

# *GeneEvolve* Documentation

Rasool Tahmasbi

Institute for Behavioral Genetics,  
University of Colorado, Boulder,  
USA

`Rasool.Tahmasbi@Colorado.edu`

Matthew C. Keller

Institute for Behavioral Genetics,  
University of Colorado, Boulder,  
USA

`matthew.c.keller@colorado.edu`

version 1.0.0  
March 2016



# Contents

<b>1</b>	<b>Introduction</b>	<b>9</b>
1.1	What is <i>GeneEvolve</i> ?	9
1.2	How to download <i>GeneEvolve</i> ?	9
1.3	How to compile <i>GeneEvolve</i> ?	9
1.4	How <i>GeneEvolve</i> works?	10
1.5	Quick start	10
1.5.1	Example 1 – Simple example	10
1.5.2	Example 2 – Scaling the variance of additive, dominance and unique environmental effects	12
1.5.3	Example 3 – No dominance effect	12
1.5.4	Example 4 – Assortative mating	13
1.5.5	Example 5 – Random mating	14
1.5.6	Example 6 – Exponential population size	15
1.5.7	Example 7 – Bottleneck population	15
1.5.8	Example 8 – Simulating several phenotypes	16
1.5.9	Example 9 – Selection	17
1.5.10	Example 10 – Comprehensive example	17
<b>2</b>	<b>Population Genetics Models</b>	<b>19</b>
2.1	Population’s basic information	19
2.1.1	The Wright–Fisher model	19
2.1.2	Population size in <i>GeneEvolve</i>	19
2.1.3	Number of offspring’s distribution	19
2.2	Haplotypes of initial population	20
2.3	Recombination and linkage	21
2.3.1	Introduction	21
2.3.2	Recombination in <i>GeneEvolve</i>	22
2.4	Simulating complex quantitative traits	23
2.4.1	Genotypic value for single locus	23
2.4.2	One phenotype	24
2.4.3	Heritability	25
2.4.4	Phenotypes in <i>GeneEvolve</i>	26
2.4.5	Scaling VA and VD	26
2.4.6	Simulating multivariate phenotypes	27
2.5	Random and non-random mating systems	27

2.6	Natural selection . . . . .	27
2.6.1	Directional selection . . . . .	28
2.6.2	Stabilizing selection . . . . .	29
2.6.3	Threshold selection . . . . .	29
2.7	Familial effect . . . . .	29
2.8	Shared sibling (common) effect . . . . .	30
2.9	Environmental effects specific to each population . . . . .	30
2.10	Simulating several populations . . . . .	30
2.11	Population structure . . . . .	31
2.11.1	Introduction . . . . .	31
2.11.2	Population structure in <i>GeneEvolve</i> . . . . .	31
<b>3</b>	<b>Parameters Reference</b>	<b>35</b>
3.1	--file_gen_info . . . . .	36
3.2	--file_hap_name . . . . .	36
3.3	--file_recom_map . . . . .	36
3.4	--file_cv_info . . . . .	37
3.5	--file_cvs . . . . .	37
3.6	--va . . . . .	37
3.7	--vd . . . . .	37
3.8	--vc . . . . .	37
3.9	--ve . . . . .	37
3.10	--vf . . . . .	37
3.11	--next_population . . . . .	37
<b>A</b>	<b>C++ Classes</b>	<b>39</b>
A.1	Human . . . . .	39
A.2	Population . . . . .	41
A.3	Simulation . . . . .	42
<b>B</b>	<b>File formats</b>	<b>43</b>
B.1	Hap – Legend – Sample . . . . .	43
B.1.1	.hap . . . . .	43
B.1.2	.legend . . . . .	43
B.1.3	.sample . . . . .	44
B.1.4	.impute.hap.indv . . . . .	44
B.2	PLINK . . . . .	44
B.2.1	.ped . . . . .	44
B.2.2	.map . . . . .	45
B.2.3	.fam . . . . .	45

# List of Figures

1.1	<i>GeneEvolve</i> command line (basic parameters)	10
2.1	The structure of file_generaions_info.txt file.	20
2.2	The structure of [file.txt] in [-file_hap_name].	21
2.3	Recombination model: three recombination occurred at different positions. Each color code is an ancestral IBD.	22
2.4	Arbitrary assigned genotypic value.	23
2.5	Additive term	24
2.6	Mating path diagram and phenotypes	28
2.7	Path diagram for migration and environmental effects specific to each population	33



# List of Tables

2.1	Values of genotypes in a two-allele system, measured as deviations from the population mean. The population mean is $a(p - q) + 2pqd$ , and $\alpha = a + d(q - p)$ . . . . .	24
2.2	The migration matrix . . . . .	32





# Chapter 1

## Introduction

### 1.1 What is *GeneEvolve*?

*GeneEvolve* is C++ code for simulating sequence-level genetic data over large genomic regions in large populations, using an object-oriented approach. This allows compiling *GeneEvolve* on any computer platform, which supports standard C++ compiler.

Computer simulations are excellent tools for understanding the evolutionary and genetic consequences of complex processes whose interactions cannot be analytically derived. Unlike coalescent [2] based simulators, *GeneEvolve* runs forward-in-time, which allows it to provide a wide range of scenarios for selection, population size and structure, migration, recombination and familial effects.

*GeneEvolve* is fast and memory efficient simulator which can handle complex life events. User-friendly and easy to work are among its other advantages.

### 1.2 How to download *GeneEvolve*?

*GeneEvolve* can be download from <https://github.com/rtahmasbi/GeneEvolve>. The source codes, examples and this documentation are available freely for downloading.

You also can run the following commands for downloading, compiling and testing *GeneEvolve*:

```
1 wget https://github.com/rtahmasbi/GeneEvolve/archive/master.zip
2 unzip master
3 cd GeneEvolve-master/
4 make
5 cd bin
6 ./GeneEvolve --help
```

### 1.3 How to compile *GeneEvolve*?

*GeneEvolve* can be run on different platforms. After downloading it, you need to uncompress the zipped file *GeneEvolve.zip* and then go to the root directory and type *make*. After successful compiling, the *GeneEvolve* simulator will be in the *bin* subdirectory.

```

1 ./GeneEvolve --file_gen_info [file] \
2 --file_hap_name [file] \
3 --file_recom_map [file] \
4 --file_cv_info [file] \
5 --file_cvs [file]

```

Figure 1.1: *GeneEvolve* command line (basic parameters)

You also need to load a standard C++ compiler, by the command

```

1 module load gcc/gcc-4.9.2

```

## 1.4 How *GeneEvolve* works?

*GeneEvolve* is a stand-alone program. Our aim was to make it simple and user friendly for different levels of user's knowledge. Users can create complex life events by adding and combining more parameters. The minimum required parameters are *generation information* (such as population size, spousal correlation for mating, offspring distribution and selection function per generation), *haplotype information* (which is haplotypes for starting generation) *recombination map*, *cv list* (position of CVs and their additive and dominance effects) and the actual *CV's haplotype*. These parameters are listed in Figure 1.1. The detailed explanation for each parameter is illustrated in the following chapters.

For different levels of variance of additive and dominance effects, user can use the following parameters

```
--va [number] --vd [number]
```

and for the unique, familial, shared sibling (common), and population specific environmental effects, user can also specify the following parameters, respectively:

```
--ve [number] --vf [number] --vc [number] --gamma [number]
```

There are more parameters for random mating, selection function, migration and so on, available in this documentation file.

## 1.5 Quick start

In this section we will present several examples with different lifetime scenarios. All the codes are available in **Examples** directory. To run *GeneEvolve*, we first need genotype data. In this directory, we have simulated genotypes for 3 chromosomes. Clearly, the structure of real genotypes are more complex than this simulated data, but this simple data is useful for educational purpose.

### 1.5.1 Example 1 – Simple example

In this example, we simulate a population of size 3000 for 10 generations. Then we run *GeneEvolve* by the following command (you should change the **path** to the appropriate directory)

```

1 /path/GeneEvolve \
2 --file_gen_info ex1.popinfo.txt \
3 --file_hap_name par.pop1.hap_sample_address.txt \
4 --file_recom_map Recom.Map.b37.50KbDiff \
5 --file_cv_info cv.info \
6 --file_cvs par.pop1.cv_hap_files.txt \
7 --prefix out.ex1

```

## Code for Example 1

The population information for each generation is in the file `ex1.popinfo.txt`

```

1 pop_size mat_cor offspring_dist selection_func selection_func_par1 selection_func_par2
2 3000 0 p logit 20 0
3 3000 0 p logit 20 0
4 3000 0 p logit 20 0
5 3000 0 p logit 20 0
6 3000 0 p logit 20 0
7 3000 0 p logit 20 0
8 3000 0 p logit 20 0
9 3000 0 p logit 20 0
10 3000 0 p logit 20 0
11 3000 0 p logit 20 0

```

`ex1.popinfo.txt`

As can be seen, the first line is header and there are 10 lines (for each generation) with a population of size 3000. The other parameters (columns 2-6) will be explained later.

For each generation, the program simulate phenotypes and reports the additive (A), dominance (D), genotypic value ( $G=A+D$ ), common sibling effect (C), environmental effect (E), parental effect (F), phenotypic value ( $P=A+D+C+E+F$ ) for each phenotype, and reports mating value (MV) and selection value (SV) (see the `out.ex1.info.pop1.gen10.txt` listing for generation 10).

```

1 ID ID_Father ID_Mother ID_Fathers_Father ID_Fathers_Mother ID_Mothers_Father ID_Mothers_Mother sex ph1_A ph1_D
  ph1_G ph1_C ph1_E ph1_F ph1_P MV SV
2 1 771 184 957 2899 1226 559 1 -26.7357 1.01329 -25.7224 0 0.359913 0 -25.3625 -25.3625 1
3 2 771 184 957 2899 1226 559 2 -15.3725 2.87523 -12.4973 0 0.767739 0 -11.7295 -11.7295 1
4 3 830 3048 1987 688 1738 1818 2 -33.0872 6.65026 -26.437 0 0.540759 0 -25.8962 -25.8962 1
5 4 2148 703 1717 1270 2638 2423 1 1.093 -0.631975 0.461025 0 -0.335767 0 0.125258 0.125258 1
6 5 2148 703 1717 1270 2638 2423 1 13.6752 -1.90592 11.7693 0 -2.08945 0 9.67983 9.67983 1
7 6 2148 703 1717 1270 2638 2423 1 4.99932 -2.44006 2.55926 0 0.591915 0 3.15117 3.15117 1
8 7 2091 191 2884 2926 2650 1232 1 -12.4267 -1.09799 -13.5246 0 0.0112131 0 -13.5134 -13.5134 1
9 8 2091 191 2884 2926 2650 1232 2 1.11746 -2.9315 -1.81404 0 1.16346 0 -0.650585 -0.650585 1
10 9 1568 890 2283 440 2925 1329 1 -24.4192 -6.27226 -30.6915 0 0.716498 0 -29.975 -29.975 1
11 ...

```

`out.ex1.info.pop1.gen10.txt`

The program reports the summary statistics for each generation at one output file:

```

1 gen ph1_var_A ph1_var_D ph1_var_G ph1_var_C ph1_var_E ph1_var_parental_effect ph1_var_phen ph1_h2 ph1_var_G_std
  var_mating_value
2 0 193.624 25.5632 214.176 0 1 0 214.703 0.997547 1 214.703
3 1 168.984 29.557 211.362 0 1 0 211.64 0.998686 0.986857 211.64
4 2 207.226 31.3312 240.297 0 1 0 241.154 0.996446 1.12196 241.154
5 3 226.008 30.6797 257.723 0 1 0 257.814 0.999646 1.20332 257.814
6 4 226.706 29.4573 260.352 0 1 0 261.553 0.995408 1.21559 261.553
7 5 227.998 32.2255 261.462 0 1 0 263.302 0.993014 1.22078 263.302
8 6 230.707 29.4809 260.103 0 1 0 260.293 0.999271 1.21443 260.293
9 7 251.877 31.5432 278.559 0 1 0 279.491 0.996665 1.30061 279.491

```

```

10 8 258.136 30.6762 290.039 0 1 0 292.24 0.99247 1.35421 292.24
11 9 247.686 31.344 272.789 0 1 0 273.748 0.996496 1.27367 273.748
12 10 232.669 30.9782 266.217 0 1 0 267.524 0.995116 1.24298 267.524

```

out.ex1.pop1.summary

By default,  $\text{var}[E] = 1$  and  $\text{var}[A]$  and  $\text{var}[D]$  are computed based on the additive and dominance effect sizes inputted by option `--file_cv_info cv.info`. A user can scale  $\text{var}[A]$ ,  $\text{var}[D]$  and  $\text{var}[E]$  to any positive number. Next example illustrates scaling the variances.

### 1.5.2 Example 2 – Scaling the variance of additive, dominance and unique environmental effects

In order to scale the variance of additive, dominance and unique environmental effects, we run the following command:

```

1 /path/GeneEvolve \
2 --file_gen_info ex1.popinfo.txt \
3 --file_hap_name par.pop1.hap_sample_address.txt \
4 --file_recom_map Recom.Map.b37.50KbDiff \
5 --file_cv_info cv.info \
6 --file_cvs par.pop1.cv_hap_files.txt \
7 --va 3 \
8 --vd 1 \
9 --ve 2 \
10 --no_output \
11 --prefix out.ex2

```

Code for Example 2

The summary statistics output file is:

```

1 gen ph1_var_A ph1_var_D ph1_var_G ph1_var_C ph1_var_E ph1_var_parental_effect ph1_var_phen ph1_h2 ph1_var_G_std
  var_mating_value
2 0 3 1 3.87665 0 2 0 5.16892 0.749992 1 5.16892
3 1 2.55356 1.21195 3.85699 0 2 0 5.70883 0.675617 0.994929 5.70883
4 2 3.02579 1.22594 4.21461 0 2 0 6.14022 0.686394 1.08718 6.14022
5 3 3.37727 1.26539 4.70582 0 2 0 6.71178 0.701128 1.21389 6.71178
6 4 3.49527 1.18061 4.69215 0 2 0 6.66297 0.704214 1.21036 6.66297
7 5 3.47466 1.23043 4.97902 0 2 0 6.88022 0.723673 1.28436 6.88022
8 6 3.54991 1.23196 4.81836 0 2 0 6.80225 0.708348 1.24292 6.80225
9 7 3.46914 1.27584 4.85821 0 2 0 6.73862 0.72095 1.2532 6.73862
10 8 3.64999 1.18515 4.79675 0 2 0 6.90517 0.69466 1.23734 6.90517
11 9 3.54014 1.22018 4.73613 0 2 0 6.59473 0.718168 1.22171 6.59473
12 10 3.60618 1.2206 4.66857 0 2 0 6.81362 0.685181 1.20428 6.81362

```

out.ex2.pop1.summary

In this listing, the variances are scaled to  $\text{var}[A] = 3$ ,  $\text{var}[D] = 1$  and  $\text{var}[E] = 2$  for generation zero (initial population) at line 2. The variances for the next generations are computed and scaled according to the initial values. Note that we use the option `--no_output` to save space by not creating the output hap files.

### 1.5.3 Example 3 – No dominance effect

In order to work just with the additive effect and not dominance effect, we run the following command by setting `--vd 0`:

```

1 /path/GeneEvolve \
2 --file_gen_info ex1.popinfo.txt \
3 --file_hap_name par.pop1.hap_sample_address.txt \
4 --file_recom_map Recom.Map.b37.50KbDiff \
5 --file_cv_info cv.info \
6 --file_cvs par.pop1.cv_hap_files.txt \
7 --va 3 \
8 --vd 0 \
9 --ve 2 \
10 --RM \
11 --avoid_inbreeding \
12 --no_output \
13 --prefix out.ex3

```

### Code for Example 3

The summary statistics output file is:

```

1 gen ph1_var_A ph1_var_D ph1_var_G ph1_var_C ph1_var_E ph1_var_parental_effect ph1_var_phen ph1_h2 ph1_var_G_std
   var_mating_value
2 0 3 0 3 0 2 0 5.24393 0.57209 1 5.24393
3 1 2.96881 0 2.96881 0 2 0 5.04622 0.588325 0.989605 5.04622
4 2 3.40633 0 3.40633 0 2 0 5.45592 0.624338 1.13544 5.45592
5 3 3.57616 0 3.57616 0 2 0 5.68366 0.6292 1.19205 5.68366
6 4 3.68187 0 3.68187 0 2 0 5.78033 0.636966 1.22729 5.78033
7 5 3.9199 0 3.9199 0 2 0 5.90376 0.663966 1.30663 5.90376
8 6 3.7646 0 3.7646 0 2 0 5.6695 0.66401 1.25487 5.6695
9 7 3.79647 0 3.79647 0 2 0 5.94024 0.639111 1.26549 5.94024
10 8 3.80745 0 3.80745 0 2 0 5.88294 0.647202 1.26915 5.88294
11 9 3.53703 0 3.53703 0 2 0 5.79052 0.610832 1.17901 5.79052
12 10 3.62544 0 3.62544 0 2 0 5.63872 0.642955 1.20848 5.63872

```

out.ex3.pop1.summary

In this listing, the variances are scaled to  $\text{var}[A] = 3$ ,  $\text{var}[D] = 0$  and  $\text{var}[E] = 2$ . The option `--avoid_inbreeding` is used to not let inbreeding.

### 1.5.4 Example 4 – Assortative mating

In assortative mating (AM) system, couples choose each other based on the mating value, which is a function of phenotypes. For example taller mates have taller spouses with some correlations. In the following example we set the mating correlation equal to 0.5. You can specify this correlation in the second column of `ex4.popinfo.txt` file, for each generation. It's known that AM, will increase the genetic variance and heritability.

```

1 /path/GeneEvolve \
2 --file_gen_info ex4.popinfo.txt \
3 --file_hap_name par.pop1.hap_sample_address.txt \
4 --file_recom_map Recom.Map.b37.50KbDiff \
5 --file_cv_info cv.info \
6 --file_cvs par.pop1.cv_hap_files.txt \
7 --va 1 \
8 --vd 0 \
9 --ve 1 \
10 --avoid_inbreeding \
11 --no_output \

```

```
12 | --prefix out.ex4
```

### Code for Example 4

The population information for each generation is in the file `ex4.popinfo.txt`

```
1 | pop_size mat_cor offspring_dist selection_func selection_func_par1 selection_func_par2
2 | 3000 0.5 p logit 20 0
3 | 3000 0.5 p logit 20 0
4 | 3000 0.5 p logit 20 0
5 | 3000 0.5 p logit 20 0
6 | 3000 0.5 p logit 20 0
7 | 3000 0.5 p logit 20 0
8 | 3000 0.5 p logit 20 0
9 | 3000 0.5 p logit 20 0
10 | 3000 0.5 p logit 20 0
11 | 3000 0.5 p logit 20 0
```

`ex4.popinfo.txt`

The summary statistics output file is:

```
1 | gen ph1_var_A ph1_var_D ph1_var_G ph1_var_C ph1_var_E ph1_var_parental_effect ph1_var_phen ph1_h2 ph1_var_G_std
   | var_mating_value
2 | 0 1 0 1 0 1 0 1.87267 0.533996 1 1.87267
3 | 1 1.10922 0 1.10922 0 1 0 2.18531 0.507578 1.10922 2.18531
4 | 2 1.28119 0 1.28119 0 1 0 2.29807 0.557504 1.28119 2.29807
5 | 3 1.43656 0 1.43656 0 1 0 2.50903 0.572554 1.43656 2.50903
6 | 4 1.44637 0 1.44637 0 1 0 2.42259 0.597034 1.44637 2.42259
7 | 5 1.62156 0 1.62156 0 1 0 2.65714 0.610264 1.62156 2.65714
8 | 6 1.60513 0 1.60513 0 1 0 2.59393 0.618801 1.60513 2.59393
9 | 7 1.57147 0 1.57147 0 1 0 2.57401 0.610515 1.57147 2.57401
10 | 8 1.56984 0 1.56984 0 1 0 2.5517 0.615215 1.56984 2.5517
11 | 9 1.65574 0 1.65574 0 1 0 2.69236 0.614976 1.65574 2.69236
12 | 10 1.71919 0 1.71919 0 1 0 2.71289 0.633712 1.71919 2.71289
```

`out.ex4.pop1.summary`

As it can be seen, the variance increased from 1 to 1.719 and heritability from 0.5 to 0.63.

## 1.5.5 Example 5 – Random mating

Another mating system is random mating (RM), where couples are chosen randomly. You can use `--RM` in *GeneEvolve*. In RM, the second column of `ex1.popinfo.txt` will not be used, since in the RM system there is no mating correlation.

```
1 | /path/GeneEvolve \
2 | --file_gen_info ex1.popinfo.txt \
3 | --file_hap_name par.pop1.hap_sample_address.txt \
4 | --file_recom_map Recom.Map.b37.50KbDiff \
5 | --file_cv_info cv.info \
6 | --file_cvs par.pop1.cv_hap_files.txt \
7 | --va 1 \
8 | --vd 0 \
9 | --ve 1 \
10 | --avoid_inbreeding \
11 | --no_output \
12 | --RM
```

```
13 --prefix out.ex5
```

Code for Example 5

### 1.5.6 Example 6 – Exponential population size

In *GeneEvolve* it is possible to simulate any scenario for population growth by modifying the first column of `ex6.popinfo.txt` file. Here is an example for exponential population size.

```
1 /path/GeneEvolve \
2 --file_gen_info ex6.popinfo.txt \
3 --file_hap_name par.pop1.hap_sample_address.txt \
4 --file_recom_map Recom.Map.b37.50KbDiff \
5 --file_cv_info cv.info \
6 --file_cvs par.pop1.cv_hap_files.txt \
7 --avoid_inbreeding \
8 --no_output \
9 --prefix out.ex6
```

Code for Example 6

The population information for each generation is in the `ex6.popinfo.txt` file:

```
1 pop_size mat_cor offspring_dist selection_func selection_func_par1 selection_func_par2
2 500 0 p logit 20 0
3 1000 0 p logit 20 0
4 2000 0 p logit 20 0
5 4000 0 p logit 20 0
6 8000 0 p logit 20 0
7 16000 0 p logit 20 0
8 32000 0 p logit 20 0
```

`ex6.popinfo.txt`

### 1.5.7 Example 7 – Bottleneck population

In this example, we simulate a bottleneck population with a sharp reduction in the size. The population size reduce to 200 in the 8th generation.

```
1 /path/GeneEvolve \
2 --file_gen_info ex7.popinfo.txt \
3 --file_hap_name par.pop1.hap_sample_address.txt \
4 --file_recom_map Recom.Map.b37.50KbDiff \
5 --file_cv_info cv.info \
6 --file_cvs par.pop1.cv_hap_files.txt \
7 --avoid_inbreeding \
8 --no_output \
9 --prefix out.ex7
```

Code for Example 7

The population information for each generation is in the `ex7.popinfo.txt` file:

```

1 pop_size mat_cor offspring_dist selection_func selection_func_par1 selection_func_par2
2 3000 0 p logit 20 0
3 3000 0 p logit 20 0
4 3000 0 p logit 20 0
5 3000 0 p logit 20 0
6 3000 0 p logit 20 0
7 3000 0 p logit 20 0
8 3000 0 p logit 20 0
9 200 0 p logit 20 0
10 250 0 p logit 20 0
11 300 0 p logit 20 0
12 350 0 p logit 20 0
13 500 0 p logit 20 0
14 700 0 p logit 20 0
15 1000 0 p logit 20 0

```

ex7.popinfo.txt

### 1.5.8 Example 8 – Simulating several phenotypes

In this example, we simulate a population with two phenotypes. In *GeneEvolve* you can simulate any number of phenotypes simply by adding CV information with the parameters `--file_cv_info` and `--file_cvs`.

```

1 /path/GeneEvolve \
2 --file_gen_info ex1.popinfo.txt \
3 --file_hap_name par.pop1.hap_sample_address.txt \
4 --file_recom_map Recom.Map.b37.50KbDiff \
5 --file_cv_info cv.info \
6 --file_cvs par.pop1.cv_hap_files.txt \
7 --file_cv_info cv2.info \
8 --file_cvs par.pop1.cv2_hap_files.txt \
9 --avoid_inbreeding \
10 --no_output \
11 --prefix out.ex8

```

Code for Example 8

It is possible to scale the additive variance for each phenotype by adding the parameter `--va`. This is also true for other variances (dominance, environmental effect and so on).

The summary statistics output file is:

```

1 gen ph1_var_A ph1_var_D ph1_var_G ph1_var_C ph1_var_E ph1_var_parental_effect ph1_var_phen ph1_h2 ph1_var_G_std
   ph2_var_A ph2_var_D ph2_var_G ph2_var_C ph2_var_E ph2_var_parental_effect ph2_var_phen ph2_h2 ph2_var_G_std
   var_mating_value
2 0 770.722 35.2391 850.553 0 1 0 844.433 1.00725 1 312.507 70.7022 383.235 0 1 0 380.49 1.00722 1 1173.35
3 1 452.149 33.1611 439.712 0 1 0 441.35 0.996287 0.516972 237.848 90.5534 347.64 0 1 0 347.214 1.00123 0.907119
   697.072
4 2 397.948 31.2319 433.186 0 1 0 434.912 0.996032 0.509299 261.977 76.4754 350.085 0 1 0 349.863 1.00063 0.913498
   724.967
5 3 360.242 30.2915 388.889 0 1 0 390.463 0.995967 0.457219 267.811 77.0184 336.895 0 1 0 337.483 0.998259 0.879082
   670.047
6 4 363.391 30.566 390.824 0 1 0 393.339 0.993607 0.459494 274.691 77.4435 358.635 0 1 0 359.433 0.997781 0.935809
   697.42
7 5 349.659 31.0504 384.623 0 1 0 386.509 0.99512 0.452203 273.84 80.2313 360.162 0 1 0 360.167 0.999985 0.939793
   682.824

```



```

8 | 6 343.226 31.4154 367.685 0 1 0 368.228 0.998525 0.432289 284.365 74.9454 369.714 0 1 0 369.673 1.00011 0.964718
   | 680.261
9 | 7 316.698 30.3888 339.883 0 1 0 340.16 0.999186 0.399602 294.572 80.8907 380.605 0 1 0 382.025 0.996282 0.993136
   | 680.898
10 | 8 317.104 29.8229 349.112 0 1 0 349.81 0.998006 0.410453 280.626 76.5021 346.795 0 1 0 347.227 0.998757 0.904914
    | 657.258
11 | 9 300.557 29.5321 327.926 0 1 0 329.548 0.995077 0.385544 289.043 81.2075 349.52 0 1 0 350.632 0.996827 0.912024
    | 661.524
12 | 10 329.317 29.4555 356.58 0 1 0 357.765 0.99669 0.419234 275.649 81.7718 365.932 0 1 0 366.569 0.998262 0.954849
    | 686.647

```

out.ex8.pop1.summary

As it can be seen, there are two phenotypes. Note that in the multi-phenotype simulations, the selection value and mating value is the average of the phenotypes. You can easily choose different weights for them using the parameters `--lambda` and `--omega` (see figure 2.6 and chapter 3)

### 1.5.9 Example 9 – Selection

In this example, we simulate a population with an under selection phenotype. For each generation, you can choose different selection functions and different parameters in the columns 4–6 of `ex9.popinfo.txt` file using the parameter `--file_gen_info`.

```

1 | /path/GeneEvolve \
2 | --file_gen_info ex9.popinfo.txt \
3 | --file_hap_name par.pop1.hap_sample_address.txt \
4 | --file_recom_map Recom.Map.b37.50KbDiff \
5 | --file_cv_info cv.info \
6 | --file_cvs par.pop1.cv_hap_files.txt \
7 | --avoid_inbreeding \
8 | --no_output \
9 | --prefix out.ex9

```

Code for Example 9

The population information for each generation is in the `ex9.popinfo.txt` file:

```

1 | pop_size mat_cor offspring_dist selection_func selection_func_par1 selection_func_par2
2 | 3000 0 p logit 20 0
3 | 3000 0 p logit 10 5
4 | 3000 0 p logit 1 1
5 | 3000 0 p probit 0 1
6 | 3000 0 p probit 0 2
7 | 3000 0 p stab 0 1
8 | 3000 0 p stab 1 4
9 | 3000 0 p thr .8 3
10 | 3000 0 p thr .8 10
11 | 3000 0 p thr .9 20

```

ex9.popinfo.txt

For the definition of `logit`, `probit`, `stab` and `thr` functions see Section 2.6.

### 1.5.10 Example 10 – Comprehensive example

For a comprehensive example, we simulate two populations, each with two phenotypes, in Chapter 3. The initial inputted genotype files are two different datasets. The variance components for each phenotype are

scaled. Populations are under selection with different selection functions. The selection functions also change over time (over generations). The selection and mating values are computed with unequal weights. The mating system is assortative mating while avoiding inbreeding. There are also familial and siblings effects. The first phenotype has more effect on the familial effect. Populations migrate with different migration rates per generation.

# Chapter 2

## Population Genetics Models

### 2.1 Population's basic information

#### 2.1.1 The Wright–Fisher model

The total population size is more often determined by external factors like availability of food or living space, or the action of predators, than by summing independent family sizes (the branching process models). His first approximation to reality is therefore a model in which the total population size is a fixed number dictated by external constraints, and the most popular is that associated with the names of Sewall Wright and R. A. Fisher (for more info see [3]).

This assumes discrete, non-overlapping generations  $G_0, G_1, G_2, \dots$  in which each generation contains a fixed number  $N$  of individuals. Each member of  $G_{i+1}$  is the child of exactly one member of  $G_i$ , but the number of children born to the  $j$ th member of  $G_i$  is a random variable  $v_i$  subject of course to the constraint

$$\sum_{j=1}^N v_i = N.$$

#### 2.1.2 Population size in *GeneEvolve*

For inputting the basic information of population, user should use the parameter

```
--file_gen_info [file_generaions_info.txt]
```

where `file_generaions_info.txt` is a text file with 6 columns and  $n + 1$  lines, where  $n$  is the number of generations to be simulated by *GeneEvolve*. This file should have a header and the first column is the population size. The structure of this file is listed in figure 2.1. *GeneEvolve* can simulate different population sizes in each generation, by specifying some numbers in the first column.

#### 2.1.3 Number of offspring's distribution

The third column of file `file_generaions_info.txt` determines the distribution of offsprings per generation. It can be `p` for Poisson or `f` for fixed number of offsprings (see the third column of figure 2.1).

The mean of Poisson distribution in each generation is obtained by

$$\frac{N_i}{C_i},$$

```

1 pop_size mat_cor offspring_dist selection_func selection_func_par1 selection_func_par2
2 3000 0 p logit 20 0
3 3000 0 p logit 20 0
4 3000 0 p logit 20 0
5 3000 0 p logit 20 0
6 3000 0 p logit 20 0
7 3000 0 p logit 20 0
8 3000 0 p logit 20 0
9 3000 0 p logit 20 0
10 3000 0 p logit 20 0
11 3000 0 p logit 20 0
12 3000 0 p logit 20 0
13 3000 0 f logit 20 0

```

Figure 2.1: The structure of file\_generaions\_info.txt file.

where  $N_i$  is the user specified population size (first column of `file_generaions_info.txt`) and  $C_i$  the number of couples at generation  $i$ .

For the fixed number of offsprings, we have

$$k = \text{round}\left(\frac{N_i}{C_i}\right).$$

The selection function has a direct effect on  $C_i$ . If the selection function allows all individuals to marry, then  $C_i$  should be  $N_i/2$  in average. So each family should have 2 offsprings in average. More stringent selection function can reduce the  $C_i$  and as the results, it increases the average family size,  $N_i/C_i$ . In Section 2.6, you can find more information about the selection function.

## 2.2 Haplotypes of initial population

*GeneEvolve* works with the haplotypes. For more information about haplotype file format, see the appendix B. Since the size of genome can be large, the haplotype of each chromosome should be in a separate file. For the initial population (generation 0), the user should specify a file containing the *address* of each chromosome. With the parameter `--file_hap_name [file.txt]`, user should address the `.hap`, `.legend`, and `.indv` files, respectively, for each chromosome per line. Figure 2.2 shows a typical example with 22 chromosomes. This file has header and the first column is chromosome number. Note that the file `.indv` has no header.

It is also possible for *GeneEvolve* to work just with few chromosome. See the following example, where the user simulate just with chromosomes 6, 10 and 22.

**Example 1** (Working with few chromosomes). If user wants to simulate a population with 3 chromosome (6, 10 and 22), he/she should prepare the following file for the parameter `--file_hap_name [file.txt]`.

```

1 chr hap legend indv
2 6 /path/chr6.hap /path/chr6.legend /path/chr6.indv
3 10 /path/chr10.hap /path/chr10.legend /path/chr10.indv
4 22 /path/chr22.hap /path/chr22.legend /path/chr22.indv

```

file.txt

```

1 chr hap legend indv
2 1 /path/chr1.hap /path/chr1.legend /path/chr1.indv
3 2 /path/chr2.hap /path/chr2.legend /path/chr2.indv
4 3 /path/chr3.hap /path/chr3.legend /path/chr3.indv
5 4 /path/chr4.hap /path/chr4.legend /path/chr4.indv
6 5 /path/chr5.hap /path/chr5.legend /path/chr5.indv
7 6 /path/chr6.hap /path/chr6.legend /path/chr6.indv
8 7 /path/chr7.hap /path/chr7.legend /path/chr7.indv
9 8 /path/chr8.hap /path/chr8.legend /path/chr8.indv
10 9 /path/chr9.hap /path/chr9.legend /path/chr9.indv
11 10 /path/chr10.hap /path/chr10.legend /path/chr10.indv
12 11 /path/chr11.hap /path/chr11.legend /path/chr11.indv
13 12 /path/chr12.hap /path/chr12.legend /path/chr12.indv
14 13 /path/chr13.hap /path/chr13.legend /path/chr13.indv
15 14 /path/chr14.hap /path/chr14.legend /path/chr14.indv
16 15 /path/chr15.hap /path/chr15.legend /path/chr15.indv
17 16 /path/chr16.hap /path/chr16.legend /path/chr16.indv
18 17 /path/chr17.hap /path/chr17.legend /path/chr17.indv
19 18 /path/chr18.hap /path/chr18.legend /path/chr18.indv
20 19 /path/chr19.hap /path/chr19.legend /path/chr19.indv
21 20 /path/chr20.hap /path/chr20.legend /path/chr20.indv
22 21 /path/chr21.hap /path/chr21.legend /path/chr21.indv
23 22 /path/chr22.hap /path/chr22.legend /path/chr22.indv

```

Figure 2.2: The structure of [file.txt] in [-file\_hap\_name].

## 2.3 Recombination and linkage

### 2.3.1 Introduction

Genetic *recombination*, also called *crossing over*, refers to genetic events that can occur during the formation of sperm and egg cells. During the early stages of cell division in meiosis, two chromosomes of a homologous pair may exchange segments, producing genetic variations in germ cells. For example, if one homologous chromosome has a *haplotype* (genetic sequence on the same chromosome) *AB*, and another homologous chromosome has a haplotype *ab*, one of the gamete cells, because of recombination, may have a chromosome with genotype *Ab*. Such gametes are called recombinants. The proportion of recombinants is called the *recombination rate* between these two loci, which is  $1/2$  if two loci are on two different chromosomes, and thus segregate independently.

The *genetic distance* (also called *map distance*) between two loci is defined as the average number of crossovers between the loci per meiosis. The unit of genetic distance is the *centiMorgan* (cM). Two loci are 1 cM apart if on average there is one crossover occurring between these two loci on a single strand for every 100 meioses. The distribution of recombination events varies between the sexes: females have on average 1.65-fold more recombination events when they make eggs than males do when they make sperms. Even on a single chromosome, recombination rate is uneven, and there exists recombination hotspots with peak recombination rate hundreds or thousands times that of the surrounding regions.

Given a reference panel of haplotypes  $H_N = \{h_1, \dots, h_{2N}\}$  as input, where each haplotype is typed at  $L$  bi-allelic sites, that is  $h_i = (h_{i,1}, \dots, h_{i,L})$  and  $h_{i,j} \in \{0, 1\}$ , the program chooses mates based on the AM or RM, and for the newly simulated child, each haplotype is a recombination of parent's haplotypes, where

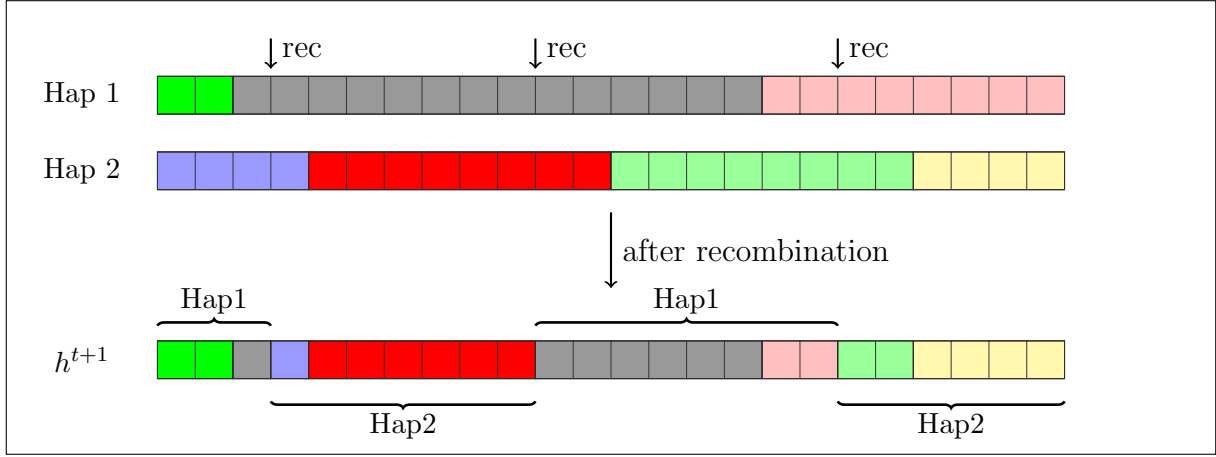


Figure 2.3: Recombination model: three recombination occurred at different positions. Each color code is an ancestral IBD.

the recombination rate is inputted by the user. Precisely, assume that at generation  $t$ , the haplotypes of one parent are  $h_i^t$  and  $h_{i+1}^t$ . These haplotypes are stored as a continues sequence of half open intervals, i.e.,  $h_i^t = \{[1, r_1) \cup [r_1, r_2) \cdots \cup [r_{v-1}, L]\}$  and  $h_{i+1}^t = \{[1, s_1) \cup [s_1, s_2) \cdots \cup [s_{u-1}, L]\}$ , for any arbitrary numbers  $u$  and  $v$ . Now assume that at a genetic phase, for a new child a recombination occurs at position  $w$ . If  $r_2 \leq w < r_3$  and  $s_3 \leq w < s_4$ , then the new recombined haplotype becomes

$$h_i^{t+1} = \{[1, r_1) \cup [r_1, r_2) \cup [r_2, w) \cup [w, s_4) \cup [s_4, s_5) \cdots \cup [s_{u-1}, L]\}.$$

It is also possible that several recombinations occurs; the idea is the same. For a graphical illustration see the figure 2.3.

Clearly, working with a continues sequence of intervals is much faster than working with real genotypes. Saving them in computer memory is more efficient, too.

Based on the user-specified mutation rate, the position of new mutations will also store for the new generated haplotypes.

### 2.3.2 Recombination in *GeneEvolve*

different positions of genome have different recombination probabilities. A high-resolution recombination map of the human genome is estimated by Kong et. al. [4] and is available online.

To locate the recombination file, you can use the parameter `--file_recom_map [file_recom.txt]`. The file `file_recom.txt` has a header with 3 columns: chromosome number, base-pair distance and cM distance.

**Example 2** (Recombination map file format). In the following listing file, you can see part of a equidistant genetic map (with 50k length). By definition, the probability of recombination in a chunk is the difference between its two cM distances divided by 100. These probabilities will be computed automatically by *GeneEvolve* program.

```

1 chr bp      cM
2 1 1128555 1.13368814337268
3 1 1138555 1.14905703198189
4 1 1148555 1.15742981154837
5 1 1158555 1.16581577645964

```

```

6 | 1 1168555 1.17462263654929
7 | 1 1178555 1.18341927550241
8 | 1 1188555 1.19221591445554
9 | 1 1198555 1.20104594872684
10 | 1 1208555 1.20989254327934
11 | 1 1218555 1.21873913783184

```

file.recom.txt

For example, the probability occurrence of one recombination in the interval [1148555, 1158555) is equal to

$$\frac{1.16581577645964 - 1.15742981154837}{100} = 0.00008385964911.$$

## 2.4 Simulating complex quantitative traits

Computer programs that can simulate genotypes with phenotypes based on user-specified disease or quantitative trait models are useful in genetic studies. They can be used to evaluate statistical power when planning a study design based on the proposed sample size, the assumed genotypic relative risks (GRR), and allele frequencies. They are also useful for evaluating type I error rates for new statistical association tests and power comparisons between the new tests and other existing tests. *GeneEvolve* can simulate several phenotypes, simultaneously, based on GWAS panels (e.g., Illumina and Affymetrix) or sequence data (UK10K, 1000 genome, and etc.).

### 2.4.1 Genotypic value for single locus

If we could replicate a particular genotype in a number of individuals and measure them under environmental conditions normal for the population, their mean environmental deviation would be zero, and their mean phenotypic value would consequently be equal to the genotypic value of that particular genotype. This is the meaning of the genotypic value of an individual. In principle it is measurable, but in practice it is not, except when we are concerned with a single single locus where the genotypes are phenotypically distinguishable, or with the genotypes represented in highly inbred lines [1].

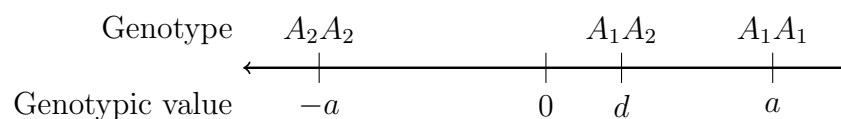


Figure 2.4: Arbitrary assigned genotypic value.

For the purposes of deduction we must assign arbitrary values to the genotypes under discussion. This is done in the following way. Considering a single locus with two alleles  $A_1$  and  $A_2$ , we call the genotypic value of one homozygote  $+a$ , that of the other homozygote  $-a$  and that of the heterozygote  $d$ . (We shall adopt the convention that  $A_1$  is the allele that increase the the value). We thus have a scale of genotypic values as in figure 2.4. The origin, or point of zero value, on this scale is mid-way between the values of the two homozygotes. The value  $d$  of the heterozygote depends on the degree of dominance. If there is no dominance,  $d = 0$ ; if  $A_1$  is dominance over  $A_2$ ,  $d$  is positive, and if  $A_2$  is dominant over  $A_1$ ,  $d$  is negative. If dominance is complete,  $d$  is equal to  $+a$  or  $-a$ , and if there is overdominance,  $d$  is greater than  $+a$  or less than  $-a$ . The degree of dominance may be expressed as  $d/a$ . The deviation from population mean is summarized in Table 2.1.

Table 2.1: Values of genotypes in a two-allele system, measured as deviations from the population mean. The population mean is  $a(p - q) + 2pqd$ , and  $\alpha = a + d(q - p)$ .

	Genotypes		
	$A_1A_1$	$A_1A_2$	$A_2A_2$
Frequencies	$p^2$	$2pq$	$q^2$
Assigned values	$a$	$d$	$-a$
Deviation from population mean	$2q(a - pd)$	$a(q - p) + d(1 - 2pq)$	$-2p(a + qd)$
or	$2q(\alpha - qd)$	$\alpha(q - p) + 2pqd$	$-2p(\alpha + pd)$
Breeding value	$2q\alpha$	$(q - p)\alpha$	$-2p\alpha$
Dominance deviation	$-2q^2d$	$2pqd$	$-2p^2d$

## 2.4.2 One phenotype

For simplicity, we assume there is one phenotype. For several phenotypes, the idea is the same. For each individual  $i$ , the phenotypic value  $P_i$  is a random variable defined as

$$P_i = A_i + D_i + F_i + E_i + C_f + \gamma_p \quad (2.1)$$

where  $A_i$  and  $D_i$  are additive and dominance genetic terms. The terms  $F_i$ ,  $E_i$ ,  $C_f$ , and  $\gamma_p$  are familial, unique, shared sibling (common), and population specific environmental effects, respectively.

For familial, shared sibling (common), and population specific environmental effects, see Sections 2.7, 2.8 and 2.9, respectively.

To simulate a phenotypes, user can specify  $m$  causal variants (CVs),  $\{cv_j\}_{j=1}^m$ , and their additive ( $a_j$ ) and dominance ( $d_j$ ) effects. The value  $a_j$  and  $d_j$  are defined in Section 2.4.1. Typically,  $a_j$  and  $d_j$  are chosen such that  $a_j \sim N(0, \sigma_g^2)$  and  $d_j \sim N(0, \sigma_d^2)$ .

For each individual  $i$ , the additive term  $A_i$  is a linear function of some causal variants multiplied by their additive effect size ( $\alpha_j$ ), i.e.,

$$A_i = \sum_{j=1}^m \frac{(x_{ij} - 2p_j)}{\sqrt{2p_jq_j}} \alpha_j, \quad (2.2)$$

where

$$\alpha_j = a_j + d_j(q_j - p_j),$$

is the additive effect size,  $p_j$  is frequency of one allele with  $q_j = 1 - p_j$  the frequency of the other,  $x_{ij} \in \{0, 1, 2\}$  is the genotypic value for  $cv_j$ , and  $m$  is the number of CVs. See Figure 2.5.

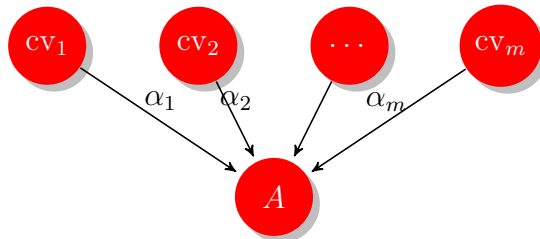


Figure 2.5: Additive term



For each individual  $i$ , the dominance term  $D_i$  is,

$$D_i = \sum_{j=1}^m t_{ij} \delta_j, \quad (2.3)$$

where

$$\delta_j = \frac{d_j}{2p_j q_j},$$

is the dominance effect size, and

$$t_{ij} = \begin{cases} -2p_j^2 & \text{if } x_{ij} = 0 \\ 2p_j q_j & \text{if } x_{ij} = 1 \\ -2q_j^2 & \text{if } x_{ij} = 2 \end{cases}.$$

For more information, see [1, 5].

Assuming no LD between CVs, for the additive term we have

$$\begin{aligned} \text{var}[A_i] &= \sum_{j=1}^m \alpha_j^2 \text{var}\left[\frac{x_{ij} - 2p_j}{\sqrt{2p_j q_j}}\right] \\ &= \sum_{j=1}^m [a_j + d_j(q_j - p_j)]^2, \end{aligned} \quad (2.4)$$

since  $\text{var}[x_{ij}] = 2p_j q_j$ .

For the dominance term,

$$\mathbb{E}[t_{ij}] = -2p_j^2 \times q_j^2 + 2p_j q_j \times 2p_j q_j - 2q_j^2 \times p_j^2 = 0,$$

and

$$\begin{aligned} \text{var}[t_{ij}] &= (-2p_j^2)^2 \times q_j^2 + (2p_j q_j)^2 \times 2p_j q_j + (-2q_j^2)^2 \times p_j^2 \\ &= 4p_j^4 \times q_j^2 + 4p_j^2 q_j^2 \times 2p_j q_j + 4q_j^4 \times p_j^2 \\ &= 4p_j^2 q_j^2 [p_j^2 + 2p_j q_j + 2q_j^2] \\ &= 4p_j^2 q_j^2. \end{aligned}$$

Therefore,

$$\begin{aligned} \text{var}[D_i] &= \sum_{j=1}^m \delta_j^2 \text{var}[t_{ij}] \\ &= \sum_{j=1}^m \frac{d_j^2}{4p_j^2 q_j^2} 4p_j^2 q_j^2 \\ &= \sum_{j=1}^m d_j^2. \end{aligned} \quad (2.5)$$

### 2.4.3 Heritability

Heritability estimates how much of the phenotypic variation can be explained by genetic, or genetic-environmental effects. Broad-sense heritability ( $H^2$ ) refers to the inclusion of all potential sources of genetic variation (additive, dominance, epistatic, maternal and paternal effects):

$$H^2 = \frac{V_G}{V_P} = \frac{V_A + V_D}{V_P},$$

where  $V_A := \text{var}[A]$ ,  $V_D := \text{var}[D]$  and so on.

To only know the ratio of additive genetic variation to the total phenotypic variation observed,  $V_A$  can be used in the equation instead of  $V_G$ , and this becomes the narrow-sense heritability ( $h^2$ ):

$$h^2 = \frac{V_A}{V_P}.$$

#### 2.4.4 Phenotypes in *GeneEvolve*

To simulate a phenotype, you should specify the parameters

```
--file_cv_info [cv_info.txt] --file_cvs [cvs.txt]
```

The file `cv_info.txt` contains CV information and has a header with 5 columns: chromosome number, base-pair distance, MAF, additive ( $a$ ) and dominance ( $d$ ) effects (see Example 3).

The file `--cvs.txt` has no header and contains the address of haplotype files contain the CVs. This file has just 2 columns: chromosome number and address of CV haplotype file (see Example 4).

**Example 3** (CV information file format). The file format for `--file_cv_info [cv_info.txt]` is illustrated in the following listing.

```
1 chr pos maf a d
2 1 4328476 0.189021 1.18585978544321 -1.54676435875486
3 1 4500436 0.450922 -0.128733310867769 1.11559331900794
4 1 7736097 0.359727 -0.379038685626258 -1.12984403982323
5 1 14418448 0.35875 0.0959318783552277 0.527555965661557
6 1 15825195 0.303792 1.89998223576411 -0.724345249882114
7 1 17889690 0.0623102 -0.368424537129164 -0.120587409943792
8 ...
9 22 42504679 0.335594 0.990334634266996 -0.698799892553054
10 22 44338134 0.418218 1.26982516250768 0.990334634266996
11 22 47450911 0.487309 1.8384670663176 -0.625110331654888
12 22 49411595 0.413211 1.44121811394195 0.878797016360936
```

cv\_info.txt

**Example 4** (CVs file format). The file format for `--file_cvs [cvs.txt]` is illustrated in the following listing.

```
1 1 /path/CVs.chr1.hap
2 2 /path/CVs.chr2.hap
3 ...
4 22 /path/CVs.chr22.hap
```

cvs.txt

#### 2.4.5 Scaling VA and VD

In Equations 2.4 and 2.5 we compute the additive and dominance variances. It is possible in *GeneEvolve* that user scale them to any arbitrary positive number using the following parameters.

```
--va [number] --vd [number]
```

If there are more than one phenotype, the user can add more scaling parameters.

### 2.4.6 Simulating multivariate phenotypes

To simulate  $k$  phenotypes (multivariate phenotypes) with different CVs, user should specify the parameters `--file_cv_info [cv_info.txt]` and `--file_cvs [cvs.txt]` in the command line,  $k$  times.

**Example 5** (Simulating 3 phenotypes). In the following listing, *GeneEvolve* will create 3 phenotypes, where their CVs are listed in different files.

```

1 GeneEvolve --file_gen_info gen.info --file_hap_name hap_add.txt \
2 --file_recom_map map.txt \
3 --file_cv_info cv_info_p1.txt --file_cvs cvs_p1.txt \
4 --file_cv_info cv_info_p2.txt --file_cvs cvs_p2.txt \
5 --file_cv_info cv_info_p3.txt --file_cvs cvs_p3.txt

```

Simulating 3 phenotypes

## 2.5 Random and non-random mating systems

Mating and reproductive systems affect the way that alleles are combined in individuals in a population. Outcrossing organisms put together new combinations of genes rapidly, leading to many different genotypes within populations (and creating high genotype diversity) and the potential for rapid adaptation in a changing environment.

For random mating, user should use the parameter `--RM`. If it is the case, then the second column in file `file_generaions_info.txt` inputted in parameter `--file_gen_info [file_generaions_info.txt]` will not be used. In random mating, offsprings will choose their parents randomly.

On the other hand, in assortative mating (mating based on some phenotype), we will create a random variable called mating value,  $MV$ , by combing all the  $k$  phenotypes as

$$MV_i = \sum_{j=1}^k \omega_j P_{ij}, \quad (2.6)$$

where  $P_{ij}$  is the  $j$ th phenotype for individual  $i$  and  $\omega_j$ 's are some coefficients (see Figure 2.6).

Mates will chose their spouses based on the mating values,  $MV$ . User can specify the correlation (which can change across time) between mates in the second column of file `file_generaions_info.txt`. Their correlation is

$$\mu = \text{corr}[MV_{FA}, MV_{MO}]. \quad (2.7)$$

## 2.6 Natural selection

*Natural selection* is the differential survival and reproduction of individuals due to differences in phenotype. It is a key mechanism of evolution, the change in heritable traits of a population over time. Natural variation occurs among the individuals of any population of organisms. Many of these differences do not affect survival or reproduction, but some differences may improve the chances of survival and reproduction of a particular individual.

Natural selection is a process that favors or induces survival and perpetuation of one kind of organism over others. Selection can be positive (or advantageous) or negative (or purifying) and has a profound impact on the evolution of the human population. In addition, selection can be balancing in which the

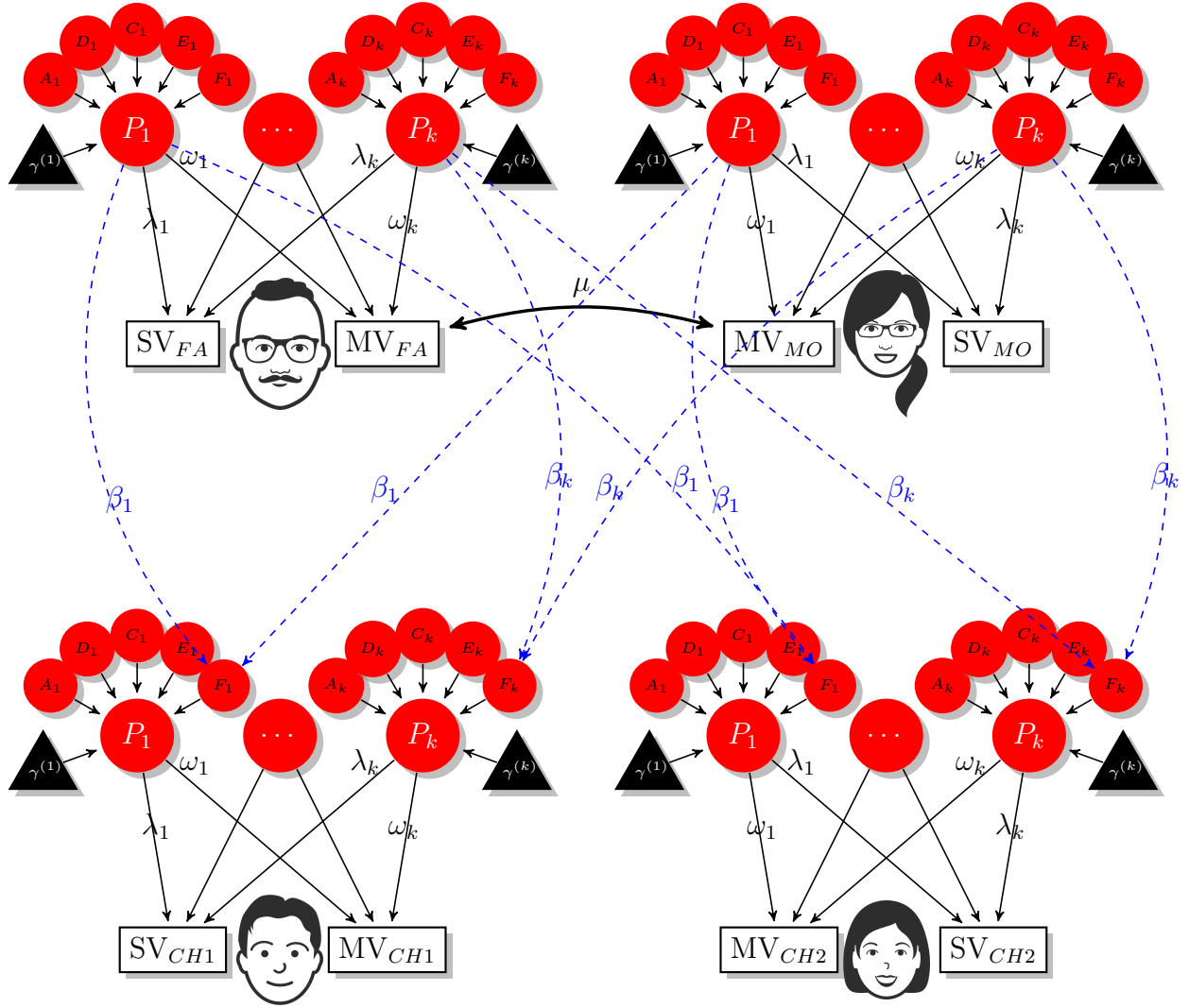


Figure 2.6: Mating path diagram and phenotypes

genotypes have a mixture of positive and negative selection pressures so that there is no net effect of selection on the individual alleles.

For each individual  $i$ , the selection value can be computed from  $k$  phenotypes as

$$SV_i = \sum_{j=1}^k \lambda_j P_{ij}, \quad (2.8)$$

where  $P_{ij}$  is the  $j$ th phenotype for individual  $i$  and  $\lambda_j$ 's are some user define coefficients (see Figure 2.6).

### 2.6.1 Directional selection

In population genetics, *directional selection* is a mode of natural selection in which an extreme phenotype is favored over other phenotypes, causing the allele frequency to shift over time in the direction of that phenotype. Under directional selection, the advantageous allele increases as a consequence of differences in survival and reproduction among different phenotypes. The increases are independent of the dominance of the allele, and even if the allele is recessive, it will eventually become fixed.

Based on the individuals selection value and selection function, *GeneEvolve* decides to let them marry or not.

User can define the following parameters as selection function in file `file_generaions_info.txt`

```
1 logit p1 p2
2 probit p1 p2
3 stab p1 p2
4 thr p1 p2
```

where  $p1$  and  $p2$  are its parameters. For more information about `file_generaions_info.txt`, see Figure 2.1.

For the `logit` function, the probability of mating can be computed from inverse *logit* function, i.e.,

$$\mathbb{P}[\text{mating}] = \frac{\exp(p_1 + p_2 \text{SV}_i)}{1 + \exp(p_1 + p_2 \text{SV}_i)}. \quad (2.9)$$

For the `probit` function, the probability of mating can be computed from inverse *probit* function, i.e., normal CDF with mean  $p_1$  and standard deviation  $p_2$ .

### 2.6.2 Stabilizing selection

For the `stab` function (stabilizing selection), the probability of mating can be computed from the normal PDF with mean  $p_1$  and standard deviation  $p_2$ .

### 2.6.3 Threshold selection

For the `thr` function (threshold selection), the probability of mating can be computed from

$$\mathbb{P}[\text{mating}] = \begin{cases} p_1 & \text{if } \text{SV}_i < p_2 \\ 1 & \text{if } \text{SV}_i \geq p_2 \end{cases}. \quad (2.10)$$

## 2.7 Familial effect

Sometimes non-genetics characters run in families and influence the phenotypes of offsprings. For example, families with higher education tend to have more educated offspring. This is also true form wealth. We can model this familial effect in *GeneEvolve*, easily.

For each offspring  $i$ , the familial effect can be computed from the following equation

$$F_i = \beta \frac{P_i(FA) + P_i(MO)}{\sqrt{2}}, \quad (2.11)$$

where  $P_i(FA)$  and  $P_i(MO)$  are parent's phenotypes (see Figure 2.6) and  $\beta$  is an arbitrary coefficient defined in `--beta [number]`. User can define the variance of familial effect in *GeneEvolve* by `--vf [number]`.

If there are more than one phenotype, user can define different  $\beta$ 's and different variances for familial effect for them.

## 2.8 Shared sibling (common) effect

Shared sibling (common) effect is defined as the environmental effects shared by groups of individuals, for example effects shared by groups of relatives that are not due to genetic effects.

For each sibling  $i$  in a family  $f$ , we add a random number  $c_f$  to its phenotype. The generated random number  $c_f$  comes from a standard Gaussian distribution with mean zero and the user specified variance VC. User can define VC in *GeneEvolve* by assigning a positive number in `--vc [VC]`.

If there are more than one phenotype, user can define different values for VC for each phenotype.

## 2.9 Environmental effects specific to each population

The aim of this subsection is to define the  $\gamma_p$  term in Equation 2.1. For simplicity, assume there is just one phenotype. For  $k$  phenotypes, the idea is the same and user should input  $\gamma^{(i)}$  for  $i \in \{1, \dots, k\}$  (see Figure 2.6).

Assume that there are  $P$  populations. For a user-specified  $\gamma$ , the environmental effects specific to a population  $p \in P$ , is  $\gamma_p$ , which can be obtained from solving the following equation

$$\text{var}(Y) = (1 + \gamma)\text{var}(X), \quad (2.12)$$

where,  $X$  is all the phenotypes obtained from the combined populations, i.e.

$$Y = \bigcup_{p \in P} \bigcup_{i \in \text{pop}(p)} S_{i,p},$$

and

$$X = \bigcup_{p \in P} \bigcup_{i \in \text{pop}(p)} T_{i,p},$$

where for each individual  $i$  in population  $p$ ,

$$S_{i,p} = A_i + D_i + F_i + E_i + C_f + b_p, \quad (2.13)$$

$$T_{i,p} = A_i + D_i + F_i + E_i + C_f, \quad (2.14)$$

where

$$b_p = \Gamma \left( \frac{2(p-1)}{P-1} - 1 \right).$$

After solving for  $\Gamma$ , the environmental effects specific to population  $p$ , becomes

$$\gamma_p = \Gamma \left( \frac{2(p-1)}{P-1} - 1 \right). \quad (2.15)$$

## 2.10 Simulating several populations

In *GeneEvolve* you can easily simulate several population. For simulating the second population, user should use the parameter `--next_population` in order to distinguish between populations parameters. There is no limit in the number of populations, but you should use the parameter `--next_population` to separate them. See Chapter 3 for more information.

## 2.11 Population structure

### 2.11.1 Introduction

A population may have substructures – different in genetic variation among its constituent parts – for several different evolutionary reasons. Exchange of individuals may not have equal probabilities throughout a population, or selection may have different effects in different parts of population. To model it, assume that a population consists of  $p$  subpopulations and that the proportion of individuals migrating from subpopulation  $j$  to subpopulation  $i$  each generation is  $m_{ij}$ . As a result there is a matrix of gene flow parameters, called the backward *migration matrix*, that describes the gene flow pattern among subpopulations (see Table 2.2). The proportion of non-migrants (or residents) for subpopulation  $i$  is given by  $m_{ii}$  (see figure 2.7). Each row of this matrix sums to unity because it describes the proportion coming from every other possible subpopulation to that particular subpopulation or

$$\sum_{j=1}^p m_{ij} = 1.$$

The columns of the matrix will generally not sum to unity. The migration matrix is denoted by

$$\mathbf{M} = \begin{bmatrix} m_{11} & m_{12} & \dots & m_{1p} \\ m_{21} & m_{22} & \dots & m_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ m_{p1} & m_{p2} & \dots & m_{pp} \end{bmatrix}.$$

Each of the subpopulations may have a different frequency of  $A_2$ , and let us indicate the allele frequency in the  $j$ th subpopulation as  $q_j$ . Therefore, the frequency of  $A_2$  in the  $i$ th subpopulation after gene flow is

$$q'_i = \sum_{j=1}^p m_{ij} q_j.$$

We can symbolize the process of allele frequency change over all the subpopulations by using matrix notation. First we can indicate the migration matrix as  $\mathbf{M}$  and the vector of allele frequency for the different subpopulations in generation  $t$  with  $\mathbf{Q}_t$ . Therefore,

$$\mathbf{Q}_{t+1} = \mathbf{M}\mathbf{Q}_t.$$

It is also possible that the migration matrix  $\mathbf{M}$  changes over generation, so we denote it by  $\mathbf{M}_t$ .

$$\mathbf{Q}_{t+1} = \mathbf{M}_t \mathbf{Q}_t.$$

### 2.11.2 Population structure in *GeneEvolve*

The migration matrix can be inputted to *GeneEvolve* using the parameter

```
--file_migration [file_migration.txt]
```

The inputted file has no header and it has  $n$  rows, where  $n$  is the number of generations. This file should have  $k^2$  columns as follows:

$$m_{11}, m_{12}, \dots, m_{1k}, m_{21}, \dots, m_{2k}, \dots, m_{k1}, \dots, m_{kk}$$

Therefore, the user can use different migration matrix  $\mathbf{M}_t$  for each generation.

Table 2.2: The migration matrix

Subpopulations in generation $t + 1$	Subpopulations in generation $t$				Total
	1	2	...	$p$	
1	$m_{11}$	$m_{12}$		$m_{1p}$	1
2	$m_{21}$	$m_{22}$		$m_{2p}$	1
$\vdots$	$\vdots$	$\vdots$		$\vdots$	$\vdots$
$p$	$m_{p1}$	$m_{p2}$		$m_{pp}$	1

**Example 6** (Migration file format). In the following listing file, there are 10 generations with 2 subpopulations.

```

1 1 .0 .05 .95
2 .9 .1 .05 .95
3 .8 .2 .05 .95
4 .9 .1 .05 .95
5 .9 .1 .05 .95
6 .9 .1 .05 .95
7 .9 .1 .05 .95
8 .9 .1 .05 .95
9 .9 .1 .05 .95
10 .9 .1 .05 .95

```

file\_migration.txt

For generation 3, we have

$$\mathbf{M}_3 = \begin{bmatrix} .8 & .2 \\ .05 & .95 \end{bmatrix}$$



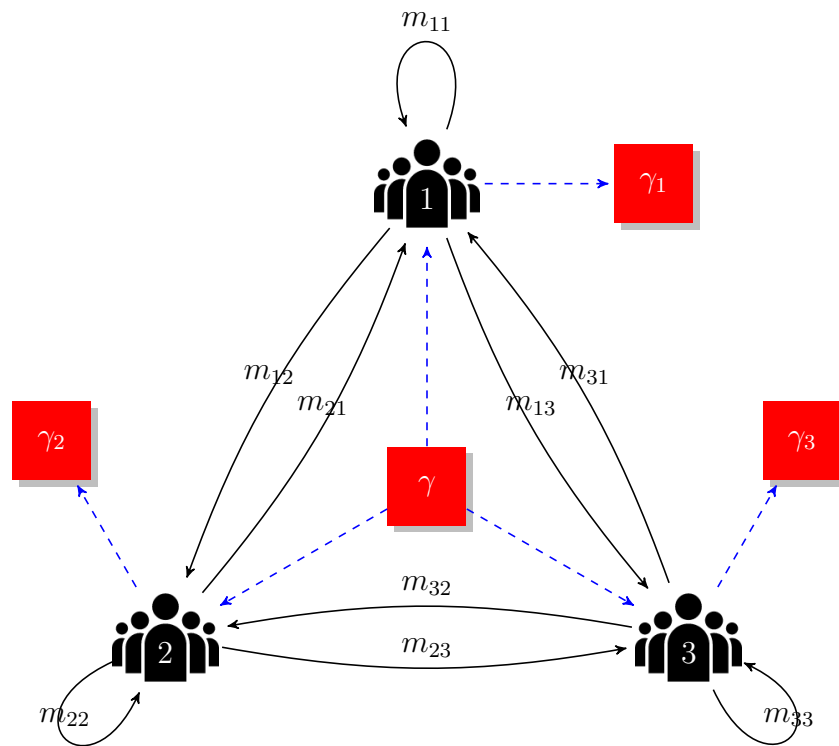


Figure 2.7: Path diagram for migration and environmental effects specific to each population



# Chapter 3

## Parameters Reference

Here is the list of all the parameters used in *GeneEvolve*.

```
1 GeneEvolve \  
2 --file_gen_info par.pop1.info.txt \  
3 --file_hap_name par.pop1.hap_sample_address.txt \  
4 --file_recom_map Recom.Map.b37.50KbDiff \  
5 --file_mutation_map mutation.Map.b37.50KbDiff \  
6 --file_cv_info par.pop1.cv1_info.txt \  
7 --file_cvs par.pop1.cv1_hap_files.txt \  
8 --va 1 \  
9 --vd .2 \  
10 --vc .1 \  
11 --ve 1 \  
12 --vf .1 \  
13 --omega .7 \  
14 --beta .7 \  
15 --lambda 1 \  
16 --file_cv_info par.pop1.cv2_info.txt \  
17 --file_cvs par.pop1.cv2_hap_files.txt \  
18 --va 2 \  
19 --va .1 \  
20 --vc .1 \  
21 --ve 1 \  
22 --vf .1 \  
23 --omega .3 \  
24 --beta .3 \  
25 --lambda 1 \  
26 --next_population \  
27 --file_gen_info par.pop2.info.txt \  
28 --file_hap_name par.pop2.hap_sample_address.txt \  
29 --file_cv_info par.pop2.cv1_info.txt \  
30 --file_cvs par.pop2.cv1_hap_files.txt \  
31 --file_cv_info par.pop2.cv2_info.txt \  
32 --file_cvs par.pop2.cv2_hap_files.txt \  
33 --file_recom_map Recom.Map.b37.50KbDiff \  
34 --file_mutation_map mutation.Map.b37.50KbDiff \  
35 --va 2 \  
36 --vd .2 \  
37 --vc 1 \  

```

```

38 --ve 1 \
39 --vf 0 \
40 --va 2 \
41 --vd .2 \
42 --vc 1 \
43 --ve 1 \
44 --vf 0 \
45 --omega 1 \
46 --omega 1 \
47 --beta 0 \
48 --beta 0 \
49 --lambda 1 \
50 --lambda 1 \
51 --file_migration par.migration.txt \
52 --prefix out1 \
53 --gamma 1 \
54 --gamma 1\
55 --avoid_inbreeding

```

The above listing we simulate 2 populations, each with 2 phenotypes. In the following subsections, we explain each of the parameters, briefly.

### 3.1 --file\_gen\_info

```
--file_gen_info par.pop1.generaions_info.txt
```

where `par.pop1.generaions_info.txt` is a text file with 6 columns and  $n + 1$  lines, where  $n$  is the number of generations to be simulated by *GeneEvolve*. This file should have a header and the first column is the population size. The structure of this file is listed in Figure 2.1. *GeneEvolve* can simulate different population sizes in each generation, by specifying some numbers in the first column.

### 3.2 --file\_hap\_name

```
--file_hap_name par.pop1.hap1_sample_address.txt
```

where `par.pop1.hap1_sample_address.txt` is a text file with 4 columns, counting the address of the initial hap, legend and sample files.

### 3.3 --file\_recom\_map

```
--file_recom_map Recom.Map.b37.50KbDiff
```

where `Recom.Map.b37.50KbDiff` is a text file with 4 columns, counting the cM distance for all the snap panel.

### 3.4 --file\_cv\_info

--file\_cv\_info par.pop1.cv1\_info.txt

where the CV information are in file par.pop1.cv1\_info.txt.

### 3.5 --file\_cvs

--file\_cvs par.pop1.cv1\_hap\_files.txt

where the actual haplotype address are saved in file par.pop1.cv1\_info.txt.

### 3.6 --va

--va 2

*GeneEvolve* will transforms the additive variance to 2, in generation zero.

### 3.7 --vd

--vd 0.1

*GeneEvolve* will transforms the dominance variance to 0.1, in generation zero.

### 3.8 --vc

--vc 0.1

*GeneEvolve* will transforms the variance of sibling environment (common) effect to 0.1, in all generations.

### 3.9 --ve

--ve 1

*GeneEvolve* will transforms the environmental effect variance to 1, in all generations.

### 3.10 --vf

--vf 0.1

*GeneEvolve* will transforms the familial effect variance to 0.1, in generation zero.

### 3.11 --next\_population

--next\_population

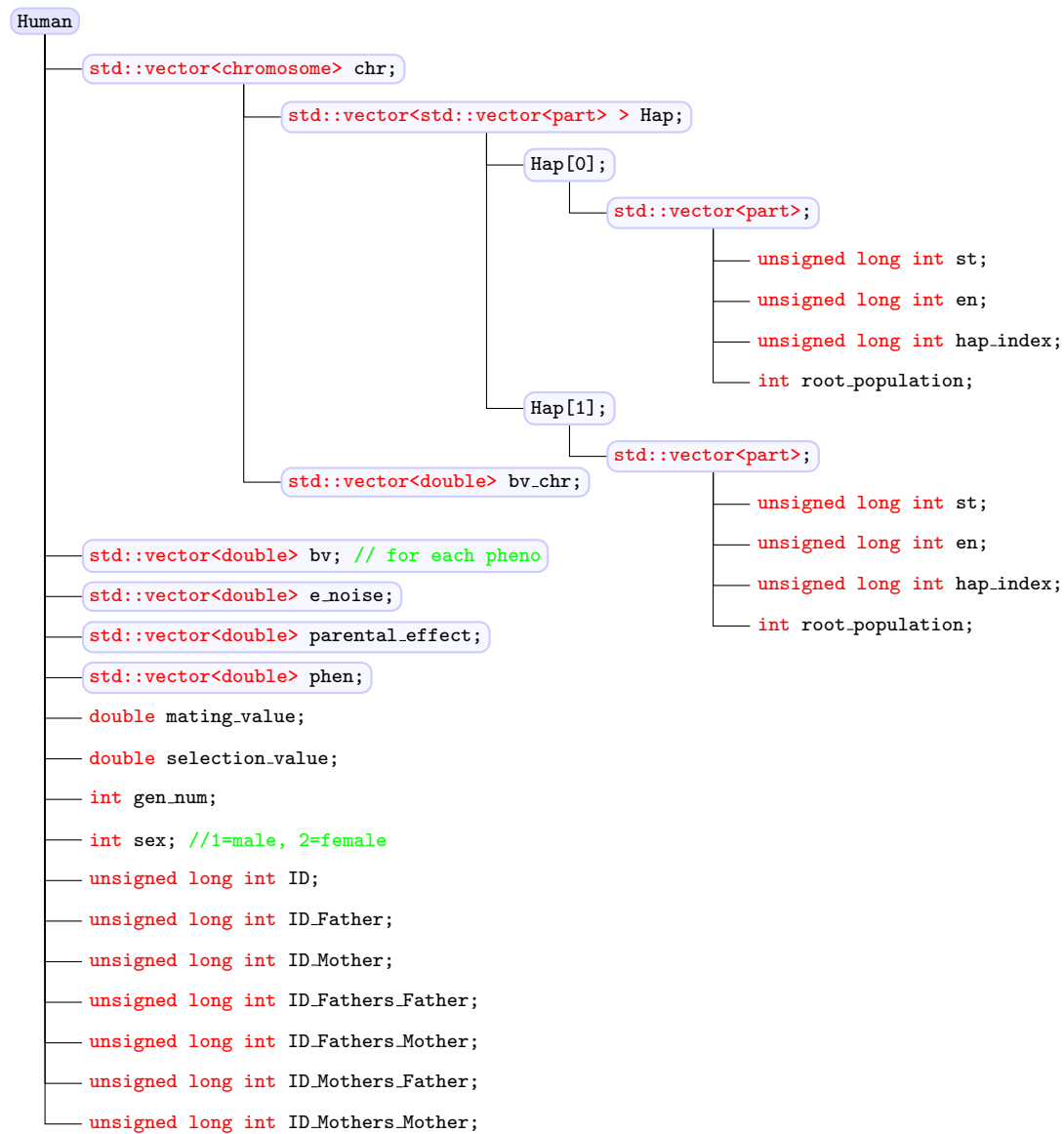
This parameter will used to give the next population information.



# Appendix A

## C++ Classes

### A.1 Human







## A.2 Population

### Population

```

int _pop_num;
int _nchr;
std::vector<Phenotype_scheme> _pheno_scheme; // for each phenotype
    std::vector<CV_INFO> _cv_info; // for each chr
        std::vector<unsigned long int> bp;
        std::vector<double> maf;
        std::vector<double> alpha;
    std::vector<CV> _cvs; // for each chr
        std::vector<std::vector<bool> > val;
    std::vector<std::string> _name_cv_hap; // for each chr
    double _va;
    double _ve;
    double _vf;
    double _alpha; // mating value coefficient
    double _beta; // transmission of environmental effects from parents to offspring
    double _delta; // selection value coefficient
std::vector<Human> h;
std::vector<Couples_Info> _couples_info;
    unsigned long int pos_male; // pos human, not pos hap
    unsigned long int pos_female; // pos human, not pos hap
    bool inbreed;
    int num_offspring;

std::vector<unsigned long int> _pop_size; // for each generation
std::vector<double> _mat_cor; // for each generation
std::vector<std::string> _offspring_dist; // for each generation
std::string _selection_func; // --logit 0 1
std::vector<double> _selection_func_pars; // --logit 0 1
std::vector<std::vector<std::string> > _hap_legend_sample_name; // for each chr, with 3 columns
std::vector<rMap> _rmap; // for each chr
std::vector<std::vector<double> > _recom_prob; // for each chr
std::vector<double> _var_bv_gen0; // for each phenotype
std::vector<int> _all_active_chrs; // for each chr
bool _avoid_inbreeding;
bool _no_output;
bool _output_all_generations;
bool _debug;
std::string _out_prefix;
std::string _format_output;
double _RM_percent; // Random mating percent (inds who have 2 spouses)
std::vector<double> ret_var_mating_value; // for each gen
std::vector<std::vector<double> > ret_var_phen; // for each pheno and gen
std::vector<std::vector<double> > ret_var_bv; // for each pheno and gen
std::vector<std::vector<double> > ret_var_parental_effect; // for each pheno and gen

```

## A.3 Simulation

### Simulation

```

int _n_pop;
int _tot_gen;
bool _debug;
std::vector<Population> population;
Parameters par;
std::vector<std::vector<double> > imigration_mat_gen;
std::string _out_prefix;
std::string _format_output;
bool _output_all_generations;
std::vector<int> _all_active_chrs;
std::vector<double> _gamma; // for each phenotype and all the populations
std::vector<Pop_phen_info> Pop_info_prev_gen; // Saving mating_value for the next generation
    std::vector<double> mating_value; // for each ind
    std::vector<double> selection_value; // for each ind
    std::vector<std::vector<double> > phen; // for each phen and ind

```

# Appendix B

## File formats

### B.1 Hap – Legend – Sample

#### B.1.1 .hap

No header.

This file is SPACE delimited. Each line corresponds to a single SNP. Each successive column pair (0, 1), (2, 3), (4, 5) and (6, 7) corresponds to the alleles carried at the 4 SNPs by each haplotype of a single individual. For example a pair "1 0" means that the first haplotype carries the B allele while the second carries the A allele as specified in the LEGEND file. The haplotypes are given in the same order than in the SAMPLE file. This file should have L lines and 2N columns, where L and N are the numbers of SNPs and individuals respectively.

	ind1.hap1	ind1.hap2	ind2.hap1	ind2.hap2	...	
snp1	0	0	1	0		
snp2	0	1	1	0		
⋮						
						nsnp × nhaps

#### B.1.2 .legend

Has header.

This file is SPACE delimited. The first line is a header line that describe the content of the file. Each line corresponds to a single SNP.

	col1	col2	col3	col4	
header	ID	pos	allele0	allele1	
snp1	rs17432784	17196300	T	C	
snp2	rs2845379	rs1807512	A	G	
snp3	rs17432784	17196300	G	T	
⋮					
					(nsnp+1) × 4

### B.1.3 .sample

Has header.

It is SPACE delimited. The first line is a header line that describe the content of the file. Then, each line corresponds to a single individual.

	col1	col2	col3	col4	
header	sample	population	group	sex	
ind1	CEU1	CEU	EUR	1	
ind2	CEU2	CEU	EUR	2	
ind3	GBR1	GBR	EUR	2	
⋮					
					$(nind+1) \times 4$

### B.1.4 .impute.hap.indv

No header.

This file has just one column.

	col1	
ind1	5659883013-R02C01	
ind2	5648551130-R02C01	
ind3	5648560075-R02C01	
⋮		
		$nind \times 1$

## B.2 PLINK

There are two formats for PLINK: Binary format (.bed, .bim and .fam) or uncompressed format (.ped and .map).

### B.2.1 .ped

No header.

Each line corresponds to a single individual.

1. Family ID
2. Sample ID
3. Paternal ID
4. Maternal ID
5. Sex (1=male; 2=female; other=unknown)
6. Affection (0=unknown; 1=unaffected; 2=affected)
7. Genotypes (space or tab separated, 2 for each marker. 0=missing)

	FID	IID	PID	MID	Sex	Aff	SNP1	SNP1	SNP2	SNP2	
IND1	fam1	ind1	0	0	1	2	A	C	C	T	
IND2	fam1	ind2	0	0	2	1	C	A	T	T	
IND2	fam2	ind1	0	0	1	1	C	C	C	T	
⋮											
											$(nind) \times (6+2*nsnp)$

### B.2.2 .map

No header.

It is SPACE delimited. Each line corresponds to a single SNP. Chromosome can be (1-22, X, Y or 0 if unplaced)

	chromosome	rs#	cM	bp	
snp1	1	rs123456	0	1234555	
snp2	1	rs234567	0	1237793	
snp3	1	rs233556	0	1337456	
⋮					
					$(nsnp) \times 4$

### B.2.3 .fam

No header.

It is SPACE delimited. Each line corresponds to a single individual.

1. Family ID ('FID')
2. Within-family ID ('IID'; cannot be '0')
3. Within-family ID of father ('0' if father isn't in dataset)
4. Within-family ID of mother ('0' if mother isn't in dataset)
5. Sex code ('1' = male, '2' = female, '0' = unknown)
6. Phenotype value ('1' = control, '2' = case, '-9'/'0'/non-numeric = missing data if case/control)

If there are any numeric phenotype values other than -9, 0, 1, 2, the phenotype is interpreted as a quantitative trait instead of case/control status. In this case, -9 normally still designates a missing phenotype;

	FID	IID	ID-Father	ID-mother	Sex	Phenotype	
ind1	1	child1	0	0	1	1	
ind2	1	child2	0	0	1	2	
ind3	2	child1	0	0	1	2	
⋮							
							$(nind) \times 6$



# Bibliography

- [1] DS Falconer and TFC Mackay. *Introduction to Quantitative Genetics*. Longman, 4 edition, 1996.
- [2] John FC Kingman. The coalescent. *Stochastic processes and their applications*, 13(3):235–248, 1982.
- [3] John FC Kingman. On the genealogy of large populations. *Journal of Applied Probability*, pages 27–43, 1982.
- [4] Augustine Kong, Daniel F Gudbjartsson, Jesus Sainz, Gudrun M Jonsdottir, Sigurjon A Gudjonsson, Bjorgvin Richardsson, Sigrun Sigurdardottir, John Barnard, Bjorn Hallbeck, Gisli Masson, et al. A high-resolution recombination map of the human genome. *Nature genetics*, 31(3):241–247, 2002.
- [5] Zulma G Vitezica, Luis Varona, and Andres Legarra. On the additive and dominant variance and covariance of individuals within the genomic selection scope. *Genetics*, 195(4):1223–1230, 2013.