

Using GWAF package to conduct association analysis with family data

Ming-Huei Chen and Qiong Yang
Departments of Neurology and Biostatistics
Boston University
Contact: qyang@bu.edu

Table of Contents

Overview	1
Methods	1
Required Files	2
Examples	4
Output	5

Overview

This package GWAF (Genome-Wide Association analyses with Family data) was designed mainly to analyze a batch of SNPs against a continuous or dichotomous phenotype measured on subjects of families for association. The number of SNPs that can be analyzed at once depends on the memory capacity of your system. For genome-wide association studies (GWAS), if the memory is not enough to analyze all SNPs together, one can split the dataset by columns into several datasets, and analyze each of them sequentially using functions in this package. In addition, GWAF also provides functions for making genome-wide p-values plot and QQ plot.

Methods

Linear mixed effects model (LME) is used in this package to analyze continuous traits, with person specific random effects correlated according to degree of relatedness (i.e. kinship coefficient) within a family to account for within family correlation. Logistic regressions via generalized estimating equations (GEE) is used

in this package to analyze dichotomous traits, treating each pedigree (i.e. individuals with the same family id) as a cluster with independent working correlation structure used in the robust variance estimator. These methods were implemented in the `lmekin()` and `gee()` functions in packages `kinship` and `gee`, and our package is a wrapper that enables users to analyze more than one SNP and automatically summarizing the results in an informative and convenient output.

Required Files

Before performing analyses with this package, following files have to be created.

1. **Pedigree file:** A file containing all the families is required. It should be comma delimited with the column names exactly the same (case sensitive) as in following example. Missing father (fa) or mother (mo) ids should be 0. Individuals who are unrelated to anyone can be included as family of size 1.

```
famid,id,fa,mo,sex
1,10,0,0,1
1,11,0,0,2
1,12,10,11,1
1,13,10,11,1
3,32,0,0,1
3,33,0,0,2
3,334,32,33,1
3,335,32,33,2
10,50,0,0,1
11,60,0,0,2
```

2. **Kinship coefficient matrix for LME:** For LME, one should create the kinship coefficient matrix as an R object and save it on disk, as shown below. Then in the LME analysis, path to this R object file is supplied to the 'kinmat' argument in the `lme` analysis function (`lme.batch`).

- Sample R code for create kinship coefficient matrix (must be named kmat in R as shown below)

```
library(kinship)
kmat<-makekinship(ped$famid,ped$id,ped$fa,ped$mo)
## using twice the kinship coefficient
kmat<-kmat*2
## save the kinship matrix to a file
save(kmat,file="fhs_unrel_comb.kinship.Rdata")
```

- In LME analyses, supply the path to the kinship coefficient matrix file to “kinmat” argument.

```
lme.batch(phenfile, genfile, pedfile, phen,
kinmat="fhs_unrel_comb.kinship.Rdata", ...)
```

3. Format of phenotype and genotype files

Phenotype file contains phenotype and covariates. It should be comma delimited with first row being id, followed by variable names. Use empty space for missing values. The first column must be personal id with “id” as column name.

Dichotomous phenotype must be coded as 0, 1 with 1 being affected.

Covariates values must be coded numerically (dichotomous covariate can have any two numeric values). Following is an example of the phenotype file:

```
id,phen1,phen2,covar1,covar2
10,100,1,1,0.2
112,,0,1,0.3
312,130,1,2,0.4
513,125,0,,0.5
```

Genotype file should also be comma delimited with first row being id and SNP names. The first column must be id. **Genotype should be coded as 0, 1, 2 representing the copies of the coded allele.** Use empty space for missing genotypes. SNP names should not contain special characters such as “-“,”/”, etc. But “.” and “_” are allowed. For example:

```
id,SNP.1,SNP_2
10,0,1
11,,
12,1,2
13,2,0
```

Examples

Install R packages: Please install **R**, and packages: GWAF, kinship, gee that can be downloaded from R website (<http://cran.r-project.org/>).

Here are example function calls for analyzing a single phenotype against all SNPs in genotype file. Suppose "phenfile.csv", "genfile.csv", "pedfile.csv", are directory paths to phenotype, genotype and pedigree files respectively; "phen1" is the name of the phenotype to be analyzed. Kinship coefficient matrix file is "fhs_unrel_comb.kinship.Rdata".

LME:

```
library(GWAF)
lme.batch(phenfile="phenfile.csv", genfile="genfile.csv",
pedfile="pedfile.csv", phen="phen1",
model="a",kinmat="fhs_unrel_comb.kinship.Rdata",covars=c(
"covar1","covar2"),outfile="lme.result.csv")
#covars argument can be omitted if no covariates need to be adjsuted
```

GEE:

```
library(GWAF)
gee.lgst.batch(phenfile="phenfile.csv",
genfile="genfile.csv", pedfile="pedfile.csv",
phen="phen2",
model="a",covars=c("covar1","covar2"),outfile="gee.result
.csv")
```

#covars argument can be omitted if no covariates need to be adjusted

Important: These functions are designed to analyze a single phenotype against all the SNP genotypes in a genotype file in a single call. To analyze multiple phenotypes, multiple calls of the functions are needed.

Output

Output information: Output from a function call is saved to the file specified in *outfile* argument in each function. Tables 1 and 2 describe the output columns for LME and GEE analyses, respectively.

Table 1: Output columns from LME analysis (Genotype should be coded as 0, 1, 2 representing the copies of the coded allele)

Column	Description
phen	Phenotype Name
snp	SNP name
n0	number of subjects with non-missing phenotype and genotype 0
n1	number of subjects with non-missing phenotype and genotype 1
n2	number of subjects with non-missing phenotype and genotype 2
h2q^s	% total phenotypic variance explained by the SNP
<i>Output fields for additive, dominant or recessive model</i>	
beta	<u>additive model</u> : beta coefficient per 1 copy increment of coded allele; <u>recessive model</u> : beta coefficient for genotype 2 vs. all other genotypes; <u>dominant model</u> : beta coefficient for 1 and 2 combined vs. genotype 0;
chisq	Chi-square statistic for testing beta equal to zero
df	degrees of freedom for the chi-square statistic
model	model used in the analysis
pval	p-value of the chi-square statistic
<i>Output fields for general model</i>	
beta10	beta coefficient for genotype 1 vs. 0. If the dominant model is used in the analysis, this is the beta coefficient for genotype 1 and 2 combined vs. genotype 0.
beta20	beta coefficient for genotype 2 vs. 0
beta21	beta coefficient for genotypes 2 vs. 1
se10	standard error of beta10

se20	standard error of beta20
se21	standard error of beta21
chisq	Chi-square statistic for testing global hypothesis that both beta10 and beta20 equal zero
df	degrees of freedom of the chi-square statistic
model	model used in the analysis
pval	p-value of the chi-square statistic

$$^s h_q^2 = \max\left(0, \frac{\sigma_{G.null}^2 + \sigma_{e.null}^2 - \sigma_{G.full}^2 - \sigma_{e.full}^2}{Var(y)}\right), \text{ where } Var(y) \text{ is the total phenotypic}$$

variance, $\sigma_{G.null}^2, \sigma_{e.null}^2$ are the polygenic variance and error variance when modeling without the tested SNP, and $\sigma_{G.full}^2, \sigma_{e.full}^2$ are the polygenic variance and error variance when modeling with the SNP.

Table 2. Output columns from GEE analyses (Genotype should be coded as 0, 1, 2 representing the copies of the coded allele)

Column	Description
phen	Phenotype Name
snp	SNP name
n0	number of subjects with non-missing phenotype and genotype 0
n1	number of subjects with non-missing phenotype and genotype 1
n2	number of subjects with non-missing phenotype and genotype 2
nd0	number of diseased subjects with genotype 0
nd1	number of diseased subjects with genotype 1
nd2	number of diseased subjects with genotype 2
miss.0	rate of missing genotypes among non-diseased subjects
miss.1	rate of missing genotypes among diseased subjects
miss.diff.p	P-value of test of differential missingness between unaffected and affected subjects
<i>Output fields when additive, dominant or recessive model specified in control file</i>	
beta	<u>additive model</u> : beta coefficient per 1 copy increment of coded allele; <u>recessive model</u> : beta coefficient for genotype 2 vs. genotypes 0 and 1 combined <u>dominant model</u> : beta coefficient for genotype 1 and 2 combined vs. genotype 0
se	standard error of beta
chisq	Chi-square statistic for testing beta equal to zero
df	degrees of freedom of the chi-square statistic
model	model used in the analysis
remark	warning or additional information for the analysis
pval	p-value of the chi-square statistic
<i>Output fields for general model</i>	

beta10	beta coefficient for genotype 1 vs. 0. If the dominant model is used in the analysis, this is the beta coefficient for genotype 1 and 2 combined vs. genotype 0.
beta20	beta coefficient of genotype with 2 copies of coded allele vs. that with 0 copy
beta21	beta coefficient of genotype with 2 copies of coded allele vs. that with 1 copy
se10	standard error of beta10
se20	standard error of beta20
se21	standard error of beta21
chisq	Chi-square statistic for testing at least one of the beta10 and beta20 not zero
df	degrees of freedom of the chi-square statistic
model	model used in the analysis
remark†	warning or additional information for the analysis
pval	p-value of the chi-square statistic

† Remark column contains warning or additional information. Here is a detailed explanation of the meaning of each remark.

Remark	Reason
"not converged"	The GEE analysis did not converge. So results are not reliable and should be discarded.
"logistic reg"	Logistic regression assuming independent observations is performed, when the number of pedigrees with 2 or more individuals is less than 10 or there are zero genotype counts in any cell of snp by phenotype (3 by 2) table.
"exp count<5"	At least one expected count is less than 5 in 2xN table, N =number of genotype categories for general model, and N=2 for other models. The test results may have a higher false positive rate.
"not converged & exp count<5"	See above
"logistic reg& exp count<5"	See above
"collinearity"	If there are any covariates highly correlated with a snp (abs(correlation)>0.9999999),no analysis is performed.