

## SAFE- $h^2$ Manual

SAFE- $h^2$  is a python script designed for estimation of SNP heritability. It provides heritability profiles and thereby improves estimation accuracy of SNP heritability. SAFE- $h^2$  not only prevents the biologically meaningless negative contributions to heritability made by SNPs with large  $p$ -values (defined as “unassociated”) but also avoids the contribution of false-positive SNP hits. It can also extract intra-locus non-additive allelic effects and estimate the heritability by considering both the additive and non-additive allelic effects. This also makes it possible to perform GWAS using linear models while covering non-additive allelic effects.

### Citation:

*Behrooz Darbani, Mogens Nicolaisen. On the Genetic Origins of Phenotypes in Genome-Wide Association Studies: The SNP Allocation For Estimating Heritability (SAFE- $h^2$ ) Tool to Explore Additive-only Allelic Effects or Additive and Non-Additive Allelic Effects. Doi: <https://doi.org/10.1101/2023.08.28.555092>*

Contact: [behroozdarbani@gmail.com](mailto:behroozdarbani@gmail.com) , [bds@rsyd.dk](mailto:bds@rsyd.dk)

Department of Agroecology, Faculty of Technical Sciences, Aarhus University, Slagelse, Denmark

### SAFE-h2 benefits from two major technical features:

- I) A  $p$ -value guided SNP heritability profiling approach to include as many causal SNPs as possible in the heritability estimation models while expelling unassociated SNPs. To do this, SAFE- $h^2$  finds a  $p$ -value threshold differentiating associated and unassociated SNPs. This approach facilitates the estimation of phenotype-specific, *i.e.*, associated SNPs-specific, genetic variances. SAFE- $h^2$  also facilitates a false-positive free estimation of SNP heritability.
- II) Exploiting an allele adjustment strategy which makes it possible to include the additive effects together with different non-additive allelic effects of dominance, overdominance, and heterosis-like in the heritability estimation models. This also makes it feasible to perform GWAS using linear models while taking both additive and non-additive allelic effects into account.

### Installation requirements:

- 1) Install pip, gawk, and some python modules:
  - A) `sudo apt update` or `sudo dnf update`
  - B) `sudo apt install python3-pip` or `sudo dnf install python3-pip`
  - C) `sudo apt install gawk` or `sudo dnf -y install gawk-devel`
  - D) `pip install pandas`
  - E) `pip install numpy`
  - F) `pip install rich`
  - G) `pip install matplotlib`

- H) pip install matplotlib
- I) pip install pathlib
- J) pip install pyfiglet

- 2) Download PLINK1.9 binary file from:  
<https://www.cog-genomics.org/plink/>
- 3) Download PLINK2 binary file from:  
<https://www.cog-genomics.org/plink/2.0/>
- 4) Download GCTA-1.94.1-linux-kernel-3 binary file from:  
<https://yanglab.westlake.edu.cn/software/gcta/#Download>
- 5) Download LDAK5.2.linux binary file from:  
<https://dougsspeed.com/downloads2/>
- 6) Download EMMAX-beta-07Mar2010 binary files from:  
<http://csg.sph.umich.edu/kang/emmax/download/index.html>
- 7) Download gemma-0.98.5-linux-static-AMD64 binary file from:  
<https://github.com/genetics-statistics/GEMMA/releases/tag/v0.98.5>
- 8) Unpack the SAFE-h2 zip file.
- 9) Move the binary files (unpack if need and copy the binary file) into the SAFE-h2 folder (rename the binaries to plink, plink2, gcta, ldak, emmax, emmax-kin, and gemma).

### Usage:

- 1) Prepare and put the bfiles and the GWAS  $p$ -value file into the SAFE-h2 folder. SAFE- $h^2$  uses these files as input. The bfiles should be named as "MAIN.fam, MAIN.bed, and MAIN.bim". The  $p$ -value file should be named MainPs. It is a tab-delimited text file with three columns as:

```
SNP1_ID  2.04562  9.5E-10
SNP2_ID  -1.87003  0.002
SNP3_ID   0.23419  0.092567
...
```

The first col. is for IDs, the second can be beta coefficient, and the third col. is for  $p$ -values. The covariates can also be included in the analyses. They should also be provided as tab-delimited text files.

For EMMAX model the covariate file should be named as "Covar\_emmax".

```
Sample1  Sample1      1    2.657
Sample2  Sample2      1      3
Sample3  Sample3      1    9.342
...
```

For LDAK model the covariate file should be named as "Covar\_ldak".

```
Sample1  Sample1      2.657
Sample2  Sample2      3
```

Sample3	Sample3	9.342
---------	---------	-------

...	...	...
-----	-----	-----

For GEMMA model the covariate file should be named as “Covar\_gemma”. The order of covariates of genotypes (individuals) should follow the .fam file.

2
7
2
1
6
1
.
.
.

For GCTA-GREML model the quantitative and qualitative covariates should be provided by separate files. Please bear in mind that, for now, one of these should be provided. The quantitative covariate file should be named as “qCovar\_greml” and the qualitative covariate as “Covar\_greml”.

Sample1	Sample1	2.657
Sample2	Sample2	3
Sample3	Sample3	9.342
...	...	...

Please bear in mind that the sample IDs should be identical to the sample IDs within your fam file.

## 2) Open the terminal from the SAFE-h2 folder and run:

python3 SAFE-h2.py or python3 SAFE-h2.py | tee output

This will perform the SNP heritability profiling by implementing the EMMAX model by default. It is possible to implement other estimation models by changing the configuration file in advance. For this:

Open the configuration file (Prog.config) and define the type of analyses you are going to use as:

Impl_status_Emmax, Impl_status_LdakGCTA, Impl_status_LdakThin, Impl_status_GCTA-GREML, Impl_status_GEMMA,
Test_ADOH, PURE_Effects
X1 X2 X3 X4 X5 X6
X7

**X1-5** can get the values of 0 or 1 to ignore or implement the corresponding heritability estimation model within the analyses. **X4** can also get 2 when dealing with inbred lines or when the inbreeding coefficient indicates some levels of inbreeding.

The run will create and save the output folders of files and figures.

## 3) To consider the non-additive allelic effects in addition to the additive effects, define **X6** as “initial\_files”, “intermediate\_files”, or “plink\_based\_all” within the

[configuration file](#). The default is “plink\_based\_all”. This allows SAFE-h2-preADOH to do the allele adjustments and association tests together. It uses the “glm” function of the PLINK for the association test.

To do this, put your bfiles (“MAIN.fam, MAIN.bed, and MAIN.bim”) within the SAFE-h2 folder. Covariates can also be provided under the name of Covar\_plink as a tab-delimited text file:

FID	IID	Cov001
Sample1	Sample1	3
Sample2	Sample2	9
...	...	...

Now, the user can open the terminal from the SAFE-h2 folder and run:

```
python3 SAFE-h2-preADOH.py or python3 SAFE-h2-preADOH.py | tee output1
```

The run will create and save the output folder of results for allele adjustments. To perform the SNP heritability profiling based on additive and non-additive effects (on the allele adjusted genotypes), run (again from the SAFE-h2 folder):

```
python3 SAFE-h2-ADOH.py or python3 SAFE-h2-ADOH.py | tee output2
```

(Before the run, **X1-5** within the configuration file can get the values of 0 or 1 to ignore or implement different heritability estimation models. **X4** can also get 2 when dealing with inbred lines or when the inbreeding coefficient indicates some levels of inbreeding)

The user can also consider providing covariates as Covar\_emmax, Covar\_ldak, and qCovar\_greml or Covar\_greml as explained at the beginning for SAFE- $h^2$ .

The run will create and save the output folders of figures and files.

3-1) Instead of the “plink\_based\_all”, user can set the option of “initial\_files” within the configuration file to perform any other association test of interest on the 5 different adjusted genotypic datasets. Here, SAFE-h2-preADOH provides 5 different genotypic datasets (Additive, Dominance [over Ref allele], dominance [over Alt allele], Overdominance [over Ref allele], overdominance [over Alt allele], and Heterosis). For this run:

```
python3 SAFE-h2-preADOH.py or python3 SAFE-h2-preADOH.py | tee output1
```

3-2) The user will need to perform association tests and put the  $p$ -value files within the SAFE-h2 folder. The format of  $p$ -value files should be the same as explained for MainPs. They should be named as “**AdditivePs**” for Additive data, “**DominancePs**” for Dominance adjusted data, “**dominancePs**” for dominance adjusted data, “**OverdominancePs**” for Overdominance adjusted data, “**overdominancePs**” for overdominance adjusted data, and “**HeterPs**” for Heterosis adjusted data.

Furthermore, the user must set the option of “intermediate\_files” instead of “initial\_files” within the configuration file and run:

```
python3 SAFE-h2-preADOH.py or python3 SAFE-h2-preADOH.py | tee output2
```

This will create and save the output folder and combinatory files that can be used for SNP heritability profiling by SAFE-h2-ADOH. Before this step, the user needs to run another association test of interest and put all *p*-value files into the folder. They should be names as MainPs\_Add, MainPs\_AddDom, MainPs\_AddDomOD, MainPs\_AddDomODHet.

For SNP heritability profiling by SAFE-h2-ADOH, the user can also include covariates files. Finally, for performing the SNP heritability profiling based on additive and non-additive effects (on the allele adjusted genotypes), run:

```
python3 SAFE-h2-ADOH.py or python3 SAFE-h2-ADOH.py | tee output3
```

(Before the run, **X1-5** within the configuration file can get the values of 0 or 1 to ignore or implement different heritability estimation models. **X4** can also get 2 when dealing with inbred lines or when the inbreeding coefficient indicates some levels of inbreeding)

The run will create and save the output folders of figures and files.

- 4) To estimate the SNP heritability after minimizing the likelihood for contributions by false-positive SNP hits, **X7** in the configuration file (Prog.config) should be 1. The rest of the settings are already described in section 2. In addition to the MAIN.bim, MAIN.bed, MAIN.fam, and MainPs files for the original phenotype of interest, the user should also provide MAIN1.fam, Main1Ps, MAIN2.fam, Main2Ps, MAIN3.fam, Main3Ps, MAIN4.fam, Main4Ps, MAIN5.fam, Main5Ps, MAIN6.fam, and Main6Ps files related to the 6 independent random phenotypes simulated within the range of phenotype-of-interest. To perform the analysis, run:

```
python3 SAFE-h2.py or python3 SAFE-h2.py | tee output
```

The run will create and save the output folders of files and figures for the original phenotype (named as SAFE\_...) and for the random phenotypes 1 to 6 (named as RP1\_... to RP6\_...). SAFE- $h^2$  also provides the user with aligned graphs and related datasets (named as RP1\_6...).

Output files:

For heritability profiling by additive-only effects (see section 2):

The heritability profiling using Emmax model is saved as outfileE, outfileY\_E, & Heritability\_Bar\_Graph\_E

The heritability profiling using LDAK GCTA\_model is saved as outfileL, outfileY\_L, & Heritability\_Bar\_Graph\_L

The heritability profiling using LDAK Thin\_model is saved as outfileLT, outfileY\_LT, & Heritability\_Bar\_Graph\_LT

The heritability profiling using GCTA-GREML model is saved as outfile\_G, outfileY\_G, & Heritability\_Bar\_Graph\_G

The heritability profiling using Gemma model is saved as outfileGe, outfileY\_Ge, & Heritability\_Bar\_Graph\_Ge

Clustered SNP hits are saved as Number\_of\_pvalues & Clustered\_SNP\_Hits\_Bar\_Graph

### For heritability profiling by additive and non-additive effects (see section 3):

The heritability profiling using Emmax model is saved as outfileE\_α, outfileY\_E\_α, & Heritability\_Bar\_Graph\_E

The heritability profiling using LDAK GCTA\_model is saved as outfileL\_α, outfileY\_L\_α, & Heritability\_Bar\_Graph\_L

The heritability profiling using LDAK Thin\_model is saved as outfileLT\_α, outfileY\_LT\_α, & Heritability\_Bar\_Graph\_LT

The heritability profiling using GCTA-GREML model is saved as outfile\_G\_α, outfileY\_G\_α, & Heritability\_Bar\_Graph\_G

The heritability profiling using Gemma model is saved as outfile\_Ge\_α, outfileY\_Ge\_α, & Heritability\_Bar\_Graph\_Ge

Clustered SNP hits are saved as Number\_of\_pvalues\_α & Clustered\_SNP\_Hits\_Bar\_Graph

α is Add for additive model.

α is AddDom for additive and dominance model.

α is AddDomOD for additive, dominance, and overdominance model.

α is AddDomODHet for additive, dominance, overdominance, and heterosis-like model.

The allele adjusted genotypic datasets also include:

MAIN\_α.bed, MAIN\_α.bim, MAIN\_α.fam, & MainPs\_α

### For heritability profiling after minimizing the false-positive contributions (see section 4):

For each of the original and random phenotypes:

The heritability profiling using Emmax model is saved as outfileE, outfileY\_E, & Heritability\_Bar\_Graph\_E

The heritability profiling using LDAK GCTA\_model is saved as outfileL, outfileY\_L, & Heritability\_Bar\_Graph\_L

The heritability profiling using LDAK Thin\_model is saved as outfileLT, outfileY\_LT, & Heritability\_Bar\_Graph\_LT

The heritability profiling using GCTA-GREML model is saved as outfile\_G, outfileY\_G, & Heritability\_Bar\_Graph\_G

The heritability profiling using GEMMA model is saved as outfile\_Ge, outfileY\_Ge, & Heritability\_Bar\_Graph\_Ge

Clustered SNP hits are saved as Number\_of\_pvalues & Clustered\_SNP\_Hits\_Bar\_Graph

And for all phenotypes together:

The heritability profiling using Emmax model is saved as outfileY\_E, outfileY\_E1, outfileY\_E2, outfileY\_E3, outfileY\_E4, outfileY\_E5, outfileY\_E6, outfileY\_E16, E16, E16\_error, & Heritability\_Bar\_Graph\_E-purifiedEffects

The heritability profiling using LDAK GCTA\_model is saved as outfileY\_L, outfileY\_L1, outfileY\_L2, outfileY\_L3, outfileY\_L4, outfileY\_L5, outfileY\_L6, outfileY\_L16, L16, L16\_error, & Heritability\_Bar\_Graph\_L-purifiedEffects

The heritability profiling using LDAK Thin\_model is saved as outfileY\_LT, outfileY\_LT1, outfileY\_LT2, outfileY\_LT3, outfileY\_LT4, outfileY\_LT5, outfileY\_LT6, outfileY\_LT16, LT16, LT16\_error, & Heritability\_Bar\_Graph\_LT-purifiedEffects

The heritability profiling using GCTA-GREML model is saved as outfileY\_G, outfileY\_G1, outfileY\_G2, outfileY\_G3, outfileY\_G4, outfileY\_G5, outfileY\_G6, outfileY\_G16, G16, G16\_error, & Heritability\_Bar\_Graph\_G-purifiedEffects

The heritability profiling using GEMMA model is saved as outfileY\_Ge, outfileY\_Ge1, outfileY\_Ge2, outfileY\_Ge3, outfileY\_Ge4, outfileY\_Ge5, outfileY\_Ge6, outfileY\_Ge16, Ge16, Ge16\_error, & Heritability\_Bar\_Graph\_Ge-purifiedEffects

## Citations:

Behrooz Darbani, Mogens Nicolaisen. SNP Allocation For Estimating Heritability (SAFE- $h^2$ ): A tool to explore genomic origins of phenotypes for estimation of SNP heritability using additive-only allelic effects or additive and non-additive allelic effects. Doi: <https://doi.org/10.1101/2023.08.28.555092>

Chang CC et al., Second-generation PLINK: rising to the challenge of larger and richer datasets. GigaScience 2015Dec;4(1):s13742-015-0047-8. Doi: <https://doi.org/10.1186/s13742-015-0047-8>.

Depending on the used algorithm users also need to cite one or more of the following studies:

Kang HM et al., Variance component model to account for sample structure in genome-wide association studies. Nat Genet.2010Apr;42(4):348-54. Doi: [10.1038/ng.548](https://doi.org/10.1038/ng.548) ... if applying EMMAX model

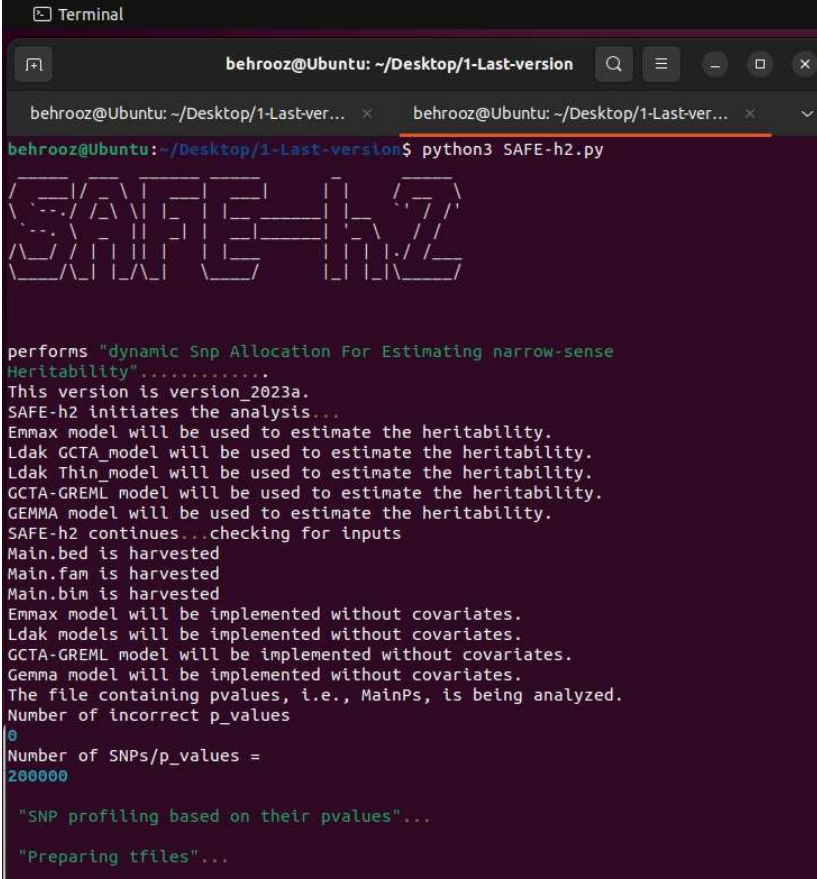
Speed D et al., Improved heritability estimation from genome-wide SNPs. Am J Hum Genet.2012Dec7;91(6):1011-21. Doi: [10.1016/j.ajhg.2012.10.010](https://doi.org/10.1016/j.ajhg.2012.10.010) ... if applying LDAK models

Jang J et al., Common SNPs explain a large proportion of the heritability for human height. Nat Genet.2010Jul;42(7):565-9. Doi: [10.1038/ng.608](https://doi.org/10.1038/ng.608) ... if applying GCTA-GREML model

Zhou, X. & Stephens, M. Genome-wide efficient mixed-model analysis for association studies. Nat Genet.2012;44:821-824. Doi: <https://doi.org/10.1038/ng.2310> ... if applying GEMMA model

## Screenshots from some example runs:

### Without covariate



```
Terminal
behrooz@Ubuntu: ~/Desktop/1-Last-version
behrooz@Ubuntu: ~/Desktop/1-Last-ver... x behrooz@Ubuntu: ~/Desktop/1-Last-ver... x
behrooz@Ubuntu: ~/Desktop/1-Last-version$ python3 SAFE-h2.py

SAFE-h2

performs "dynamic Snp Allocation For Estimating narrow-sense
Heritability".....
This version is version_2023a.
SAFE-h2 initiates the analysis...
Emmax model will be used to estimate the heritability.
Ldak GCTA_model will be used to estimate the heritability.
Ldak Thin_model will be used to estimate the heritability.
GCTA-GREML model will be used to estimate the heritability.
GEMMA model will be used to estimate the heritability.
SAFE-h2 continues...checking for inputs
Main.bed is harvested
Main.fam is harvested
Main.bim is harvested
Emmax model will be implemented without covariates.
Ldak models will be implemented without covariates.
GCTA-GREML model will be implemented without covariates.
Gemma model will be implemented without covariates.
The file containing pvalues, i.e., MainPs, is being analyzed.
Number of incorrect p_values
0
Number of SNPs/p_values =
200000

"SNP profiling based on their pvalues"...

"Preparing tfiles"...
```



## With covariate

```
Terminal
behrooz@Ubuntu: ~/Desktop/1-Last-version
behrooz@Ubuntu: ~/Desktop/1-Last-ver... x behrooz@Ubuntu: ~/Desktop/1-Last-ver... x
behrooz@Ubuntu: ~/Desktop/1-Last-version$ python3 SAFE-h2.py

SAFE-h2

performs "dynamic Snp Allocation For Estimating narrow-sense
Heritability".....
This version is version_2023a.
SAFE-h2 initiates the analysis...
Emmax model will be used to estimate the heritability.
Ldak GCTA_model will be used to estimate the heritability.
Ldak Thin_model will be used to estimate the heritability.
GCTA-GREML model will be used to estimate the heritability.
GEMMA model will be used to estimate the heritability.
SAFE-h2 continues...checking for inputs
Main.bed is harvested
Main.fam is harvested
Main.bim is harvested
The covariate file will be included in the analysis by Emmax model.
The covariate file will be included in the analysis by Ldak models.
The categorical-covariate file will be included in the analysis by GCTA-GREML
model.
The covariate file will be included in the analysis by Gemma model.
The file containing pvalues, i.e., MainPs, is being analyzed.
Number of incorrect p_values
0
Number of SNPs/p_values =
200000
The file containing covariates, i.e., Covar_emmax, is being analyzed.
Number of covariates for EMMAX =
1000
The file containing covariates, i.e., Covar_ldak, is being analyzed.
```

## Additive and non-additive

```
Terminal
behrooz@Ubuntu: ~/Desktop/1-Last-version
behrooz@Ubuntu: ~/Desktop/1-Last-ver... x behrooz@Ubuntu: ~/Desktop/1-Last-ver... x
behrooz@Ubuntu: ~/Desktop/1-Last-version$ python3 SAFE-h2-preADOM.py

SAFE-h2

performs "dynamic Snp Allocation For Estimating narrow-sense
Heritability".....
This version is version_2023a.
SAFE-h2 continues...checking for inputs
Main.bed is harvested
Main.fam is harvested
Main.bim is harvested
"SAFE-h2 builds VCFs and bfiles for 2 Dominance, 2 Overdominance, and 1
heterosis scenarios"...
```



```
Terminal
behrooz@Ubuntu: ~/Desktop/1-Last-version
behrooz@Ubuntu: ~/Desktop/1-Last-ver... x behrooz@Ubuntu: ~/Desktop/1-Last-ver... x
behrooz@Ubuntu: ~/Desktop/1-Last-version$ python3 SAFE-h2-AD0H.py

SAFE-h2

performs "dynamic Snp Allocation For Estimating narrow-sense
Heritability".....
This version is version_2023a.

SAFE-h2 continues...checking for inputs
Main.bed is harvested
Main.fam is harvested
Main.bim is harvested
Main.bed is harvested
Main.fam is harvested
Main.bim is harvested
Main.bed is harvested
Main.fam is harvested
Main.bim is harvested
Main.bed is harvested
Main.fam is harvested
Main.bim is harvested
SAFE-h2 initiates the analysis...
Emmax model will be used to estimate the heritability.
Ldak GCTA_model will be used to estimate the heritability.
Ldak Thin_model will be used to estimate the heritability.
GCTA-GREML model will be used to estimate the heritability.
GEMMA model will be used to estimate the heritability.
Emmax model will be implemented without covariates.
Models by Ldak will be implemented without covariates.
GCTA-GREML model will be implemented without covariates.
Gemma model will be implemented without covariates.
The file containing pvalues, i.e., MainPs_Add, is being analyzed.
```

## End of run

```
Terminal
behrooz@Ubuntu: ~/Desktop/1-Last-version
behrooz@Ubuntu: ~/Desktop/1-Last-ver... x behrooz@Ubuntu: ~/Desktop/1-Last-ver... x
SAFE-h2

"has completed the analyses. The SAFE-h2 version is 2023a."
SAFE-h2 Citation: Behrooz Darbani, Mogens Nicolaisen. SNP Allocation For
Estimating Heritability (SAFE-h2): A tool to explore genomic origins of
phenotypes for estimation of SNP heritability using additive-only allelic
effects or additive and non-additive allelic effects. doi:
https://doi.org/10.1101/2023.08.28.555092
PLINK Citation: GigaScience 2015Dec;4(1):s13742-015-0047-8
(doi:https://doi.org/10.1186/s13742-015-0047-8)
Also cite Nat Genet.2010Apr;42(4):348-54.doi:10.1038/ng.548 ... if using EMMAX
model
Also cite Am J Hum Genet.2012Dec7;91(6):1011-21.doi:10.1016/j.ajhg.2012.10.010
... if using LDAK models
Also cite Nat Genet.2010Jul;42(7):565-9.doi:10.1038/ng.608 ... if using
GCTA-GREML model
Also cite Nat Genet.2012;44:821-824.doi.org/10.1038/ng.2310 ... if using GEMMA
model
behrooz@Ubuntu:~/Desktop/1-Last-version$
```