

December 17, 2016

cophesim
User Manual

Ilya Y. Zhbannikov¹, Konstantin G. Arbeev¹, Anatoliy I. Yashin¹

¹Biodemography of Aging Research Unit (BARU) at Social Sciences Research Institute (SSRI),
Duke University, Durham, NC,

Keywords:

Bioinformatics, Genomics, Data Simulation, Artificial Data, Synthetic Data

Correspondence:

Ilya Y. Zhbannikov
Biodemography of Aging Research Unit (BARU) at Social Sciences Research Institute (SSRI)
Duke University
Durham, NC, 27705
Email: ilya.zhbannikov@duke.edu

Running Title:

cophesim: a user manual

1 Introduction

`cophesim` - a comprehensive phenotype simulator for genetic data, i.e. `cophesim` adds phenotype to provided genotype files.

2 Installation

2.1 Prerequisites

- Python v2.7.10
- plinkio v0.9.6
- R v3.2.4 (for the examples)
- PLINK v1.7 (for the examples)

2.2 How to install

`cophesim` is an open-source software application available from the Bitbucket for free under this link: <http://bitbucket.org/izhbannikov/cophesim>.

1. Install prerequisites: Python, plinkio (if you were not able to install plinkio, just skip it. The “cophesim” will still work, but not able to handle binary (.bed, .bim, .fam) PLINK files, only .ped and .map files will be handled), R, PLINK.
2. Download or clone the `cophesim` repository: <https://bitbucket.org/izhbannikov/cophesim>. Save under some name you wish and unzip. The software is ready to use.
3. Additionally, you may download the data repository: https://bitbucket.org/izhbannikov/cophesim_data. It provides some simulated data and code examples.

21 Save the file under some name you wish and unzip. The software is ready to use.

22 **2.3 Usage**

23 `cophesim.py -i <path to genotype> -o <output prefix> [options]`

24 **2.4 Options**

25 Input option described in Table 1

26 **2.5 Description of input files**

27 Input genotype data can be in one of the formats generated from the following applications: Plink
28 (.bed, .bim, .fam); ms, msms, msHot (plain text file); Genome (plain text file). Plink format is used
29 by default. In this case you have to provide a path to the files (i.e. full prefix without file extension).

30 **2.6 Description of output**

31 `cophesim` generates the following output files:

32 1. Phenotype file. This file is in text format and has the following suffices depending on the
33 simulated phenotype trait: `_pheno_bin.txt`, `_pheno_cont.txt`, `_pheno_surv.txt`
34 representing dichotomous (binary), quantitative (continuous) and survival phenotype. Below
35 we show description of columns.

36 `_pheno_bin.txt`:

- 37 • Individual ID
- 38 • Sex
- 39 • Phenotype

Table 1: Input options

Option	Extended option	Description
-h	-help	Show the help message and exit.
-i IDATA	-input IDATA	Path input file(s). Extension should not be used in itype = plink.
-o OUTPUT_PREFIX	-output OUTPUT_PREFIX	Output prefix.
-itype ITYPE		Input format: <code>plink</code> (for Plink, default), <code>ms</code> (for <code>ms</code> , <code>msms</code> , <code>msHot</code>), <code>genome</code> (for Genome).
-otype OTYPE		Indicates output format, by default <code>OTYPE=plink</code> . Other possible output format: <code>blossoc</code> (for BLOSSOC), <code>qtdt</code> (for QTDT), <code>tassel</code> (for Tassel), <code>emmax</code> (for EMMAX).
-d	-dichotomous	A flag for dichotomous phenotype, True by default.
-c	-continuous	A flag for continuous phenotype, False by default.
-s	-survival	A flag to simulate survival phenotype, False by default.
-ce CEFF		A path to the file with effect of each causal SNP. Must be in format: <code>snp_index:effect</code> . One snp per line.
-alpha ALPHA		An 'alpha' parameter for inverse probability equation for the Gompertz hazard (see Bender at al., Generating survival times to simulate Cox proportional hazards models), 2005. Default ALPHA = 0.2138
-epi EPIFILE		File with interacting SNPs. One pair per line. Format: <code>snp1_index, snp2_index, effect</code>
-weib		A flag to use Weibull distribution for survival phenotype. True by default.
-gomp		A flag to use Gompertz distribution for survival phenotype. False by default.
-LD LDFILE		File with collinear SNPs. One pair per line. Format: <code>snp1_index, snp2_index, corr_cf([-1,1])</code>

40 `_pheno_cont.txt:`

- 41 • Individual ID
- 42 • Sex
- 43 • Phenotype

44 `_pheno_surv.txt:`

- 45 • Individual ID
- 46 • Sex
- 47 • Age (phenotype)
- 48 • Case

49 2. Genotype file(s). Can be in the following formats: Plink (`.bed`, `.bim`, `.fam`).

50 Other possible output format: `blossoc` (for BLOSSOC, suffices `.blossoc_pos`,
51 `.blossoc_geno`), `qtdt` (for QTDT, suffices `.ped`, `.map`, `.dat`), `tassel` (for Tassel,
52 suffices `.poly`, `.trait`), `emmax` (for EMMAX, suffices `.emma_geno`, `.emma_pheno`).

53 3. Summary statistics file. This is a plain text file which keeps the information about the run.

54 **2.7 Examples**

55 Below we show several examples of usage of `cophesim`.

56 *2.7.1 Quick start*

```
57 plink --simulate-ncases 5000 --simulate-ncontrols 5000 --simulate wgas.sim \  
58 --out sim.plink --make-bed  
59  
60 python cophesim.py -i sim.plink -o testout
```

61 The first command runs the data simulation. Here we simulate genetic dataset of 10k individ-
62 uals, 5k cases and 5k controls. SNPs defined in `wgas.sim` (should be in the `cophesim` home
63 directory). Then with next command we add a phenotype (dichotomous by default) to simulated
64 genetic data.

65 To simulate continuous phenotypic trait, add the `'-c'` flag:

```
66 python cophesim.py -i sim.plink -o testout -c
```

67 This will simulate both continuous and dichotomous traits. To simulate survival trait, add `'-s'`
68 flag:

```
69 python cophesim.py -i sim.plink -o testout -s
```

70 2.7.2 Specifying causal variants

71 Causal variants are specified in the file `effects.txt` and the option `'-ce'` is used:

```
72 python /Users/ilya/Projects/cophesim/cophesim.py -i sim.plink -o testout \  
73 -ce effects.txt
```

74 The file `effects.txt` if a plain text file and causal SNPs are specified in the following format:

```
75 snp_index:effect
```

76 Here `snp_index` is the index of causal SNP and `effect` is the effect size. Example:

```
77 19:-0.82.
```

78 2.8 Specifying epistatic interactions

79 Epistatic interaction are specified in the `'epifile.txt'` with the `'-epi'` flag:

```
80 python /Users/ilya/Projects/cophesim/cophesim.py -i sim.plink -o testout \  
81 -ce effects.txt -epi epifile.txt
```

82 The file `epifile.txt` if a plain text file and a pair of interacting SNPs is specified in the following
83 format:

84 `snp1_index, snp2_index, effect`

85 Here `snp1_index` is the index of the first interacting SNP (`snp1`), `snp2_index` is the index
86 of the second interacting SNP (`snp2`) and `effect` is the corresponding effect size of this interact-
87 ing pair. Example: `12, 16, 1.57`.

88 **3 Acknowledgements**

89 This work was supported by the National Institute on Aging of the National Institutes of Health
90 (NIA/NIH) under Award Numbers P01AG043352, R01AG046860, and P30AG034424. The con-
91 tent is solely the responsibility of the authors and does not necessarily represent the official views
92 of the NIA/NIH.