

GADMA user manual

Corresponding to version 1.0.0

Ekaterina Noskova

Contents

1	Introduction	3
1.1	Getting help	3
2	Parameter file	3
3	Hands on	3
4	Input data	7
4.1	VCF data format	8
4.2	Frequency spectrum file format	8
4.3	SNP data format	8
4.4	Unlinked SNP's, AIC and CLAIC	9
5	Specifying a scheme	10
6	Specifying a model	10
6.1	Dynamics of population size change	10
6.2	Specifying the structure of the model	11
6.2.1	Initial structure	12
6.2.2	Final structure	12
6.3	Specifying a model in details	13
7	Several repeats and parallel computing	14
8	Output	15
8.1	Stdout and log file	15
8.2	Model representation	15
8.3	Output directory content	17
8.4	Generated code of models	18
9	Plotting model	18
9.1	Time units on model's plot	19
9.2	Plotting after GADMA finished	19

10 Usefull options	19
10.1 Time for generation	19
10.2 Theta0	20
10.2.1 Estimating Theta0	20
10.2.2 Changing Theta0	20
10.3 Examples of Theta0 and time for generation	21
10.4 Relative to N_A parameters	21
10.5 Disable migrations	21
10.6 Upper bound for time split	21
10.7 Resume launch	22
10.7.1 Only models	22
11 Extra parameter file	22
12 Units	22
13 Decreasing the number of model's parameters	22
14 Confidence intervals	23
14.1 Run local search on bootstrapped data	23
14.2 Get Confidence Intervals from table	24
15 Installation	25
15.1 Dependencies	25
15.2 $\partial a \partial i$	25
15.3 <i>moments</i>	26
15.4 Pillow	26
15.5 GADMA	26
15.5.1 Verifying installation	27
A Example of the parameter file	28
B Example of generated code	34

1 Introduction

Welcome to GADMA!

GADMA implements methods for automatic inferring the joint demographic history of multiple populations from genetic data.

GADMA is based on two open source packages: *∂a∂i* developed by Ryan Gutenkunst [<https://bitbucket.org/gutenkunstlab/dadi/>] and *moments* developed by Simon Gravel [<https://bitbucket.org/simongravel/moments/>].

In contrast, GADMA is a command-line tool. It presents a series of launches of the genetic algorithm and infers demographic history from the Allele Frequency Spectrum of multiple populations (up to three).

1.1 Getting help

If you have any problems or questions, please don't hesitate to contact me by email: ekaterina.e.noskova@gmail.com. I will be glad to answer your questions.

2 Parameter file

A parameter file is a simple text file (created in a text editor, such as Notepad), which contains a list of parameters and variables with their assigned values. Create a parameter file that reflect your particular parameters. Hereinafter, as an example, the name of your parameter file will be defined as `param_file`.

3 Hands on

Test case

GADMA has a test case for a simple demographic model for one population: just a constant size of 10000 individuals in population. To run a test case:

```
$ gadma --test
```

Example 2

Suppose we have SNP data for two populations. Data is in *∂a∂i*'s SNP file format in the `snp_data.txt`. Suppose we want to get all output in some `gadm_output` directory:

```
$ gadma -i snp_data.txt -o gadm_output
```

Example 3

We didn't specify AFS size or labels for populations, they are taken automatically from the input file. We can see a parameter file of our run in the `gadm_output/param_file`.

```
# gadma_output/param_file
...
Population labels : pop_name_1, pop_name_2
Projections: 18, 20
...
```

But we know that spectrum should be 20×20 ! To specify size of AFS we need to create a parameter file and set Projections:

```
# param_file
Projections : 20,20
```

Order of populations can be changed as:

```
# param_file
Projections : 20,20
Population labels : pop_name_2, pop_name_1
```

If we want to rename populations, we should change names on `snp_data.txt` file.

Now assume we want to get the simplest demographic model as fast as we can. We will tell GADMA that we need no other dynamics of population sizes except sudden (constant) population size change and that we want to use *moments* library.

```
# param_file
Projections : 20,20
Population labels : pop_name_2, pop_name_1
Only sudden : True
Use moments or dadi : moments
```

To run GADMA we need to specify `-p/--params` command-line option in cmd:

```
$ gadma -i snp_data.txt -o gadma_output -p params_file
```

Example 4

Consider some AFS file `fs_data.fs`. There is a spectrum for three populations: YRI, CEU, CHB. However axes are mixed up: CHB, YRI, CEU. To run GADMA we should order them from most ancient to most recent:

```
# param_file
Population labels : YRI, CEU, CHB
```

We want to allow exponential growth (it is the default behaviour) and have some extra changes in size of the ancient population. To do so we should specify **Initial structure**. It is list of three numbers: first — number of time intervals before first split (we want here 2); second — number of time periods between first and second split events (at least 1); third — number of time periods after second split.

```
# param_file
Population labels : YRI, CEU, CHB
Initial structure : 2,1,1
```

Also we can put information about input file and output directory to our parameter file:

```
# param_file
Input file : fs_data.fs
Output directory : gadma_output
Population labels : YRI, CEU, CHB
Initial structure : 2,1,1
```

Now we can run GADMA in the following way:

```
$ gadma -p params
```

Example 5

We have our GADMA launch interrupted for some reasons. We want to resume it:

```
$ gadma --resume gadma_output
```

where `gadma_output` is output directory of previous run. We can find resumed run in `gadma_output_resumed`

Example 6

Our launch was finished, we used $\partial a \partial i$ with a default grid size which GADMA determines automatically if it is not specified by user. We found out that it would be better to find some models using greater number of grid points in $\partial a \partial i$ scheme, but we want to take final models from previous run:

```
# param_file
Pts : 40, 50, 60 #Greater value of grid size than it was
```

And run GADMA:

```
$ gadma --resume gadma_output --only_models -p params
--only_models tell GADMA to take from gadma_output final models only.
```

There is another way to do the same:

```
# param_file
Resume from : gadma_output
Only models : True
Pts : 40, 50, 60 #Greater value of grid size than it was
```

And run GADMA in the following way:

```
$ gadma -p params
```

Example 7

We can add a custom model using a parameter `Custom filename` in the parameter file:

```
# param_file
Custom filename : YRI_CEU_demographic_model.py
```

Our custom file need to contain a function with a fixed name `model_func` (see Appendix A). For example:

```
# YRI_CEU_demographic_model.py
def model_func((nu1F, nu2B, nu2F, m, Tp, T), (n1,n2), pts)
    xx = yy = dadi.Numerics.default_grid(pts)

    phi = dadi.PhiManip.phi_1D(xx)
    phi = dadi.Integration.one_pop(phi, xx, Tp, nu=nu1F)

    phi = dadi.PhiManip.phi_1D_to_2D(xx, phi)
    nu2_func = lambda t: nu2B*(nu2F/nu2B)**(t/T)
    phi = dadi.Integration.two_pops(phi, xx, T, nu1=nu1F,
                                    nu2=nu2_func, m12=m, m21=m)

    sfs = dadi.Spectrum.from_phi(phi, (n1,n2), (xx,yy))
    return sfs
```

In addition, we can easily specify values for lower and upper bounds through a parameter file. Let's set lower and upper bounds for the model we defined above:

```
# param_file
Lower bounds : 1e-2, 1e-2, 1e-2, 0, 0, 0
Upper bounds : 100, 100, 100, 10, 3, 3
```

Example 8

Also, we can get the values of lower/upper bounds, both, or none of them in the parameter file automatically. For this, each identifier in the parameter file must be declared through a parameter `Parameter` identifiers. Below is an identifier list:

```
# param_file
# An identifier list:
# T - time
# N - size of population
# m - migration
# s - split event, proportion in which population size
# is divided to form two new populations.
```

For example, we set a lower bound for the model we defined above (see Example 7) and we want to get an upper bound automatically.

```
# param_file
Lower bound : 1e-2, 1e-2, 1e-2, 0, 0, 0
Upper bound : None
```

```
Parameter identifiers : n, n, n, m, t, t
```

Example YRI, CEU

GADMA has an example of the parameter file `example_params` that is presented at the

end of the manual. To run GADMA with this parameters one should just run from the GADMA's home directory:

```
$ gadma -p example_params
```

4 Input data

GADMA supports several types of input data, which are familiar to anyone who has used *∂a∂i* or *moments* in the past:

- Frequency spectrum file format (.fs or .sfs),
- SNP data format (.txt).

Input file can be specified to GADMA in two ways:

1. Through command-line option `-i/--input`:

```
$ gadma -i fs_file.fs -o out_dir
```

or

```
$ gadma --input snp_file.txt -o out_dir
```

2. Use a parameter Input file in the parameter file:

```
# param_file

# Input file path
Input file : fs_file.fs
...
```

Extra information about input AFS can also be put in the parameter file. For example, AFS can be projected to a smaller size with **Projections** option, populations can be named or their order can be changed with **Population labels** option. If parameter file does not contain some options, they are automatically pulled out from the input file.

```
# param_file

# Input file path
Input file : fs_file.fs

# (New) size of the AFS
Projections : 20,20

# Labels of populations
Population labels : CEU, YRI
...
```

GADMA can be launched with a parameter file in the following way:

```
$ gadma -p params_file -o out_dir
```

4.1 VCF data format

To convert a VCF (.vcf) file into a SFS (.sfs) file use <https://github.com/isaacovercast/easySFS>.

4.2 Frequency spectrum file format

Each file begins with any number of comment lines beginning with **#**. The first non-comment line contains P integers giving the dimensions of the FS array, where P is the number of populations represented. For a FS representing data from 4x4x2 samples, this would be 5 5 3. (Each dimension is one larger than the number of samples, because the number of observations can range, for example, from 0 to 4 if there are 4 samples, for a total of 5 possibilities.) On the same line, the string **folded** or **unfolded** denoting whether or not the stored FS is folded.

The actual data is stored in a single line listing all the FS elements separated by spaces, in the order `fs[0,0,0] fs[0,0,1] fs[0,0,2]... fs[0,1,0] fs[0,1,1]...`. This is followed by a single line giving the elements of the mask in the same order as the data, with 1 indicating masked and 0 indicating unmasked.

4.3 SNP data format

Human	Chimp	Allele1	YRI	CEU	Allele2	YRI	CEU	Gene	Position
ACG	ATG	C	29	24	T	1	0	abcb1	289
CCT	CCT	C	29	23	G	3	2	abcb1	345

Listing 1: Example of SNP file format

The data file begins with any number of comment lines that begin with **#**. The first parsed line is a column header line. Whitespace is used to separate entries within the table, so no spaces are allowed within any entry. Individual rows may be commented out using **#**.

The first column contains the in-group reference sequence at that SNP, including the flanking bases. If the flanking bases are unknown, they can be denoted by **-**. The header label is arbitrary.

The second column contains the aligned outgroup reference sequence at that SNP, including the flanking bases. Unknown entries can be denoted by **-**. The header label is arbitrary.

The third column gives the first segregating allele. The column header must be exactly **Allele1**.

Then follows an arbitrary number of columns, one for each population, each giving the number of times Allele1 was observed in that population. The header for each column should be the population identifier.

The next column gives the second segregating allele. The column header must be exactly **Allele2**.

Then follows one column for each population, each giving the number of times Allele2 was observed in that population. The header for each column should be the population identifier, and the columns should be in the same order as for the Allele1 entries.

Then follows an arbitrary number of columns which will be concatenated with _ to assign a label for each SNP.

The **Allele1** and **Allele2** headers must be exactly those values because the number of columns between those two is used to infer the number of populations in the file.

4.4 Unlinked SNP's, AIC and CLAIC

By default, SNP's that were used to build AFS are considered to be linked. In this case it is possible to compare demographic models with different number of parameters by Composite Likelihood Akaike Information Criterion (*CLAIC*). This procedure can be necessary as a model with a large number of parameters will be better able to find parameter values corresponding to the observed data than a model with a smaller number of parameters, but at the same time it will correspond less to reality, for example, due to data errors. It is called overfitting and we do not want it to happen.

Actually, *CLAIC* is modification of usual Akaike Information Criterion (*AIC*), but *AIC* can be used only for AFS built from unlinked SNP's. The smaller the *AIC* or *CLAIC* score is, the better the model fits the observed data.

It is possible to inform GADMA about linkage of SNP's and unlock the usage of *AIC*:

```
# param_file

# Inform if SNP's are not linked
Linked SNP's : False
...
```

If SNP's are linked and *CLAIC* should be evaluated (by default it is not), then the bootstrapped data should be set. In order to receive reliable correct bootstrapped data, the bootstrap should be performed on the original SNP data over the unlinked regions of genome. For example, in case of exome data one could make it over genes. Then when bootstrap is done, it is required to set the directory with it in the parameters file for *CLAIC* evaluation:

```
# param_file

# Inform if SNP's are not linked
Linked SNP's : True

# Tell where bootstrapped data is located
Directory with bootstrap: /home/dadi/examples/YRI_CEU/bootstraps/
...
```

This kind of bootstrap is called block-bootstrap and it is very important if one want to do some model selections for data with linked SNPs. **Please, be careful if it is your case.**

5 Specifying a scheme

GADMA uses either $\partial a \partial i$ or *moments* to simulate expected AFS from the demographic model. $\partial a \partial i$ is used by default. To use $\partial a \partial i$ it is recommended to check the value of the `Pts` option in the `params_file`. `Pts` is sequence of three numbers, each of which is equal to the number of points in grid size. The greater the numbers are, the more accurate $\partial a \partial i$ numerical solution of partial differential equation is. However, finding such a solution take more time. By default, GADMA takes `Pts : n, n + 10, n + 20`, where `n` — is the largest sample size among populations of interest.

moments library does not need `Pts` to be specified. To change $\partial a \partial i$ to *moments* scheme, specify option in the parameters file:

```
# param file
...
Use moments or dadi : moments
...
```

6 Specifying a model

6.1 Dynamics of population size change

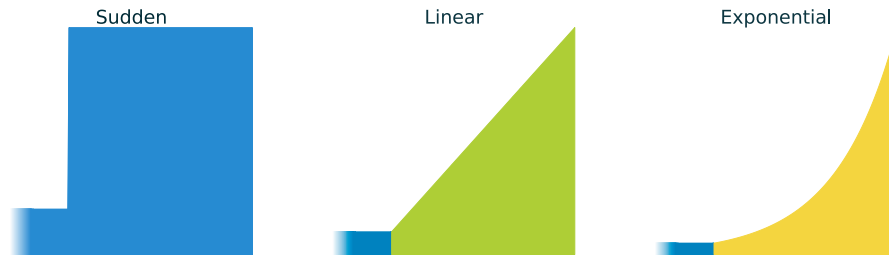


Figure 1: Three main demographic dynamics of population size change.

In GADMA the size of population can be changed due to one of three dynamics: sudden change, linear change and exponential change of the effective population size (Figure 1).

In order to infer a demographic model with sudden changes of populations sizes only, option `Only sudden` in the parameters file should be set to `True`:

```
# param_file
...
Only sudden : True
...
```

By default, this option is set to `False` and dynamics are found like other parameters of the demographic model.

6.2 Specifying the structure of the model

GADMA infers a demographic model from an AFS with nothing required from the user, except simple information, that determines how much the detailed model is required — **the structure of the model**.

Assume a division of one population into two new subpopulations and a fixed temporal order of the populations: from the most anciently to the most recently formed population.

We can divide time of our model into split events and time intervals, during which a certain dynamics of change of effective size is maintained for each population and migration rates are constant. The number of split events is one less than the number of populations under consideration. Now we can define the concept of the structure of the demographic model:

Structure of the demographic model is:

- the number of time intervals in case of one population;
- the number of time intervals as those that occur before and after a single split, in the case of two populations;
- the number of time intervals prior to the first split, those between the first and second split, and the ones after the second split, in the case of three populations.

For example, we can divide into time of the model on the Figure 2 to four time intervals: T_1 , T_2 , T_3 and T_4 , and two populations' splits: S_1 and S_2 . The structure of this model is 2, 1, 1, because two intervals (T_1 as T_2) before first split S_1 , one interval (T_3) between first and second splits and one interval (T_4) after second split S_2 .

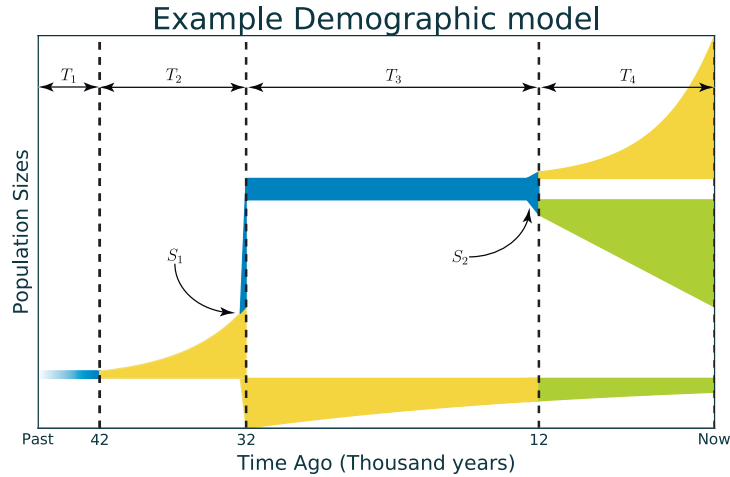


Figure 2: Example of representation of demographic model. Time is on the axis of abscissa and population size is on the axis of ordinates. The structure of that model is (2,1,1). The colours reflect different demographic dynamics.

6.2.1 Initial structure

To specify the structure of the inferred model one should set **Initial structure** in the parameter file:

```
# param_file
...
Initial structure : 2
...
```

or

```
# param_file
...
Initial structure : 2,1
...
```

or

```
# param_file
...
Initial structure : 2,1,1
...
```

By default the simplest structure is used (1 or 1,1 or 1,1,1).

6.2.2 Final structure

It is also possible to start with a simpler structure in order to get to a more complex one. To do so one should specify option **Final structure** in the parameter file. For example:

```
# param_file
Input file : some_2d_fs.fs

Initial structure : 1,1
Final structure : 2,1
...
```

Within this parameter file GADMA will find parameters for demographic model with 1,1 structure, then increase structure 2,1 and find parameters for the model with this structure. Parameters identify within a more simple structure (in this case it is 1,1) are used further to define the parameters of a more complex structure (2,1).

Note: Structure increases random, so if one specifies initial structure to 1,1 and final to 2,2, it is not guaranteed to final optimal parameters for demographic models with structures between 1,1 and 2,2, i.e. intermediate state can be either 1,2 or 2,1.

Recommendation: Use scheme with a more complex structure, as it produce more stable solutions.

Recommendation: Choose recommended values for model's structure. Final structure should differ by one element from initial structure, for example, 1,1 and 2,1; 1,2,1 and 2,2,1.

6.3 Specifying a model in details

It is also possible to use Genetic algorithm from GADMA to a proposed model that is defined as Python function using *dadi* or *moments*. It is the way that *dadi* and *moments* work with demographic models inference. To understand how to specify a model like that one can read manuals to the packages.

For example, consider simple bottleneck model for one population: at time $T_F + T_B$ in the past, an equilibrium population goes through a bottleneck of depth nu_B , recovering to relative size nu_F :

```
def bottleneck(params, ns, pts):
    nuB, nuF, TB, TF = params
    xx = dadi.Numerics.default_grid(pts)

    phi = dadi.PhiManip.phi_1D(xx)
    phi = dadi.Integration.one_pop(phi, xx, TB, nuB)
    phi = dadi.Integration.one_pop(phi, xx, TF, nuF)

    fs = dadi.Spectrum.from_phi(phi, ns, (xx,))
    return fs
```

To run optimization from GADMA one needs to run optimization function just like in *dadi* and *moments*:

```
# Import GADMA's optimization function:
import gadma

# Specify input data and its parameters:
data = dadi.Spectrum.from_file("fs_filename.fs")
ns = data.sample_sizes # size of AFS
pts = [40,50,60] # grid size for dadi

# Wrap our bottleneck function:
func_ex = dadi.Numerics.make_extrap_log_func(bottleneck)

# Specify upper and lower bounds for parameters:
upper_bound = [100, 100, 3, 3]
lower_bound = [1e-2, 1e-2, 0, 0]

# Run optimizations:
run the optimization
# Beginning GADMA optimization
popt = gadma.Inference.optimize_ga(4, data, func_ex, pts_1,
                                   p_ids = ['n', 'n', 't', 't'],
                                   lower_bound=lower_bound,
                                   upper_bound=upper_bound)
# Beginning local optimization from dadi
```

```
popt = dadi.Inference.optimize_log(popt, data, func_ex, pts_l,
                                   lower_bound=lower_bound,
                                   upper_bound=upper_bound,
                                   verbose=len(p0), maxiter=3)
```

```
print('Found parameters: {0}'.format(popt))
```

Note: As GADMA optimization is a global search algorithm, no initial parameters `p0` are set in `gadma.Inference.optimize_ga` function. However, it is possible to specify `p0`:

```
# Initial parameters can be set too:
```

```
p0 = [0.01, 1.5, 0.2, 0.2]
```

```
# Beginning GADMA optimization
```

```
popt = gadma.Inference.optimize_ga(4, data, func_ex, pts_l, p0=p0,
                                   p_ids = ['n', 'n', 't', 't'],
                                   lower_bound=lower_bound,
                                   upper_bound=upper_bound)
```

If one wants to find other parameters for `gadma.Inference.optimize_ga` function:

```
>>> import gadma
>>> help(gadma.Inference.optimize_ga)
```

7 Several repeats and parallel computing

By default, GADMA runs optimization that uses Genetic algorithm once. However, it is recommended to run optimization several times and choose the best model among the resieved ones. Option **Number of repeats** in the parameter file tells GADMA how many times an optimization should be executed. Moreover, there is another option **Number of processes** that allows GADMA to run all the processes in parallel.

Note: GADMA parallelize exactly several runs, not every of them. So if one asks to repeat optimization twice and parallelize in more than 2 processes, only two processes will be used eventually.

Recommendation: **Number of processes** should be less than **Number of repeats** and it will be better if it is aliquot to the **Number of repeats**.

Note: **Number of processes** shouldn't be greater than the number of kernels on one's computer. Otherwise there will be no sense in parallelization.

```
# param_file
```

```
...
```

```
Number of repeats : 6
```

```
Number of processes : 2 # or 3 or 6
```

```
...
```

8 Output

GADMA puts all files to the directory that user set through `-o/--output` command-line option:

```
$ gadma -o output_dir -i input_fs.fs
```

or through `Output directory` option in the parameter file:

```
# param_file
...
Output directory : output_dir
...
```

8.1 Stdout and log file

GADMA prints its progress about every minute in stdout and in `output_dir/GADMA.log` file:

```
[hhh:mm:ss]
All best logLL models:
GA number      logLL      AIC      Model

All best AIC models:
GA number      logLL      AIC      Model

--Best model by log likelihood--
Log likelihood:      Best logLL
with AIC score:      AIC_score
Model:  representation of best logLL model

--Best model by AIC score--
Log likelihood:      logLL
with AIC score:      Best AIC score
Model:  representation of best AIC model
```

Note: One can set `Silence` option in the parameter file to `True` to disable output in stdout, file `output_dir/GADMA.log` will still have it.

8.2 Model representation

Every model is printed as a line of parameters. All model parameters, except mutation rates, have precision equaled to 2.

Consider designations:

- T — time,
- P — percent of split,

- N_i — size of population number i ,
- d_i — dynamic of changing of the size of population number i ,
- m_{ij} — mutation rate from population i to population j .

Dynamic of population size change has numerical values:

- $d_i = 0$ — sudden change of size for population number i ;
- $d_i = 1$ — linear change of size for population number i ;
- $d_i = 2$ — exponential change of size for population number i .

Model is printed as sequence of time intervals and splits that are represented in the following way:

- First period (NA — size of ancestry population):
[NA]
- Split:
 - If there are one population before split event, so there will be two populations after it:
[$P\%$, [$N1$, $N2$]]
 - If there are two population before split event, so last formed population is divided and there will be three populations after it:
[$P\%$, [$N1$, $N2$, $N3$]]
- Usual time period:
 - If there is one population:
[T , [$N1$], [$d1$]]
 - If there are two populations:
[T , [$N1$, $N2$], [$d1$, $d2$], [[$None$, $m12$], [$m21$, $None$]]]
 - If there are three populations:
[T , [$N1$, $N2$, $N3$], [$d1$, $d2$, $d3$],
[[$None$, $m12$, $m13$], [$m21$, $None$, $m23$], [$m31$, $m32$, $None$]]]

At the end of the string, that corresponds to the model, there is an information about model's ancestry in the genetic algorithm:

- c — for model, that is child of crossover,
- r — if it was formed random way,
- m — if it was mutated,

- **f** — final model of genetic algorithm.

Note: **m** is added as many times as model is mutated.

Example of the demographic model for two populations:

```
[ [144.38] ][ 16.00%, [23.10, 121.28] ][ 375.77, [143.31, 30.07],
                                     [2, 2], [[None, 3.33e-03][7.53e-04, None]] ]
```

8.3 Output directory content

For every repeat of Genetic algorithm GADMA creates a new folder in output directory with corresponding number.

In every folder there is `GADMA_GA.log`, where every iteration of algorithm is saved, folders `pictures` and `python_code`.

When genetic algorithm finishes GADMA saves picture and python code of received model in the corresponding folder.

When all GA's are executed, the codes are saved in the root directory.

```
- <output_dir>
  - 1
      GADMA_GA.log
      - pictures
      - python_code
          - dadi
          - moments
      result_model.png
      result_model_code.py
  - 2
      GADMA_GA.log
      - pictures
      - python_code
          - dadi
          - moments
      result_model.png
      result_model_code.py
params
extra_params
GADMA.log
best_logLL_model.png
best_logLL_model_dadi_code.py
best_logLL_model_moments_code.py
best_aic_model.png
best_aic_model_dadi_code.py
best_aic_model_moments_code.py
```

8.4 Generated code of models

By default, GADMA generates python code only for final models both for *∂a∂i* and *moments*. However, it can do it every N iteration of genetic algorithm. In this case option `draw_models_every N iteration` should be set in the parameter file. GADMA saves files with code to the `output_dir/<GA_number>/python_code` directory. Both *∂a∂i* and *moments* code are generated and saved in different folders there.

Each code contains function of the model, which takes values of the parameters as input, and strings that load observed AFS, simulates expected AFS from the model's function and calculates log likelihood between two AFS'. The result log likelihood is printed to stdout. For the *moments* code picture is also drawn.

All codes can be run in the following way:

```
$ python file_with_code.py
```

Example of generated code one can find at the end of this manual in the corresponding section.

9 Plotting model

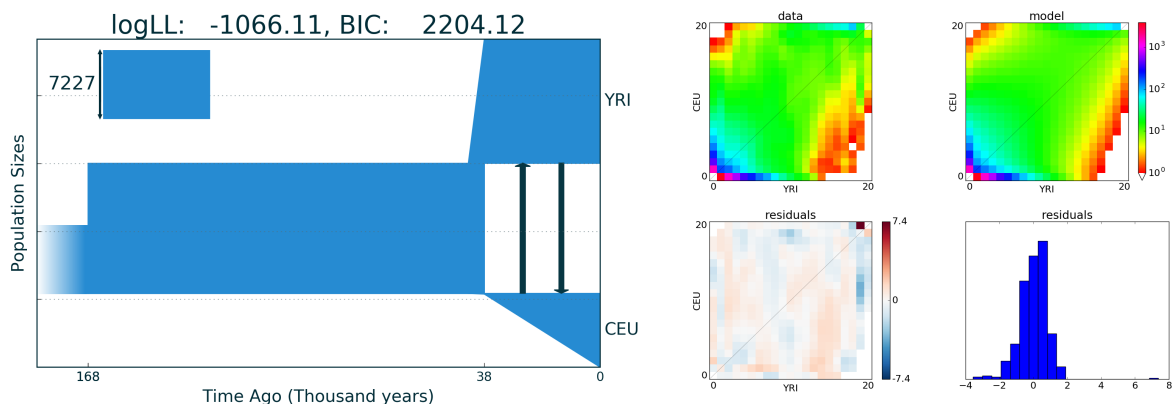


Figure 3: Example of demographic model plot that GADMA draws during run.

GADMA always draws final best models during genetic algorithms search for best solutions.

However, models can be drawn during the pipeline every N 's iteration of each genetic algorithm. N is equal to the value of `Draw models every N iteration` option in the parameter file. So to enable drawing one should set this option in file.

Recommendation: Don't draw models very often, because changes are usually not very significant and drawing takes some time, so optimization will be slower.

Note: One can disable drawing by setting `Draw models every N iteration : 0` in the parameter file. This is also default behaviour.

Models are drawn with *moments* library, so it should be installed if one wants to have pictures. In the top left corner there is a size of ancestry population. Other parameters one can find in string representation of the model in the log files.

9.1 Time units on model's plot

Time on the demographic model's plot can be drawn in units of years, thousand years or in genetic units. By default choice depends on the `Time for generation` option in the parameter file: if it is set to some value (in years) then time will be shown in years, otherwise it will be in genetic units.

But, of course, it is possible to tell GADMA which units are preferable. For example, if one wants time to be in thousands of years on the pictures, as it is a big value in years:

```
# param_file
...
Units of time in drawing : thousand years
...
```

9.2 Plotting after GADMA finished

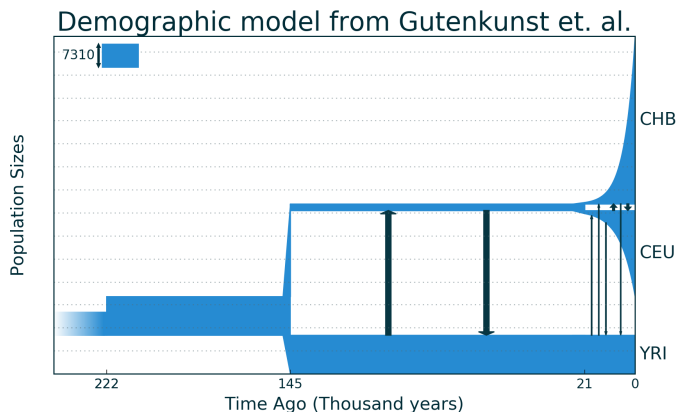


Figure 4: Example of model's plot that was drawn with generated python script.

Sometimes final pictures aren't satisfying or one didn't draw for some reasons (didn't want to install *moments*, wanted fast launch) it is possible to plot demographic model (and only it) again.

To do it one should run corresponding generated python script for the model. For example, final model can be drawn again in the following way:

```
$ python best_logll_model_moments_code.py
```

Note: Again it's possible only if one has *moments* installed.

Note: One can change code inside the file and drawn again if picture isn't satisfying.

10 Usefull options

10.1 Time for generation

Option `Time for generation` in the parameter file is corresponding to the time per one generation in Wright-Fisher model. It is responsible basically for one thing: time on the

model's plots. If it is specified, then time on the pictures will be converted from genetic units by scaling with this value. Otherwise, if it is not set, time will be shown in genetic units.

Note: If `Time for generation` is specified, it should be consistent with another option: `Theta0`, which is described in the corresponding section.

10.2 Theta0

`Theta0` is equal to the expected number of mutations that occur in one chromosome in one generation in the infinite-sited model. GADMA can scale all values of demographic model parameters due to known value of `Theta0`. However, it is not always possible to find it. There is a way to solve this problem: one can set `Theta0` to `None` or just not specify it at all, so GADMA will take it as 1.0 and after launch one can scale result values due to found `Theta0` (how to do it is described in the corresponding section).

10.2.1 Estimating Theta0

If μ — is the neutral mutation rate per site per generation and L — is the length of the sequence, then:

$$\theta_0 = 4 \cdot \mu \cdot L$$

Note: L is the effective sequence length, which accounts for losses in alignment and missed calls.

Note: μ should be estimated based on generation time. One can leave `Time per generation` option in the parameter file unspecified (then time on the model's plots will be in the genetic units), but one should recalculate μ !

For example (Gutenkunst et al, 2009 [1]):

We estimate the neutral mutation rate μ using the divergence between human and chimp. Comparing aligned sequences in our data, we estimate the divergence to be 1.13%. Assuming a divergence time of 6 million years and a mean generation time for human and chimp over this interval of 25 years, we have

$$\mu = 0.0113 \cdot 25 / (2 \cdot 6 \cdot 10^6) = 2.35 \times 10^{-8} \text{ per generation.}$$

10.2.2 Changing Theta0

If GADMA was launched with one `Theta0` and now one wants to use another or if it was launched with default `Theta0` = 1.0 and now one has estimated its real value, model's parameters can be simply scaled:

$$\text{Let } a = \frac{\theta_0^{NEW}}{\theta_0^{OLD}},$$

- Size of population $/a$,

- Time $/a$,
- Migration rates $\cdot a$,
- Split percent stay the same.

10.3 Examples of Theta0 and time for generation

The following tables produce different possible values for the demographic model inference for three populations of modern people: YRI, CEU, CHB.

FS filename	Gen. time (years)	μ (per site per gen.)	L (base pair)	Theta0 (per chr. per gen.)
YRI_CEU_CHB.fs	25 [1]	$2.35 \cdot 10^{-8}$ [1]	$4.04 \cdot 10^6$ [1]	0.37976
YRI_CEU_CHB.fs	24 [3]	$2.26 \cdot 10^{-8}$ [1]	$4.04 \cdot 10^6$ [1]	0.36521
YRI_CEU_CHB.fs	29 [2]	$1.44 \cdot 10^{-8}$ [2]	$4.04 \cdot 10^6$ [1]	0.23270
YRI_CEU_CHB.fs	24 [3]	$1.2 \cdot 10^{-8}$ [2]	$4.04 \cdot 10^6$ [1]	0.19392

Table 1: Examples of different values of generation time and its influence on μ and Theta0

In Gutenkunst et al. [1] generation time for human populations was equal to 25 years and mutation rate μ was estimated as $2.35 \cdot 10^{-8}$. If one wants to change time for one generation to 24 years, one needs to scale μ : $\mu/25 \cdot 24 = 2.26 \cdot 10^{-8}$.

In Jouganous et al. [2] generation time was grater — 29 years and mutation rate was equal to $1.44 \cdot 10^{-8}$. To change generation time to 24, one needs to change value of the mutation rate: $\mu_{NEW} = \mu/29 \cdot 24 = 1.2 \cdot 10^{-8}$. Theta0 is calculated then by the formula above.

Note: One can find FS files in `fs_examples` directory.

10.4 Relative to N_A parameters

Sometimes it is more important to see parameters scaled to $N_{ref} = N_A$. To tell GADMA shows models with scaled parameters, option `Relative parameters` should be set to `True`. By default, it is `False`. It is conveniently when Theta0 is unknown.

10.5 Disable migrations

GADMA can to exclude migrations rates from optimization and consider them be equal to zero. In that case all migrations are disabled. One should set option `No migrations` to `True` in the parameter file.

10.6 Upper bound for time split

To limit time of some split one should specify option in the parameter file. Splits are numerated from the most ancient. So split 1 is split that occurred with ancient population and

split 2 is next division of second population (exist only for three populations). There are two appropriate options: `Upper bound of first split` and `Upper bound of second split`.

One should translate time from years into genetic units, therefore divide it by $2 \cdot T_g$, where T_g is time (in years) for one generation. For example, one wants to limit last split with 2000 years. Time for one generation is estimated as 24 years, then one should specify in the parameter file:

```
# param_file
...
Upper bound of second split : 83.333
...
```

10.7 Resume launch

To resume interrupted launch one can use `--resume` command-line option or set `Resume` from in the parameter file. One needs to set output directory of previous run.

If neither `Output directory` or `-o/--output` is not specified, GADMA will continue evaluation in the directory: `<previous_output_dir>_resumed`.

10.7.1 Only models

GADMA can resume launch taking final models only from previous run. This means, that it is not resumption, but run from some initial values. It is useful, for example, when one has to run GADMA with some small grid size for `dadi` and then wants to restart it with greater number of grid points. To do so, one should set command-line option `--only_models` with `--resume` or specify `Only models` option in the parameter file to `True`.

11 Extra parameter file

GADMA take an extra parameter file as input. **However one probably do not need them.** Nevertheless, if one is interested, `extra_params_template` with all options and their descriptions can be found in `gadma` folder.

12 Units

GADMA shows model parameters in genetic units. To scale them from one should multiply migration rates by 2 only. Other units are as usual. In case when option `Relative parameters` is set to `True` one should first rescale from units of N_{ref} : sizes of populations and time — multiply by `NA`, migration rates — divide by `NA`.

13 Decreasing the number of model's parameters

There are several ways to decrease the number of demographic model's parameters in GADMA. The first way as it was discussed above, is to consider only sudden dynamics

of population size changes (`Only sudden : True`). The second way to disable migration by setting `No migrations` to `True` in the parameter file.

The last and more reliable way to decrease the of model's parameters is to use GADMA optimization on the custom model written with *∂a∂i* or *moments* (section Specifying model in details). However, it requires some extra work for user.

14 Confidence intervals

GADMA contains special scripts for Confidence Intervals evaluation. To get CI one will need correctly bootstrapped data (There is discussion in section ??). If SNP's that were used for AFS are unlinked, then usual bootstrap over them should be performed. However, if they are linked then block bootstrap should be used. It is done over the unlinked regions of genome.

When bootstrapped data is ready, one should run two scripts `gadma-run_ls_on_boot_data` and `gadma-get_confidence_intervals` in order to get CI. One can find example here: https://bitbucket.org/noscode/gadma_results/src/master/YRI_CEU/model_1/.

14.1 Run local search on bootstrapped data

First script `gadma-run_ls_on_boot_data` runs local search from known optimum for initial AFS (the one that GADMA found) for each AFS from bootstrap. The usage is following:

```
$ gadma-run_ls_on_boot_data --help
usage: GADMA module for runs of local search on bootstrapped data.
Is needed for calculating confidence intervals.
    [-h] -b <dir> -d <filename> -o <dir> [-j N]
    [--opt log/powell] [-p <filename>]
```

optional arguments:

```
-h, --help                show this help message and exit
-b <dir>, --boots <dir>   Directory where bootstrapped data is
                           located.
-d <filename>, --dem_model <filename>
                           File with demographic model. Should
                           contain 'model_func' or 'generated_model'
                           function. One can put there several extra
                           parameters and they will be taken
                           automatically, otherwise one will need to
                           enter them manually. Such parameters are:
                           1) p0 (or popt) - initial parameters values
                           2) lower_bound - list of lower bounds for
                           parameters values
                           4) upper_bound - list of upper bounds for
                           parameters values
```

```

5) par_labels/param_labels - list of
string names for parameters
6) pts - pts for dadi (if there is no pts
then moments will be run automatically).
-o <dir>, --output <dir>
Output directory.
-j N, --jobs N
Number of threads for parallel run.
--opt log/powell
Local search algorithm, by now it can be:
1) 'log' - Inference.optimize_log
2) 'powell' - Inference.optimize_powell.
-p <filename>, --params <filename>
Filename with parameters, should be valid
python file.
Parameters are presented in -d/--dem_model
option description upper.

```

After the run, there will be pandas table in output directory with two different extensions - `result_table.pkl` and `result_table.csv`. It contains parameters for each bootstrap. At this point of time it is possible to change its units and add new parameters with additional manipulations in Python and then run `gadma-get_confidence_intervals` to get CI.

14.2 Get Confidence Intervals from table

After the `gadma-run_ls_on_boot_data` the result table will be in output directory. One can change its columns due to what parameters should be used for CI. For example, it is possible to translate units from genetic and add such parameter as size of ancestral population (`N_A`). To do it user should write script.

To calculate confidence intervals:

```

$ gadma-get_confidence_intervals --help
usage: GADMA module for calculating confidence
       intervals from the result table of local
       search runs on bootstrapped data.
       [-h] [--log] [--tex] [--acc N] <filename>

```

positional arguments:

```

<filename>  Filename (.csv or .pkl) with result from
             local search runs on bootstrapped data.
             Output of gadma-run_ls_on_boot_data.

```

optional arguments:

```

-h, --help  show this help message and exit
--log       If log then logarithm will be used to
             calculate confidence intervals.
--tex       Tex output.
--acc N     Accuracy of output (dafault: 5).

```


15 Installation

First version of GADMA was for Python2, but current version supports either Python2 or Python3. The last is recommended.

One can try to install GADMA and all its possible dependencies via shell script from GADMA home directory the following way (script will use Python3):

```
$ ./install
```

if you want to install with 2 version of Python:

```
$ ./install --py2
```

If some problems occur please try to install GADMA by steps that are described below.

15.1 Dependencies

GADMA requires the following dependencies:

- Python (3 is recommended)
- NumPy ($\geq 1.2.0$)
- Scipy ($\geq 0.6.0$)
- $\partial a \partial i$ ($\geq 1.7.0$) or *moments* ($\geq 1.0.0$)

To draw demographic models install the following:

- matplotlib ($\geq 0.98.1$)
- Pillow ($\geq 4.2.1$)
- *moments* ($\geq 1.0.0$)

To calculate Confidence Intervals:

- pandas

NumPy and SciPy can be installed via pip:

```
$ python -m pip install numpy
```

```
$ python -m pip install scipy
```

15.2 $\partial a \partial i$

To install $\partial a \partial i$, run the following command in the work directory:

```
$ git clone https://bitbucket.org/gutenkunstlab/dadi/
```

Then go to the dadi directory and install package:

```
$ cd dadi
```

```
$ python setup.py install
```

To check *dadi*'s installation, use the command:

```
$ python
>>> import dadi
```

15.3 *moments*

To install *moments*, go to the work directory and run:

```
$ git clone https://bitbucket.org/simongravel/moments/
```

Cython can be installed using `pip`:

```
$ python -m pip install --upgrade Cython
```

Also *moments* requires `mpmath`, which can be installed via `pip`:

```
$ python -m pip install mpmath
```

Then go to the `moments` directory and install package:

```
$ cd moments
$ python setup.py install
```

To check *moments*'s installation, use the command:

```
$ python
>>> import moments
```

15.4 Pillow

Install Pillow with `pip`:

```
$ python -m pip install Pillow
```

To check Pillow's installation, use the command:

```
$ python
>>> import PIL
```

15.5 GADMA

To download GADMA, go to the work directory and run:

```
$ git clone https://github.com/ctlab/GADMA
```

Then go to the `GADMA` directory and install GADMA:

```
$ cd GADMA
$ python setup.py install
```

Now one can run it like this:

```
$ gadma --help
```

15.5.1 Verifying installation

To verify the installation, run a test:

```
$ gadma --test
```

If the installation was successful, one will find the following information at the end:

```
--Finish pipeline--
```

```
--Test passed correctly--
```

```
Thank you for using GADMA!
```

Appendix A Example of the parameter file

```
# It is a parameter file for GADMA software.

# Lines started from # are ignored.
# Also comments at the end of a line are ignored too.
# Every line contains: Identifier of parameter : value.
# If one wants to change some default parameters, one needs to
# remove # at the beginning of a line and change corresponding
# parameter.

# Output directory to write all GADMA out.
# One needs to set it to not existing or empty directory.
# If it is resumed from other directory and output directory
# isn't set, GADMA will add '_resumed' for previous output
# directory.
Output directory : my_example_run

# One can resume from some other launch of GADMA by setting
# output directory of it to 'Resume from' parameter.
# You can set again new parameters of resumed launch.
#Resume from : another_output_dir

# If one wants to take only models from previous run set this
# flag. Then iterations of GA will start from 0 and values of
# mutation rate and strength will be initial.
# Default : None
Only models : None

# Input file can be sfs file (should end with .fs) or
# file of SNP's in dadi format (should end with .txt).
Input file : fs_examples/YRI_CEU.fs

# 'Population labels' is sequence of population names (the same
# names as in input file)
# If .fs file is in old format then it would rename population
# labels that are absent.
# It is necessary to put them in order from most ancient to less.
# (In case of more than two populations)
```

```

# It is important, because the last of formed populations take
# part in next split.
# For example, if we have YRI - African population,
# CEU - European population and CHB - Chinese population,
# then we can write YRI, CEU, CHB or YRI, CHB, CEU
# (YRI must be at the first place)
# Default: from input file
Population labels : CEU, YRI # we change populations order
                             # (in input file YRI is first)

# Also one can project your spectrum to less size.
# For example, we have 80 individuals in each of three
# populations, then spectrum will be 81x81x81 and one can
# project it to 21x21x21 by set 'Projections' parameter
# to 20, 20, 20.
# Default: from input file
Projections : None # will be 20, 20

# Are SNP's linked or unlinked?
# If they are linked, then Composite Likelihood Akaike
# Information Criterion (CLAIC) will be used to compare models.
# If they are unlinked, then usual Akaike Information Criterion
# (AIC) will be used.
# Default: True
Linked SNP's : True

# Now all main parameters:
#
# Total mutation flux - theta.
# It is equal to:
#  $\theta = 4 * \mu * L$ 
# where  $\mu$  - mutation rate per site per generation and
# L - effective sequenced length, which accounts for losses
# in alignment and missed calls.
# Note: one should estimate  $\mu$  based on generation time.
# Default: 1.0
Theta0 : 0.37976 # the same as in Gutenkunst et al 2009

# Time (years) for one generation. Can be float.

```

```

# It is important for drawing models. If one doesn't want to draw,
# one can pass it.
# Default: 1.0
Time for generation : 25 # the same as in Gutenkunst et al 2009

# Parameters for demographic models:
#
# Use moments or dadi
# Default: dadi
Use moments or dadi : moments

# Use multinom scheme: N_A is not parameter for search,
# it is calculated through optimal_sfs_scaling.
# Multinom scheme decrease number of parameter by one and
# is usually faster, however non multinom scheme usually
# finds better solutions.
# Default: False
Multinom : False

# If you choose to use dadi, please set pts parameter - number
# of points in grid
# Default: Let n = max number of individuals in one population,
# then pts = n, n+10, n+20
#Pts : 20, 30, 40

# Using a custom demographic model.
# Please, specify file with function named 'model_func' in it.
# So file should contain:
# def model_func(params, ns, pts) in case of dadi
# or
# def model_func(params, ns) in case of moments
# Default: None
Custom filename : None

# Now one should specify either bounds or identifications
# of custom model's parameters. All values are in Nref units.
# Lower and upper bounds - lists of numbers.
# List of usual bounds:
# N: 1e-2 - 100
# T: 0 - 5
# m: 0 - 10

```

```

#   s: 0 - 1
#   This bounds will be taken automatically if identifications are set.
#   Default: None
Lower bounds : None
Upper bounds : None
#   An identifier list:
#   T - time
#   N - size of population
#   m - migration
#   s - split event, proportion in which population size
#   is divided to form two new populations.
#   Default: None
Parameter identifiers : None

#   Structure of model for one population - number of time periods
#   (e.g. 5).
#   Structure of model for two populations - number of time periods
#   before split of ancestral population and after it (e.g. 2,2).
#   Structure of model for three populations - number of time periods
#   before first split, between first and second splits and after
#   second split (e.g. 2,1,2).
#
#   Structure of initial model:
#   Default: all is ones - 1 or 1,1 or 1,1,1
Initial structure : 2,1

#   Structure of final model:
#   Default: equals to initial structure
#Final structure : 2,2

#   It is possible to limit time of splits.
#   Split 1 is the most ancient split.
#   !Note that time is in genetic units (2 * time for 1 generation):
#   e.g. we want to limit by 150 kya, time for one generation is
#   25 years, then bound will be 150000 / (2*25) = 3000.
#
#   Upper bound for split 1 (in case of 2 or 3 populations).
#   Default: None
#Upper bound of first split : None

#   Upper bound for split 2 (in case of 3 populations).

```

```

#   Default: None
#Upper bound of second split : None

#       Print parameters of model in units of  $N_{ref} = N_A$ .
#    $N_A$  will be placed in brackets at the end of string.
#   Default: False
Relative parameters : False

#   Disable migrations in demographic models.
#   Default: False
No migrations : false

#   Local optimization.
#
#   Choice of local optimization, that is launched after
#   each genetic algorithm.
#   Choices:
#
#   *   optimize (BFGS method)
#
#   *   optimiza_log (BFGS method)
#
#   *   optimize_powell (Powell's conjugate direction method)
#       (Note: is implemented in moments: one need to have moments
#       installed.)
#
#   (If optimizations are often hitting the parameter bounds,
#   try using these methods:)
#   *   optimize_lbfgsb
#   *   optimize_log_lbfgsb
#   (Note that it is probably best to start with the vanilla BFGS
#   methods, because the L-BFGS-B methods will always try parameter
#   values at the bounds during the search.
#   This can dramatically slow model fitting.)
#
#   *   optimize_log_fmin (simplex (a.k.a. amoeba) method)
#
#   *   hill_climbing
#
#   Default: optimize_powell
Name of local optimization : optimize_log

```



```

# Parameters of pipeline
#
# One can automatically draw models every N iteration.
# If 0 then never.
# Pictures are saved in GA's directory in picture folder.
# Default: 0
Draw models every N iteration : 100

# One can automatically generate dadl and moments code for models.
# If 0 then only current best model will be printed in GA's
# working directory.
# Also result model will be saved there.
# If specified (not 0) then every N iteration model will be saved
# in python code folder.
# Default: 0
Print models' code every N iteration : 100

# One can choose time units in models' plots: years or thousand
# years (kya, KYA). If time for one generation isn't specified
# then time is in genetic units.
# Default: years
Units of time in drawing : thousand years

# No std output.
# Default: False
Silence : False

# How many times launch GADMA with this parameters.
# Default: 1
Number of repeats : 3

# How many processes to use for this repeats.
# Note that one repeat isn't parallelized, so increasing number
# of processes doesn't effect on time of one repeat.
# It is desirable that the number of repeats is a multiple of
# the number of processes.
# Default: 1
Number of processes : 3

```

Appendix B Example of generated code

For example, *moments* generated code for 2D AFS for YRI, CEU populations from [1]:

```
#current best params = [7194.792822462478, 13410.251542201073,
#                        0.9582544565961783, 13542.979276844108,
#                        12114.968575519626, 2683.3787253409746,
#                        846.6668954957415, 0.00014172779289593632,
#                        0.00012195685425105608]

import matplotlib
matplotlib.use("Agg")
import moments
import numpy as np

def generated_model(params, ns):
    Ns = params[:5]
    Ts = params[5:7]
    Ms = params[7:]
    theta1 = 0.37976
    sts = moments.LinearSystem_1D.steady_state_1D(sum(ns),
                                                    theta=theta1, N=Ns[0])

    fs = moments.Spectrum(sts)

    before = [Ns[0]]
    T = Ts[0]
    after = Ns[1:2]
    growth_funcs = [lambda t: after[0]]
    list_growth_funcs = lambda t: [ f(t) for f in growth_funcs]
    fs.integrate(Npop=list_growth_funcs, tf=T, dt_fac=0.1,
                 theta=theta1)

    before = after
    fs = moments.Manips.split_1D_to_2D(fs, ns[0], sum(ns[1:]))

    before.append((1 - Ns[2]) * before[-1])
    before[-2] *= Ns[2]
    T = Ts[1]
    after = Ns[3:5]
    growth_funcs = [lambda t: after[0],
                    lambda t:
                        before[1] * (after[1] / before[1]) ** (t / T)]
    list_growth_funcs = lambda t: [ f(t) for f in growth_funcs]
    m = np.array([0, params[7],[params[8], 0])
    fs.integrate(Npop=list_growth_funcs, tf=T, m=m, dt_fac=0.1,
                 theta=theta1)
```

```

        before = after
        return fs
data = moments.Spectrum.from_file('data/YRI_CEU.fs')

popt = [7194.792822462478, 13410.251542201073, 0.9582544565961783,
        13542.979276844108, 12114.968575519626, 2683.3787253409746,
        846.6668954957415, 0.00014172779289593632,
        0.00012195685425105608]
ns = [20, 20]
model = generated_model(popt, ns)
ll_model = moments.Inference.ll(model, data)
print('Model log likelihood (LL(model, data)): {0}'.format(ll_model))
#now we need to norm vector of params so that first value is 1:
popt_norm = [1.0, 1.8638829321580523, 0.9582544565961783,
             1.8823306815120924, 1.6838523185401713, 0.3729612223111333,
             0.1176777311575127, 1.0197021070711312, 0.8774542996156008]
print('Drawing model to model_from_GADMA_from_simple.png')
model = moments.ModelPlot.generate_model(generated_model,
                                         pop_norm, ns)
moments.ModelPlot.plot_model(model,
                             save_file='model_from_GADMA_from_simple.png',
                             fig_title='', pop_labels=['YRI', 'CEU'],
                             nref=7194, draw_scale=True, gen_time=0.025,
                             gen_time_units="Thousand years",
                             reverse_timeline=True)

```

References

- [1] Ryan N Gutenkunst, Ryan D Hernandez, Scott H Williamson, and Carlos D Bustamante, *Inferring the joint demographic history of multiple populations from multidimensional snp frequency data*, PLoS genetics **5** (2009), no. 10, e1000695.
- [2] Julien Jouganous, Will Long, Aaron P. Ragsdale, and Simon Gravel, *Inferring the joint demographic history of multiple populations: Beyond the diffusion approximation*, Genetics **206** (2017), no. 3, 1549–1567.
- [3] Marguerite Lapierre, Amaury Lambert, and Guillaume Achaz, *Accuracy of demographic inferences from the site frequency spectrum: the case of the yoruba population*, Genetics (2017), genetics–116.