

Genome-wide Association Analysis - Data Quality Control

Copyright © 2024 Merry-Lynn McDonald, Isabelle Schrauwen & Suzanne M. Leal

Introduction

In this exercise, you will learn how to perform data quality control (QC) by removing markers and samples that fail QC quality control criteria. You will also examine your samples for individuals that are related to each other and/or are duplicate samples. Each sample will also be tested for excess homozygosity and heterozygosity of genotype data. Each SNP will be tested for deviations from Hardy-Weinberg Equilibrium. These exercises will be carried out using PLINK1.9 and R.

1. Using PLINK

PLINK can upload data in different formats please see the PLINK documentation (<https://www.cog-genomics.org/plink/1.9/input>) for additional details. The data for this exercise is in PLINK/LINKAGE file format. There are two files: a pedfile (GWAS.ped) and a map file (GWAS.map). Please examine these files and the PLINK documentation. Please note the commands must be given in the directory where the data resides.

Navigate via the command prompt to the directory which contains the files for the exercise. Type **plink** in the command prompt and make note of the output. Next type:

```
plink --file GWAS
```

Note, that PLINK outputs a file called **plink.log** that contains the same output which you see on the screen. To see all options, type `plink --help` for more information. Determine how many samples there are in your data set and fill in Oval 1 of the flowchart below.

2. Data Quality Control

a. Removing Samples and SNPs with Missing Genotypes.

You will exclude samples that are missing more than 10% of their genotype calls. These samples are likely to have been generated using low quality DNA and can also have higher than average genotyping error rates.

```
plink --file GWAS --mind 0.10 --recode --out GWAS_clean_mind
```

Examine **GWAS_clean_mind.log** to see how many samples are excluded based on this criterion and fill in Box 1.

Create two versions of your dataset, one with SNPs with a minor allele frequencies (MAFs) $\geq 5\%$ and the other with SNPs with a MAFs $< 5\%$.

You will now remove SNPs with MAFs $\geq 5\%$ that are missing $> 5\%$ of their genotypes and then remove SNPs with MAFs $< 5\%$ that are missing $> 1\%$ of their genotypes. SNPs which are missing genotypes can have higher error rates than those SNP markers without missing data.

```
plink --file GWAS_clean_mind --maf 0.05 --recode --out MAF_greater_5
```

```
plink --file GWAS_clean_mind --exclude MAF_greater_5.map --recode --out MAF_less_5
```

```
plink --file MAF_greater_5 --geno 0.05 --recode --out MAF_greater_5_clean
```

Fill in Box 2a.

```
plink --file MAF_less_5 --geno 0.01 --recode --out MAF_less_5_clean
```

Fill in Box 2b.

Merge the two files.

```
plink --file MAF_greater_5_clean --merge MAF_less_5_clean.ped MAF_less_5_clean.map --
recode --out GWAS_MAF_clean
```

A more stringent criterion for missing data is used, samples missing >3% of their genotypes are removed.

```
plink --file GWAS_MAF_clean --mind 0.03 --recode --out GWAS_clean2
```

Fill in Box 3.

b. Checking Sex

Error of the reported sex of an individual can occur. Information from the SNP genotypes can be used to verify the sex of individuals, by examining homozygosity (F) on the X chromosome for every individual. F is expected to be <0.2 in females and >0.8 in males. To check sex run

```
plink --file GWAS_clean2 --check-sex --out GWAS_sex_checking
```

Use R to examine the GWAS_sex_checking.sexcheck file and determine if there are individuals whose recorded sex is inconsistent with genetic sex.

```
R
sexcheck = read.table("GWAS_sex_checking.sexcheck", header=T)
names(sexcheck)
sex_problem = sexcheck[which(sexcheck$STATUS=="PROBLEM"),]
sex_problem
q()
```

NA20530 and NA20506 were coded as a female (2) and from the genotypes appear to be males (1). In addition, 3 individuals (NA20766, NA20771 and NA20757) do not have enough information to determine if they are males or females and PLINK reports sex = 0 for the genotyped sex. Fill in the table below:

Table 1: Sex check

FID	IID	PEDSEX	SNPSEX	STATUS	F
NA20506	NA20506				
NA20530	NA20530				
NA20766	NA20766				
NA20771	NA20771				
NA20757	NA20757				

Reasons for these kinds of discrepancies, include the records are incorrect, incorrect data entry, sample swap, unreported Turner or Klinefelter syndromes. Additionally, if a sufficient number of SNPs have not been genotyped on the X chromosome it can be difficult to accurately predict the sex of an individual. In this dataset, there are only 194 X chromosomal SNPs. If you cannot validate the sex of the individual they should be removed. For this exercise, we are going to assume that when the sex was checked, we found it was incorrectly recorded (i.e. these samples were male). Therefore, this error could simply be corrected.

Question 1: Why do you expect the homozygosity rate to be higher on the X chromosome in males than females?

c. Duplicate Samples

The following PLINK command can be used to check for duplicate samples:

```
plink --file GWAS_clean2 --genome --out duplicates
```

Open the **duplicates.genome** file in R with the following command:

```
dups = read.table("duplicates.genome", header = T)
```

We are interested in the Pi-Hat (the estimated proportion IBD sharing) value. You may notice that there is more than one duplicate (Pi-Hat= ~ 1). Also, examine the output for pairs of individuals with high Pi-Hat values which can indicate they are related. The amount of allele sharing [Z(0), Z(1) and Z(2)] across all SNPs provides information on the type of relative pair.

```
problem_pairs = dups[which(dups$PI_HAT > 0.4),]  
problem_pairs
```

Table 2: Duplicate and Related Individuals

FID1	IID1	FID2	IID2	Z(0)	Z(1)	Z(2)	PI_HAT
FID1- Family ID for 1st individual; IID1 - Individual ID for 1st individual; FID2- Family ID for 2nd individual; IID2 - Individual ID for 2nd individual; Z(0)- P(IBD=0); Z(1)- P(IBD=1); Z(2)- P(IBD=2); PI_HAT-P(IBD=2)+0.5*P(IBD=1) (proportion IBD)							

Question 2: How many duplicate pairs do you find (**hint: Pi-Hat = ~ 1**)? Do pairs with a **Pi-Hat = ~ 1** have to be duplicate samples? What is another explanation? What proportion would you expect a parent/ child to share IBD? Can you find any such relationship?

Note: Pi-hat can be inflated and individuals appear to be related to each other if you have samples from different populations. This explains why we observe pairs of individuals with Pi-hat > 0.05 since three distinct populations were analyzed. Additionally, this phenomenon can be observed if a subset(s) of samples have higher genotyping/sequencing error rates, which creates two or more “populations” and the individuals within these “populations” incorrectly appear to be related.

Using this R script please observe how many sample pairs have pi-hat > 0.05 :

```
problem_pairs = dups[which(dups$PI_HAT > 0.05),]  
myvars = c("FID1", "IID1", "FID2", "IID2", "PI_HAT")  
problem_pairs[myvars]
```

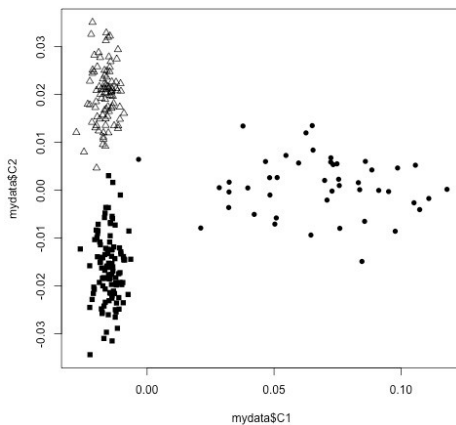
Create the following txt file:

```
1344 NA12057  
1444 NA12739  
M033 NA19774
```

name it ‘IBS_excluded.txt’ and save it to the folder with your PLINK data. Give the command:

```
plink --file GWAS_clean2 --remove IBS_excluded.txt --recode --out GWAS_clean3
```

Fill in Box 4 and Oval 3.



As part of QC usually the data is examined for outliers by plotting the first and second principal or multidimensional scaling (MDS) components. Using a subset of markers that have been trimmed to remove LD ($r^2 < 0.5$). Principal components analysis (PCA) and MDS will be performed in the second part of the exercise to detect outliers and control for populations substructure. Outlier can be due to study subjects coming from different populations e.g. European- and African-Americans or batch effects. If it is suspected that outliers are due to study subjects having been sampled from different populations than data from HapMap can be included to elucidate population membership, e.g. for a study of European-Americans if African-American study subjects are included they would cluster between the European and African HapMap samples. If you perform this type of analysis you should remove the HapMap samples and re-estimate the MDS or PC components before adjusting for population substructure

or stratification. For this exercise data **is used** from HapMap Phase III which consists of CEU (Europeans from Utah), MEX (Mexicans from Los Angeles) and TSI (Tuscans from Italy). Three clusters can be observed that consist of the three data sets but no extreme outliers are observed. This data set is being used for demonstration purposes. Different populations should be analyzed separately and the results can be combined using meta-analysis. In part two of this exercise MDS and PC components will be constructed and analyzed.

d. Hardy-Weinberg Equilibrium (HWE):

To test for HWE we will test separately in each ancestry group and by case-control status. Therefore, we will need to use information on ancestry and cases-control status. Please note that this should be tested in the 3 different populations separately (CEU, MEX, TSI), but due to the small sample sizes, we tested it in the 3 populations together for example purposes. It should also be noted if the sample sizes are small it is difficult to detect a deviation from HWE.

```
plink --file GWAS_clean3 --pheno pheno.txt --pheno-name Aff --hardy
```

Using R examine the file **plink.hwe** and look for SNPs with p-values of 10^{-7} or smaller.

```
hardy = read.table("plink.hwe", header = T)
names(hardy)
hwe_prob = hardy[which(hardy$P < 0.0000009),]
hwe_prob
```

Using a criterion of $p < 10^{-7}$ to reject the null hypothesis of HWE, how many SNPs fail HWE in the controls? Fill out Oval 5 and Box 4. Using the same criteria, how many SNPs fail HWE in the controls? Complete Table 2 with this information.

Table 3: Hardy-Weinberg Equilibrium

Cases			Controls		
SNP	Pvalue	Population(s)	SNP	Population(s)	Pvalue

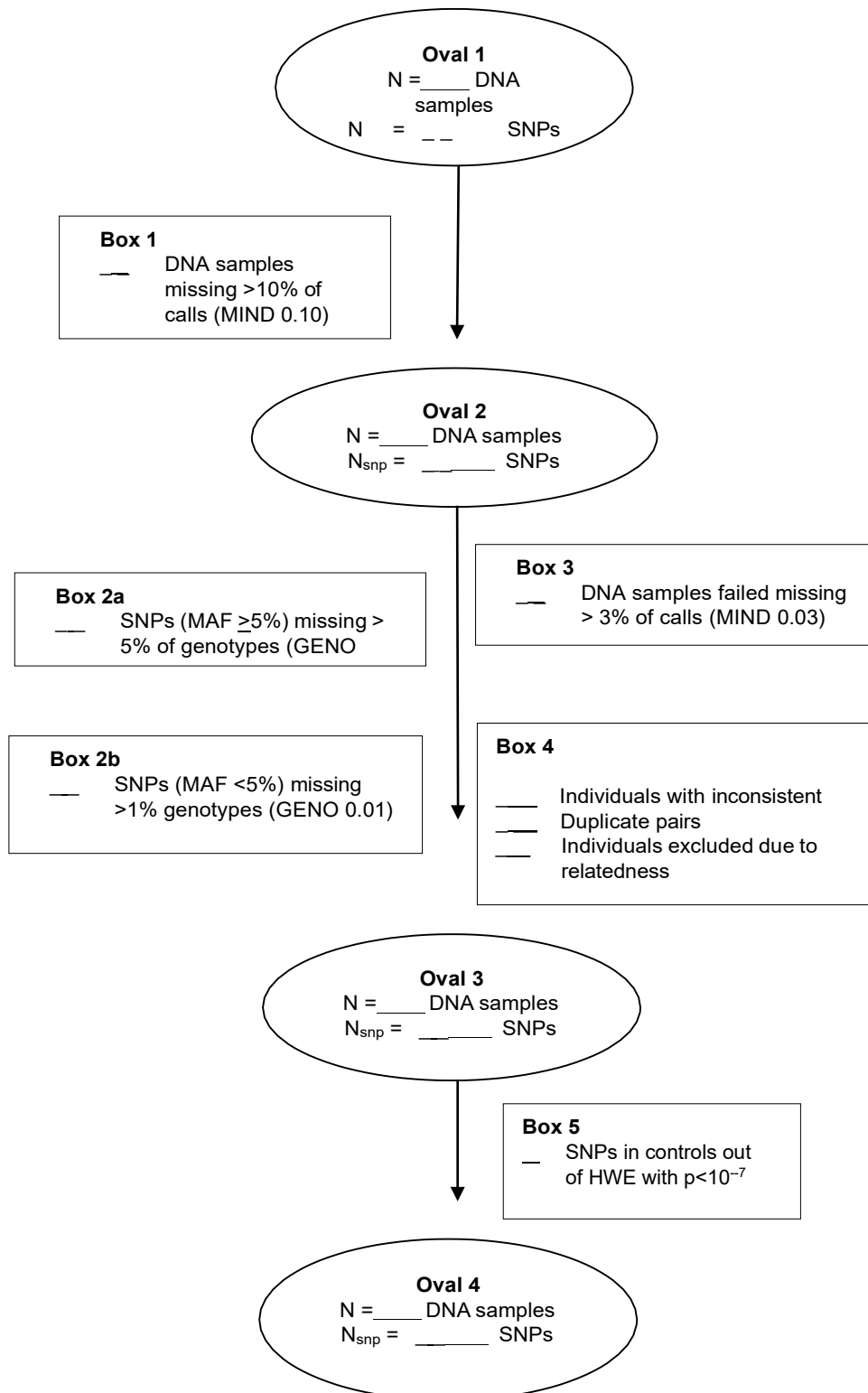
Create a text file called HWE_out.txt with the following SNP in it:

rs2968487

and type the following command:

```
plink --file GWAS_clean3 --exclude HWE_out.txt --recode --out GWAS_clean4
```

There are a number of SNPs with HWE p-values in the range of 10^{-5} to 10^{-6} in the controls. Based on above criterion they will not be excluded however, if they reach genome-wide significance during association testing they SNPs should be further investigated to ensure there is no genotyping error. You can now fill in Box 5 and Oval 4.



Answers to Questions:

Oval 1 and 2 also and Box 1 information:

```
PLINK v1.90b4.9 64-bit (13 Oct 2017)
Options in effect:
  --file GWAS
  --mind 0.10
  --out GWAS_clean_mind
  --recode

Random number seed: 1515434515
16384 MB RAM detected; reserving 8192 MB for main workspace.
Scanning .ped file... done.
Performing single-pass .bed write (6424 variants, 248 people) [Oval 1].
--file: GWAS_clean_mind-temporary.bed + GWAS_clean_mind-temporary.bim +
GWAS_clean_mind-temporary.fam written.
6424 variants loaded from .bim file.
248 people (125 males, 123 females) loaded from .fam.
1 person removed due to missing genotype data (--mind) [Box 1].
ID written to GWAS_clean_mind.irem .
Using 1 thread (no multithreaded calculations invoked).
Before main variant filters, 247 founders and 0 nonfounders present.
Calculating allele frequencies... done.
Warning: 6 het. haploid genotypes present (see GWAS_clean_mind.hh ); many
commands treat these as missing.
Total genotyping rate in remaining samples is 0.996863.
6424 variants and 247 people pass filters and QC [Oval 2].
Note: No phenotypes present.
--recode ped to GWAS_clean_mind.ped + GWAS_clean_mind.map ... done.
```

Box 2a information:

```
PLINK v1.90b4.9 64-bit (13 Oct 2017)
Options in effect:
  --file MAF_greater_5
  --geno 0.05
  --out MAF_greater_5_clean
  --recode

Random number seed: 1515435189
16384 MB RAM detected; reserving 8192 MB for main workspace.
Scanning .ped file... done.
Performing single-pass .bed write (5868 variants, 247 people).
--file: MAF_greater_5_clean-temporary.bed + MAF_greater_5_clean-temporary.bim +
MAF_greater_5_clean-temporary.fam written.
5868 variants loaded from .bim file.
247 people (125 males, 122 females) loaded from .fam.
Using 1 thread (no multithreaded calculations invoked).
Before main variant filters, 247 founders and 0 nonfounders present.
Calculating allele frequencies... done.
Warning: 6 het. haploid genotypes present (see MAF_greater_5_clean.hh ); many
commands treat these as missing.
Total genotyping rate is 0.996858.
2 variants removed due to missing genotype data (--geno) [Box2a].
5866 variants and 247 people pass filters and QC.
Note: No phenotypes present.
--recode ped to MAF_greater_5_clean.ped + MAF_greater_5_clean.map ... done.
```

Box 2b information:

```
PLINK v1.90b4.9 64-bit (13 Oct 2017)
Options in effect:
  --file MAF_less_5
  --geno 0.01
  --out MAF_less_5_clean
  --recode

Random number seed: 1515435255
16384 MB RAM detected; reserving 8192 MB for main workspace.
Scanning .ped file... done.
```

```
Performing single-pass .bed write (556 variants, 247 people).
--file: MAF_less_5_clean-temporary.bed + MAF_less_5_clean-temporary.bim +
MAF_less_5_clean-temporary.fam written.
556 variants loaded from .bim file.
247 people (125 males, 122 females) loaded from .fam.
Using 1 thread (no multithreaded calculations invoked).
Before main variant filters, 247 founders and 0 nonfounders present.
Calculating allele frequencies... done.
Total genotyping rate is 0.996913.
59 variants removed due to missing genotype data (--geno) [Box2b].
497 variants and 247 people pass filters and QC.
Note: No phenotypes present.
--recode ped to MAF_less_5_clean.ped + MAF_less_5_clean.map ... done.
```

Box 3 information:

PLINK v1.90b4.9 64-bit (13 Oct 2017)

Options in effect:

```
--file GWAS_MAF_clean
```

```
--mind 0.03
```

```
--out GWAS clean2
```

```
--recode
```

Random number seed: 1515435827

```
16384 MB RAM detected; reserving 8192 MB for main workspace.
```

```
Scanning .ped file... done.
```

Performing single-pass .bed write (6363 variants, 247 people).

```
--file: GWAS_clean2-temporary.bed + GWAS_clean2-temporary.bim +
```

GWAS_clean2-temporary.fam written.

```
6363 variants loaded from .bim file.
```

247 people (125 males, 122 females) loaded from .fam.

0 people removed due to missing genotype data (--mind) **[Box 3]**.

Using 1 thread (no multithreaded calculations invoked).

Before main variant filters, 247 founders and 0 nonfounders present.

Calculating allele frequencies... done.

```
Warning: 6 het. haploid genotypes present (see GWAS_clean2.hh ); many commands
treat these as missing.
```

Total genotyping rate is 0.99716.

6363 variants and 247 people pass filters and QC.

Note: No phenotypes present.

```
--recode ped to GWAS clean2.ped + GWAS clean2.map ... done.
```

Answer to Question 1: Why do you expect the homozygosity rate to be higher on the X chromosome in males than females?

Because males only have one allele for each SNP on the X chromosome they will appear homozygous.

Table 1: Sex check

FID	IID	PEDSEX	SNPSEX	STATUS	F
NA20506	NA20506	2	1	PROBLEM	1
NA20530	NA20530	2	1	PROBLEM	1
NA20766	NA20766	2	0	PROBLEM	0.2292
NA20771	NA20771	2	0	PROBLEM	0.2234
NA20757	NA20757	2	0	PROBLEM	0.2141

Table 2: Duplicate and Related Individuals

FID1	IID1	FID2	IID2	Z(0)	Z(1)	Z(2)	PI_HAT
M033	NA19774	M041	NA25000	0.0000	0.0000	1.0000	1.00
1344	NA12057	13291	NA25001	0.0000	0.0025	0.9975	1.00
1444	NA12739	1444	NA12749	0.0026	0.9807	0.0168	0.51
1444	NA12739	1444	NA12748	0.0026	0.9949	0.0025	0.50

F1D1- Family ID for 1st individual; ID1 - Individual ID for 1st individual; F1D2- Family ID for 2nd individual; ID2 - Individual ID for 2nd individual; Z(0)- P(IBD=0); Z(1)- P(IBD=1); Z(2)- P(IBD=2); PI HAT- $P(\text{IBD}=2)+0.5*P(\text{IBD}=1)$ (proportion IBD)

Question 2: How many duplicate pairs do you find (**hint: $\text{Pi-Hat} \approx 1$**)? Do pairs with a **$\text{Pi-Hat} \approx 1$** have to be duplicate samples? What is another explanation? What proportion would you expect a parent/ child to share IBD? Can you find any such relationship?

There are two duplicate pairs and also a trio (two parents and a child). Parent/child relationships will have a Pi-Hat value of ≈ 0.5 , but so will sibpairs. We can tell that this is a parent child relationship by examine $Z(0)$, $Z(1)$ and $Z(2)$. We will retain only one sample from each duplicate pair and the parents NA12749 and NA12748. If you perform mixed-model analysis related individuals can be retained in the sample.

Oval 3 information

```
PLINK v1.90b4.9 64-bit (13 Oct 2017)
Options in effect:
  --file GWAS_clean2
  --out GWAS_clean3
  --recode
  --remove IBS_excluded.txt
Random number seed: 1515440989
16384 MB RAM detected; reserving 8192 MB for main workspace.
Scanning .ped file... done.
Performing single-pass .bed write (6363 variants, 247 people).
--file: GWAS_clean3-temporary.bed + GWAS_clean3-temporary.bim +
GWAS_clean3-temporary.fam written.
6363 variants loaded from .bim file.
247 people (125 males, 122 females) loaded from .fam.
--remove: 244 people remaining.
Using 1 thread (no multithreaded calculations invoked).
Before main variant filters, 244 founders and 0 nonfounders present.
Calculating allele frequencies... done.
Warning: 6 het. haploid genotypes present (see GWAS_clean3.hh ); many commands
treat these as missing.
Total genotyping rate in remaining samples is 0.997225.
6363 variants and 244 people pass filters and QC [Oval 3].
Note: No phenotypes present.
--recode ped to GWAS_clean3.ped + GWAS_clean3.map ... done.
```

Table 3: Hardy Weinberg Equilibrium

Fail Cases		Fail Controls	
SNP	pvalue	SNP	pvalue
None		rs2968487	2.262e-007

```
PLINK v1.90b4.9 64-bit (13 Oct 2017)
Options in effect:
  --exclude HWE_out.txt
  --file GWAS_clean3
  --out GWAS_clean4
  --recode

Random number seed: 1515442367
16384 MB RAM detected; reserving 8192 MB for main workspace.
Scanning .ped file... done.
Performing single-pass .bed write (6363 variants, 244 people).
--file: GWAS_clean4-temporary.bed + GWAS_clean4-temporary.bim +
GWAS_clean4-temporary.fam written.
6363 variants loaded from .bim file.
244 people (123 males, 121 females) loaded from .fam.
--exclude: 6362 variants remaining.
Using 1 thread (no multithreaded calculations invoked).
Before main variant filters, 244 founders and 0 nonfounders present.
Calculating allele frequencies... done.
Warning: 6 het. haploid genotypes present (see GWAS_clean4.hh ); many commands
treat these as missing.
Total genotyping rate is 0.997229.
6362 variants and 244 people pass filters and QC [Oval 4].
Note: No phenotypes present.
--recode ped to GWAS_clean4.ped + GWAS_clean4.map ... done
```