
Unfavorable Haplotype Finder

Software Tool

Jeremy T. Howard¹
Francesco Tiezzi¹
Christian Maltecca¹

¹North Carolina State University, Raleigh, NC, USA

The logo for North Carolina State University, featuring the words "NC STATE" in white, bold, sans-serif capital letters on a red rectangular background.

NC STATE

Contents

Introduction	1
Overview of Program	1
Disclaimer	3
Computing Environment	4
Running the Program	5
Program Parameters	6
Optimization Parameters	7
Files to Read In	7
Location of Variables	9
Variance Components	10
Running Parameters	11
Output Files	14
Walkthrough Example	15
Data Setup	15
Null Model Variance Components	16
Double Check with AS-REML	17
Running Program	19

Introduction

The Unfavorable Haplotype Finder program is method that is designed to identify haplotypes contained within a run of homozygosity (ROH) that give rise to an unfavorable phenotype. The overall objective is to identify these haplotypes across a variety of traits and the ones that are have a consistent unfavorable effect across multiple traits would be great candidates to use in mating designs in tandem with lethal haplotypes.

Overview of Program

The program is ran in three stages and each one is described below:

Stage 1

- Step 1: Tabulate Means of non-ROH and unique ROH for sliding windows (start with largest desired length).
- Step 2: Combine nested windows

Before (each haplotype contains same set of animals):

Start End Haplotype

132	182	22002220002200222200000222220220202222022000200202
133	183	200222000220022220000022222022022220220220002002020
134	184	00222000220022220000022222022020222202202200020020200
135	185	02220002200222200000222220220202222022022000200202000
136	186	22200022002222000002222202202022220220220002002020000
137	187	220002200222200000222220220202222022202200020020200002
138	188	20002200222200000222220220202222022022000200202000022
139	189	00022002222000002222202202022220220220002002020000220
140	190	00220022220000022222022020222202202200020020200002202

After:

132	190	2200222000220022220000022222022020222202200020020200002202
-----	-----	--

- Step 3: Reduce window size by 5 until lowest desired length is reached.
- Step 4: Combine nested windows .

Before (each haplotype contains same set of animals):

Start End Haplotype

614	656	200020220022002222200000200022020002000020
614	646	20002022002200222220000020002202
614	651	2000202200220022222000002000220200020
614	661	20002022002200222220000020002202000200002002022
614	671	2000202200220022222000002000220200020000200202220002
614	666	2000202200220022222000002000220200020000200202220202
614	641	200020220022002222200000200

After:

614	641	200020220022002222200000200
-----	-----	-----------------------------

Stage 2

- Determine the significance of each window that passed Stage 2 using a model that allows:
 - 1.) Fixed Environmental Effects
 - 2.) Additive effect of animal based on pedigree
 - 3.) Permanent effect of animal
- This step can be done with any program, but is parallelized in this program to make it faster.
- A window is kept if a contrast between a unique ROH versus non-ROH is greater than certain cutoff.

Stage 3

- Remove nested windows

Example: (only keep Window 2)

<u>Window 1</u>	<u>Window 2</u>
Animal	Animal
1	1
2	2
3	3
4	4
5	5
6	6
-	7
-	8
-	9

Disclaimer

This document outlines how to run the Unfavorable Haplotype Finder program and describes the parameters utilized within the program. The software is free for academic and non-commercial use. The authors accept no responsibility for the accuracy of results obtained by using Geno-Diver software.

Please notify jthoward@ncsu.edu or cmaltec@ncsu.edu if you think the results are not correct or you have encountered a bug.

The software was profiled using the Valgrind software ([valgrind](#)). In a few instances, calling functions within the Intel MKL library did result in "possibly lost" errors. Although, as outlined in the Intel MKL [user manual](#), new buffers that the library allocates when an Intel MKL function gets called are not deallocated until the program ends. The use of `mkl_free_buffers()` and `mkl_thread_free_buffers()` were called at the end of the program and this fixed the majority of the issues.

The reference for the software:

-

Computing Environment

The code is written in the C++11 language using object-oriented techniques. The application has been tested to run on Linux platforms. The software makes use of two external libraries, Intel MKL and Eigen.

EIGEN Library:

EIGEN is freely available at: [Eigen Site](#)

Once at the site, download the latest stable release and uncompress it. For example, the current downloaded package is called “eigen-eigen-07105f7124f9.tar”. To use it you just have to place it in the file where all of the other Geno-Diver files are located and uncompress the file. Once you uncompressed the file it will be a folder (e.g. “eigen-eigen-07105f7124f9”). The uncompressed file serves as your path in the makefile outlined below.

Intel MKL Library:

Intel MKL is a commercial library and is available for purchase. However, there is an opportunity to obtain the Intel MKL library (for Linux) free of charge for non-commercial use at the following website: [Intel MKL Site](#)

The Intel MKL can sometimes be tricky to download and link, but there is a step-by-step protocol within the folders that are downloaded or instructions are at [Intel MKL Guide](#). Depending on the computing system you are running, a guide to linking the Intel MKL libraries are at [Intel MKL Linking Guide](#).

C++11 Version:

Versions of gcc 4.7 or newer support the C++11 standards. You need to install or update to the correct version of gcc using the standard package manager or installer, depending on what type of OS you are using. Some helpful websites include:

[gcc helper 1](#) [gcc helper 2](#)

Compiling:

Once both EIGEN, Intel MKL libraries and gcc version 4.7 or newer have been correctly installed and the folders placed in the directory where all of the Geno-Diver source code files are located the last thing you have to do is change the path for EIGEN and Intel MKL libraries.

In the makefile change lines 13 and 14 to the path which aligns to the updated EIGEN and MKL path.

After changing the path, type “make” on the command line. An executable file called “ROH_Finder” is now in your working directory.

Running the Program

At the current time executable files are available only for linux operating environments. To run the program place the following file in a folder:

- ROH_Finder.

Prior to running the program you will need to check the permissions of the files and if so make them executable (i.e. type “chmod 755 ROH_Finder”). Once the permissions have been changed you need to generate a parameter file and place it in the same folder as ROH_Finder. If you are new to the program an example is outlined. The parameter file is read by searching for key words that are capitalized and then followed by a colon. Therefore any other phrase will not be utilized.

To run the program type in “./ROH_Finder” and then the name of your parameter file.

Program Parameters

A parameter file that specifies all the parameters is outlined below. Not all of them are required for the program to run, which was done to reduce the complexity of running the software. The key words with a “!!” are not mandatory for the program to run. All key words are in capital letters and the parameter(s) specified are separated by spaces. Only parameters after the key words impact the simulation.

-----	Parameter File	-----
THREAD: 4		
-----	Files to Read In	-----
MAP_FILE: haplo_map_.txt		
PHENO_FILE: haplo_phenotype.txt		
GENO_FILE: haplo_genotype.txt		
PEDIGREE: Pedigree_File.txt.SRT		
-----	Location of Variables	-----
ID: 1		
CLASS: 2,3 !!		
COV: 4 !!		
PHENOTYPE: 5		
-----	Variance Components	-----
VAR_RESIDUAL: 0.05		
VAR_ANIMAL: 0.025		
VAR_PERMANENT: 0.0 !!		
-----	Running Parameters	-----
WINDOW: 50,45,40,35,30,25,20,15 !!		
CUTOFF: data 1000		
UNFAV_DIRECTION: low		
ONE_SIDED_T_CUTOFF: 2.326 !!		
SUBTRACT_MEAN: no !!		
MINIMUM_FREQ: 0.0075 !!		
ASREML_CHECK: no !!		

A description of all the parameters that can be specified in the program is outlined below along with helpful hints and suggestions. A complete example going from start to finish is outlined in Example 1.

Optimization Parameters

THREAD

Description: - Declares the number of threads used for parallel processing.

Value: - Integer value based on number of cores available.

Usage: - "THREAD: 4".

Type: - Mandatory.

Note:

- Running the program on a single thread will take a large amount of time. The program is designed to run multiple models at a time and therefore the running time is greatly reduced.

Files to Read In

MAP_FILE

Description:

- The file where SNP location is stored and ordered by chromosome then position. Column 1 is chromosome and column 2 is SNP position with a space delimiter.

Value: - String.

Usage: - "MAP_FILE: haplo_map_.txt".

Type: - Mandatory.

Note:

- Example files are provided as a template and is described in Example 1.

PHENO_FILE

Description:

- The file where ID, phenotype and the variables that are fixed in the model are located. Separator is space delimiter. No animals can have missing datapoints.

Value: - String.

Usage: - "PHENO_FILE: haplo_pheno.txt".

Type: - Mandatory.

Note:

- Example files are provided as a template and described in Example 1. A log file, referred to as "LogFile" is created that prints out how each line is interpreted.

GENO_FILE:Description:

- The file where ID and genotype are located. Separator is space delimiter. The genotype is represented as a string of 0, 1, or 2 (i.e. 021012). Can't have any missing genotypes. It is assumed that the first row of map file is the first genotype, second row is the second genotype etc.....

Value:

- String.

Usage:

- "GENO_FILE: haplo_geno.txt".

Type:

- Mandatory.

Note:

- Example files are provided as a template and described in Example 1.
-

PEDIGREEDescription:

- The pedigree file where ID, Sire and Dam is outlined and the delimiter is a space. The ID's can be any format (i.e. alpha-numeric or numeric), but they have to be sorted so that parents come before progeny.

Value:

- String.

Usage:

- "PEDIGREE: Pedigree_File.txt.SRT".

Type:

- Mandatory.

Note:

- Example files are provided as a template and described in Example 1. If one is generating variance components with AS-REML the sorted pedigree file using the !SORT command can be utilized.

Location of Variables

ID

Description: - Location in phenotype file that has the animal ID variable.

Value: - Integer.

Usage: - "ID: 1".

Type: - Mandatory.

CLASS

Description: - Location in phenotype file that has the class type fixed effects.

Value: - Integer.

Usage: - "CLASS: 2,3".

Type: - Optional.

Note:

- For each class variable the first level is zeroed out to make the coefficient matrix invertible. At the current time the program doesn't check to see if more dependencies exist. Therefore small CG should be removed or aggregated. It is advisable that when you get the variance component estimates you check to see how many levels got zeroed out for each class variable. If more than one effect is zeroed out for each variable this program WILL NOT GENERATE THE CORRECT RESULTS!!! When generating least-square means (LSM) for each unique ROH and non-ROH the class effects are averaged over. No animals can have missing observations.

COV

Description: - Location in phenotype file that has the covariate type fixed effects..

Value: - Integer.

Usage: - "COV: 4".

Type: - Optional.

Note:

- When generating least-square means (LSM) for each unique ROH and non-ROH the average covariate value is utilized in the LSM estimate

PHENOTYPE

Description: - Location in phenotype file that has the phenotype.

Value: - Integer.

Usage: - “PHENOTYPE: 5”.

Variance Components

VAR_RESIDUAL

Description: - Residual variance of null model without haplotype effect.

Value: - Double.

Usage: - “VAR_RESIDUAL: 0.05”.

Type: - Mandatory.

VAR_ANIMAL

Description: - Additive genetic variance of null model without haplotype effect.

Value: - Double.

Usage: - “VAR_ANIMAL: 0.025”.

Type: - Mandatory.

VAR_PERMANENT

Description: - Permanent animal variance of null model without haplotype effect. If absent assumed not included.

Value: - Double.

Usage: - “VAR_PERMANENT: 0.01”.

Type: - Optional.

Running Parameters

WINDOW

Description:

- Window sizes to scan across the genome and starts at largest number.

Value:

- Integers separated by a comma.

Usage:

- "WINDOW: 50,45,40,35,30,25,20,15".

Type:

- Optional. Default is 50,45,40,35,30,25,20,15.

Note:

- I have played around with adjusting these such as making the step size smaller as the window size decreases and reducing the smallest window and the results stay relatively consistent. This has only been utilized using medium density genotypes (i.e. 60k). My expectation that if you use higher density genotypes it should do a better job at finding the optimal haplotype that best represents the ROH haplotype.

CUTOFF

Description:

- Used to define the mean phenotype of the unfavorable haplotype that you would like to investigate with the full model. You can either specify a specific value or perform multiple (i.e. 1,000) full models with randomly chosen windows across the genome to determine the mean phenotype value that results in a t-value of at least 1.96 and any mean phenotype above or below this depending on the unfavorable direction will be investigated with the full model.

Value:

- Method:

- value: the users specifies a value that is deemed unfavorable.

- data: the software determines the threshold by performing the multiple models across a random region of the genome.

- Value:

- If value was chosen: The users specifies the minimum/maximum mean phenotypic value

- If data was chosen: The number of full models that the initial scan runs to chose a minimum/maximum mean phenotypic value.

Usage: - "CUTOFF: data 1000".
Type: - Mandatory

UNFAV

Description: - Direction of the phenotype that is unfavorable. (Example: For the trait Number of Piglets Born Alive the direction is low).
Value: - low or high
Usage: - "DIRECTION: low".
Type: - Mandatory.

ONE_SIDED_T_CUTOFF

Description: - One sided T-Statistic after running the full model that is regarded as being significant.
Value: - Double
Usage: - "ONE_SIDED_T_CUTOFF: 2.326".
Type: - Optional. Default is 2.326.

SUBTRACT_MEAN

Description: - Useful for simulation purposes with no fixed effects.
Value: - yes or no
Usage: - "SUBTRACT_MEAN: no".
Type: - Optional. Default is no.

MINIMUM_FREQ

Description: - Minimum number of times a unique ROH is required in order to enter into full model.
Value: - ranging from 0.0 to 1.0
Usage: - "MINIMUM_FREQ: 0.0075".
Type: - Optional. Default is 0.0075.

ASREML_CHECK

Description:

	- This produces the output that is generated by the software and a phenotype file that can be used as input in ASREML. If worked properly the results will be very similar.
<u>Value:</u>	- yes or no
<u>Usage:</u>	- "ASREML_CHECK: no".
<u>Type:</u>	- Optional. Default is no.

Output Files

The program output a file called Stage2_Regions that contains all the information on regions that were declared as significant. Below is a description on what the column heading are describing.

Chromosome: Chromosome haplotype is located.

StartPos: Nucleotide start position of haplotype.

EndPos: Nucleotide end position of haplotype.

StartIndex: Start index in full genotype string of haplotype.

EndIndex: End index in full genotype string of haplotype.

Genotype: Genotype string of unfavorable haplotype.

PhenoMean: Mean phenotype of unfavorable haplotype.

BetaEffect: Beta estimate of unfavorable haplotype.

LSM: Least square mean of unfavorable haplotype.

T-Stat: T-statistic of ROH haplotype versus non-ROH haplotype.

Walkthrough Example

Data Setup

In order to make the process easier I have created a simulated population with all of the datasets already pre-made. The pre-made datasets are:

- haplo_map.txt
 - A file that contains chromosome and position.
- haplo_geno.txt
 - A file that contains ID then genotype as a string.
 - Each row of map file is matched up with genotype string.
 - If multiple traits are being run on a similar set of individuals it is advisable to make a single genotype file. Animals that are not used in the program can still be in the genotype file, but if an animal isn't in the genotype file the program will exit.
- haplo_pheno.txt
 - A file that contains ID then phenotype.
 - It is advisable to make a phenotype file for each unique trait because their can't be any missing datapoints.
 - In the current situation low values are unfavorable.
- pedigreefile.txt
 - A file that contains ID, sire, dam information. It is already ordered so parents come before progeny.
 - If the pedigree needs to be ordered there are multiple options available including the !SORT option in AS-Reml.

When class fixed effects are included in the model the software still is not perfect when issues with small contemporary groups exist and therefore the issues related to confounding between variables may produce wrong results. Therefore as a check it is advisable to ensure that only a single level of each class fixed effect is zeroed out when estimating variance components for null model.

Null Model Variance Components

Whichever program you are comfortable with you need to estimate variance components for the dataset to use in the full models in order to derive significance for a unique ROH genotype. Below is an example of an input file for AS-Reml (“asreml_null.as” in folder).

```
!Workspase 16384
File
Animal * !P
Pheno
pedigreefile.txt !ALPHA !SORT
haplo_pheno.txt !MAXIT 100 !EXTRA 2 !FCON !DDF 2
Pheno mu !r Animal
0 0 1
Animal 1
0 0 AINV 1 !GP
```

The residual variance was estimated to be 0.722969 and the additive genetic component was estimated to be 0.204449.

These are the values that will get utilized in the full model.

Double Check with AS-REML

Prior to running the complete analysis and if AS-Reml is available it is preferred to double check with what AS-Reml would get. Below is an example parameter file that outputs a dataframe to use as input in ASREML and estimates to compare with AS-REML (“Parameter_File” in folder).

-----	Parameter File	-----
THREAD: 4		
-----	Files to Read In	-----
MAP_FILE: haplo_map_.txt		
PHENO_FILE: haplo_pheno.txt		
GENO_FILE: haplo_geno.txt		
PEDIGREE: pedigreefe.txt.SRT		
-----	Location of Variables	-----
ID: 1		
PHENOTYPE: 2		
-----	Variance Components	-----
VAR_RESIDUAL: 0.722969		
VAR_ANIMAL: 0.204449		
-----	Running Parameters	-----
CUTOFF: data 1000		
UNFAV_DIRECTION: low		
ASREML_CHECK: yes		

At the bottom of the LogFile it outlines the estimates that are derived from the software and a file called TestDFa.txt is outputted that has appended the unique ROH genotype string for each animal. The results are outlined below:

ROH ID	Beta Estimate	Least Square Mean	T-stat (non-ROH versus ROH
1	-0.308	1.078	-1.242
2	-0.183	1.203	-0.826
3	-0.924	0.461	-3.602

These results can then be compared with the true model where variance components are re-estimated (“asreml.check.as” in folder). More complex models with class fixed effects can be averaged over in ASReml by adding “!PRESENT” qualifier and then the associated class variables after it.

```
!Workspase 16384
File
Animal * !P
Haplotype * !A !SORT Pheno
pedigreefile.txt !ALPHA !SORT
TestDFa.txt !MAXIT 100 !EXTRA 2 !FCON !DDF 2
Pheno mu Haplotype !r Animal
0 0 1
Animal 1
0 0 AINV 1 !GP
predict Haplotype !TDIFF
```

The results using AS-Reml are outlined below and are similar:

ROH ID	Beta Estimate	Least Square Mean	T-stat (non-ROH versus ROH
1	-0.307	1.083	-1.24
2	-0.182	1.208	-0.83
3	-0.926	0.464	-3.62

Running Program

Once the program has been tested on the dataset the full analysis can be run. Below is the parameter file. In the log file on around line 174 the phenotype cutoff from the previous run can be utilized (may be slightly different on your run).

-----	Parameter File	-----
THREAD: 4		
-----	Files to Read In	-----
MAP_FILE: haplo_map_.txt		
PHENO_FILE: haplo_pheno.txt		
GENO_FILE: haplo_geno.txt		
PEDIGREE: pedigreefe.txt.SRT		
-----	Location of Variables	-----
ID: 1		
PHENOTYPE: 2		
-----	Variance Components	-----
VAR_RESIDUAL: 0.722969		
VAR_ANIMAL: 0.204449		
-----	Running Parameters	-----
CUTOFF: value 1.42227		
UNFAV_DIRECTION: low		

Once the program has finished it will output the regions that were declared as significant into the Stage2_Regions File. The columns have heading to aid in understanding.