele count data

#### 5.1.1 Analysis under the core model mode: MCMC is run under the core model

The following command allows to analyse the data under the core model (this should take ca. 4 min with the `i_baypass` executable and ca. 7 min with the `g_baypass` executable):

```
baypass -npop 18 -gfile geno.bta14 -outprefix anacore
```

To visualize the results, one may open an R session and proceed as follows:

```
require(corrplot) ; require(ape)
#source the baypass R functions (check PATH)
source("utils/baypass_utils.R")

#upload estimate of omega
omega=as.matrix(read.table("anacore_mat_omega.out"))
pop.names=c("AUB","TAR","MON","GAS","BLO","MAN","MAR","LMS","ABO",
            "VOS","CHA","PRP","HOL","JER","NOR","BRU","SAL","BPN")
```

```

dimnames(omega)=list(pop.names,pop.names)

#Compute and visualize the correlation matrix
cor.mat=cov2cor(omega)
corrplot(cor.mat,method="color",mar=c(2,1,2,2)+0.1,
          main=expression("Correlation map based on"~hat(Omega)))

#Visualize the correlation matrix as hierarchical clustering tree
bta14.tree=as.phylo(hclust(as.dist(1-cor.mat**2)))
plot(bta14.tree,type="p",
      main=expression("Hier. clust. tree based on"~hat(Omega)~("d[ij]*=1-"*rho[ij]*"))))

#Compare the estimate of omega based on the whole genome
#and based on the BTA14 SNPs

wg.omega <- as.matrix(read.table("examples/omega.bta")) #check the PATH
plot(wg.omega,omega) ; abline(a=0,b=1)
fmd.dist(wg.omega,omega)

#Estimates of the XtX differentiation measures
anacore.snp.res=read.table("anacore_summary_pi_xtx.out",h=T)
plot(anacore.snp.res$M_XtX)

```

One may further wish to calibrate the XtX estimates using a POD sample. For instance, to produce a (small) POD sample with 1,000 SNPs (continuing the R example above) using the `simulate.bypass()` function (see 4):

```

#get estimates (post. mean) of both the a_pi and b_pi parameters of
#the Pi Beta distribution
pi.beta.coef=read.table("anacore_summary_beta_params.out",h=T)$Mean
#upload the original data to obtain total allele count
bta14.data<-geno2YN("geno.bta14")
#Create the POD
simu.bta<-simulate.bypass(omega.mat=omega,nsnp=1000,sample.size=bta14.data$NN,
                           beta.pi=pi.beta.coef,pi.maf=0,suffix="btapods")

```

Then, one may analyse the newly created POD (data file named "G.btapods" in the example) giving another prefix for the output files:

```
baypass -npop 18 -gfile G.btapods -outprefix anapod
```

Continuing the above example in R, calibration of the XtX and visualisation of the results might be carried out as follows:

```

#####
#Sanity Check: Compare POD and original data estimates
#####
#get estimate of omega from the POD analysis
pod.omega=as.matrix(read.table("anapod_mat_omega.out"))
plot(pod.omega,omega) ; abline(a=0,b=1)
fmd.dist(pod.omega,omega)

#get estimates (post. mean) of both the a_pi and b_pi parameters of
#the Pi Beta distribution from the POD analysis
pod.pi.beta.coef=read.table("anapod_summary_beta_params.out",h=T)$Mean
plot(pod.pi.beta.coef,pi.beta.coef) ; abline(a=0,b=1)

```



```
#####
#XtX calibration
#####
#get the pod XtX
pod.txt=read.table("anapod_summary_pi_xtx.out",h=T)$M_XtX
#compute the 1% threshold
pod.thresh=quantile(pod.txt,probs=0.99)
#add the thresh to the actual XtX plot
plot(anacore.snp.res$M_XtX)
abline(h=pod.thresh,lty=2)
```

### 5.1.2 Analysis under the IS covariate mode: MCMC is run under the core model

This analysis allows to perform association study (under the STD covariate model) by estimating for each SNP the Bayes Factor, the empirical Bayesian P-value (and the underlying regression coefficient) using an Importance Sampling algorithm (Gautier, 2015). As a consequence, the MCMC samples of the parameters of interest are obtained by running the core model as above (5.1.1). Hence, if covariables are available, one may rather used this mode as a default mode. The example below corresponds to an association analysis with the SMS covariable measured on the 18 French cattle breeds (Gautier, 2015).

```
baypass -npop 18 -gfile geno.bta14 -efile bta.pc1 -outprefix anacovis
```

Providing the same seed and the same options have been used, one may verify that exactly the same estimates for  $\Omega$  (e.g., files `anacovis_mat_omega.out` and `anacore_mat_omega.out`) and other parameters in common are obtained than under the previous analysis (5.1.1). Continuing the above example in R, one may plot the Importance Sampling estimates of the Bayes Factor, the empirical Bayesian P-value and the underlying regression coefficient as follows:

```
covis.snp.res=read.table("anacovis_summary_betai_reg.out",h=T)
graphics.off()
layout(matrix(1:3,3,1))
plot(covis.snp.res$BF.dB.,xlab="SNP",ylab="BFis (in dB)")
plot(covis.snp.res$eBPis,xlab="SNP",ylab="eBPis")
plot(covis.snp.res$Beta_is,xlab="SNP",ylab=expression(beta~"coefficient"))
```

Recall that in the example only a subset of SNPs mapping to BTA14 are considered. To improve precision in this example, one may rather provide the program with a more accurate estimate of the matrix  $\Omega$  as obtained from the original study on the complete data sets (with 40 times as many SNPs):

```
baypass -npop 18 -gfile geno.bta14 -efile bta.pc1 \
        -omegafile omega.bta -outprefix anacovis2
```

The resulting Importance Sampling estimates of the Bayes Factor, the empirical Bayesian P-value and the underlying regression coefficient <sup>9</sup> might be plotted as follows:

```
covis2.snp.res=read.table("anacovis2_summary_betai_reg.out",h=T)
graphics.off()
layout(matrix(1:3,3,1))
plot(covis2.snp.res$BF.dB.,xlab="SNP",ylab="BFis (in dB)")
plot(covis2.snp.res$eBPis,xlab="SNP",ylab="eBPis")
plot(covis2.snp.res$Beta_is,xlab="SNP",ylab=expression(beta~"coefficient"))
```

### 5.1.3 Analysis under the MCMC covariate mode: MCMC is run under the STD model

This analysis allows to perform association study under the STD covariate model by estimating the empirical Bayesian P-value and the underlying regression coefficient using parameters values sampled from MCMC run under the STD model (Gautier, 2015). Although one may also estimate  $\Omega$  under the STD model, this option has been inactivated in BAYPASS. As a consequence, an estimate of  $\Omega$  (e.g., as obtained by a first analysis under the core model of IS covariate mode) must be provided. The example below corresponds to an association analysis with the SMS covariable measured on the 18 French cattle breeds (Gautier, 2015).

```
baypass -npop 18 -gfile geno.bta14 -efile bta.pc1 \
        -covmcmc -omegafile omega.bta -outprefix anacovmcmc
```

The resulting estimates of the empirical Bayesian P-values, the underlying regression coefficients (posterior mean) and the corrected XtX might be plotted as follows:

```
covmcmc.snp.res=read.table("anacovmcmc_summary_betai.out",h=T)
covmcmc.snp.xtx=read.table("anacovmcmc_summary_pi_xtx.out",h=T)$M_XtX
graphics.off()
layout(matrix(1:3,3,1))
plot(covmcmc.snp.res$eBPmc,xlab="SNP",ylab="eBPmc")
plot(covmcmc.snp.res$M_Beta,xlab="SNP",ylab=expression(beta~"coefficient"))
plot(covmcmc.snp.xtx,xlab="SNP",ylab="XtX corrected for SMS")
```

<sup>9</sup>Note that one may carry out a calibration of these different measures (as detailed for the XtX in 5.1.1) by analysing a POD together with the covariables:

```
baypass -npop 18 -gfile G.pods -efile bta.pc1 -omegafile omega.bta -outprefix podcovis
```

#### 5.1.4 Analysis under the AUX covariate mode: MCMC is run under the AUX model

This analysis allows to perform association study under the AUX covariate model by estimating the Bayes Factor (and the underlying regression coefficient) using parameters values sampled from MCMC run under the AUX model (Gautier, 2015). Although one may also estimate  $\Omega$  under the AUX model, this option has been inactivated in BAYPASS. As a consequence, an estimate of  $\Omega$  (e.g., as obtained by a first analysis under the core model of IS covariate mode) must be provided. The example below corresponds to an association analysis with the SMS covariable measured on the 18 French cattle breeds (Gautier, 2015).

```
baypass -npop 18 -gfile geno.bta14 -efile bta.pci \  
-auxmodel -omegafile omega.bta -outprefix anacovaux
```

The resulting estimates of the Bayes Factor, the underlying regression coefficients (posterior mean) and the corrected XtX might be plotted as follows:

```
covaux.snp.res=read.table("anacovaux_summary_betai.out",h=T)  
covaux.snp.txt=read.table("anacovaux_summary_pi_txt.out",h=T)$M_XtX  
graphics.off()  
layout(matrix(1:3,3,1))  
plot(covaux.snp.res$BF.dB.,xlab="SNP",ylab="BFmc (in dB)")  
plot(covaux.snp.res$M_Beta,xlab="SNP",ylab=expression(beta~"coefficient"))  
plot(covaux.snp.txt,xlab="SNP",ylab="XtX corrected for SMS")
```

One may further introduce spatial dependency among the SNPs by setting  $b_{is} = 1$  in the Ising prior (Figure 1C) to refine the associated region:

```
baypass -npop 18 -gfile geno.bta14 -efile bta.pci -auxmodel \  
-isingbeta 1.0 -omegafile omega.bta -outprefix anacovauxisb1
```

The resulting estimates of the posterior mean of the each auxiliary variable  $\delta_i$  under both models (AUX model with no SNP spatial dependency and AUX model Bayes Factor, the underlying regression coefficients (posterior mean) and the corrected XtX might be plotted as follows:

```
covauxisb1.snp.res=read.table("anacovauxisb1_summary_betai.out",h=T)  
graphics.off()  
layout(matrix(1:2,2,1))  
plot(covaux.snp.res$M_Delta,xlab="SNP",ylab=expression(delta[i]),main="AUX model")  
plot(covauxisb1.snp.res$M_Delta,xlab="SNP",ylab=expression(delta[i]),  
main="AUX model with isb=1")
```

## 5.2 Littorina Pool–Seq read count data

The Littorina Pool–Seq data may be analysed in a similar fashion as the cattle data above except that one needs to specify the (haploid pool) size file using the `-poolsizefile` option to activate the Pool–Seq mode. Because the haploid pool sizes are relatively large ( $= 100$ ) in the example, one may also increase the initial  $\delta$  of the  $y_{ij}$  proposal distribution (as a rule of thumbs, one may set it to a fifth of the minimum pool size). Here is an example of command to run BAYPASS under the IS covariate mode (MCMC run under the core model):

```
i_baypass -npop 12 -gfile lsa.geno -efile lsa.ecotype \  
          -poolsizefile lsa.poolsize -d0yij 20 -outprefix lsacovis
```

## 6 Credits

BAYPASS makes use of several functions and subroutines that were previously developed by other authors. These include:

- the Fortran code for the multiple streams MT19937 Mersenne–Twister (parallel) Random Number Generator was adapted from the subroutines available in the `mt_stream_f90-1.11.tar.gz` program written by Ken-Ichi Ishikawa and available under the New BSD License<sup>10</sup> at [http://theo.phys.sci.hiroshima-u.ac.jp/~ishikawa/PRNG/mt\\_stream\\_f90-1.11/](http://theo.phys.sci.hiroshima-u.ac.jp/~ishikawa/PRNG/mt_stream_f90-1.11/)).
- Various functions and subroutines for random number generations that were adapted from the Alan Miller Fortran module `random.f90` available at: <http://jblevins.org/mirror/amiller/> available under the GNU GPL license.
- the Wishart sampler utilities derived from the fortran `wishart` library written by John Burkhardt and available at [http://people.sc.fsu.edu/~%20jburkardt/f\\_src/wishart/wishart.html](http://people.sc.fsu.edu/~%20jburkardt/f_src/wishart/wishart.html) under the GNU GPL license.
- the `kracken(3f)` Fortran module developed by John S. Urban to parse command line arguments (available at <http://home.comcast.net/~urbanjost/LIBRARY/libCLI/arguments/krackenhelp.html>) under the GNU GPL license.

---

<sup>10</sup>[http://theo.phys.sci.hiroshima-u.ac.jp/~ishikawa/PRNG/mt\\_stream\\_f90-1.11/LICENSE](http://theo.phys.sci.hiroshima-u.ac.jp/~ishikawa/PRNG/mt_stream_f90-1.11/LICENSE)

I also wish to thank Renaud Vitalis for providing the L<sup>A</sup>T<sub>E</sub>Xtemplate for this manual and Andrew Beckerman for reporting bugs and advices that helped to improve the program.

## 7 Copyright

BAYPASS is a free software under the GPL- and BSD-compatible CeCILL-B licence (see [http://www.cecill.info/licences/Licence\\_CeCILL-B\\_V1-en.html](http://www.cecill.info/licences/Licence_CeCILL-B_V1-en.html)), and © INRA.

## 8 Contact

If you have any question, please feel free to contact [me](#). However, I strongly recommend you read carefully this manual first.

## Bibliography

- Coop, G., D. Witonsky, A. D. Rienzo, and J. K. Pritchard, 2010 Using environmental correlations to identify loci underlying local adaptation. *Genetics* 185: 1411–1423.
- Duforet-Frebourg, N., E. Bazin, and M. G. B. Blum, 2014 Genome scans for detecting footprints of local adaptation using a bayesian factor model. *Mol Biol Evol* 31: 2483–2495.
- Förstner, W., and B. Moonen, 2003 A metric for covariance matrices. In *Geodesy-The Challenge of the 3rd Millennium*. Springer Berlin Heidelberg, 299–309.
- Gautier, M., 2015 Genome-wide scan for adaptive divergence and association with population-specific covariates. *Genetics* 201: 1555–1579.
- Genz, A., F. Bretz, T. Miwa, X. Mi, F. Leisch, *et al.*, 2015 *mvtnorm: Multivariate Normal and t Distributions*. R package version 1.0-3.
- Günther, T., and G. Coop, 2013 Robust identification of local adaptation from allele frequencies. *Genetics* 195: 205–220.
- Hasselmann, B., 2015 *geigen: Calculate Generalized Eigenvalues of a Matrix Pair*. R package 1.5.
- Nicholson, G., A. V. Smith, F. Jonsson, O. Gustafsson, K. Stefansson, *et al.*, 2002 Assessing population differentiation and isolation from single-nucleotide polymorphism data. *J Roy Stat Soc B* 64: 695–715.
- Paradis, E., J. Claude, and K. Strimmer, 2004 Ape: Analyses of phylogenetics and evolution in r language. *Bioinformatics* 20: 289–290.
- Pickrell, J. K., and J. K. Pritchard, 2012 Inference of population splits and mixtures from genome-wide allele frequency data. *PLoS Genet* 8: e1002967.
- R Core Team, 2015 *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Riebler, A., L. Held, and W. Stephan, 2008 Bayesian variable selection for detecting adaptive genomic differences among populations. *Genetics* 178: 1817–1829.

- Spiegelhalter, D. J., N. G. Best, B. P. Carlin, and A. v. d. Linde, 2002  
Bayesian measures of model complexity and fit. *Journal of the Royal  
Statistical Society. Series B (Statistical Methodology)* 64: 583–639.
- Wei, T., 2013 *corrplot: Visualization of a correlation matrix*. R package  
version 0.73.
- Westram, A. M., J. Galindo, M. A. Rosenblad, J. W. Grahame, M. Panova,  
*et al.*, 2014 Do the same genes underlie parallel phenotypic divergence in  
different *littorina saxatilis* populations? *Mol Ecol* 23: 4603–4616.

## Appendix A Comparisons of the computational efficiency of the different version of BayPass

For illustration purposes, Table 1 gives computational times obtained when running the same example (Pool-Seq) analysis for the Littorina read count example data set (12 pools, 2,500 SNPs) considered in paragraph 5.2. The multi-threaded version (`baypass 2.0`) is compared to the initial non-parallel version (`baypass 1.01`) with the same options as in paragraph 5.2 (except for the number of threads). For both versions, both a `gfortran` and `ifort` compiled binaries were considered. All comparisons were done under the same standard computer running under Linux Debian 7.1.

Version	nthreads	gfortran binary	ifort binary
BAYPASS 1.01	1	38 min 59 s	12 min 11 s
BAYPASS 2.0	1 (default)	39 min 48 s	14 min 54 s
BAYPASS 2.0	4 ( <code>-nthreads 4</code> )	11 min 34 s	5 min 40 s

**Table 1:** Comparisons of the computational efficiency of different version of BAYPASS for the analysis of the Littorina Pool-Seq read count example data set (12 pools, 2,500 SNPs) considered in paragraph 5.2.

Note that when running on a single core (and for very short runs as well, not shown), as expected, the `baypass 1.01` is slightly more efficient computationally than the `baypass 2.0`. Also, the `ifort` binaries always outperform the `gfortran` binaries.