

Practical 4 - Genome wide complex trait analysis (GCTA)

Data

We will again use the genotype and phenotype data from practical 3, which includes the height and serum transferrin levels. These data will come again in PLINK binary format. We will attempt to use the SNP marker data to build a genetic relationship matrix and estimate the proportion of phenotypic variance explained by genome wide SNPs. Chapters 26 and 27 of Lynch et al. (1998) go over the underlying methodology of the method used in this practical.

GCTA introduction

In this practical we will use the GCTA software (Yang et al., 2011, 2010) (developed at Peter Visscher's lab <cnsgenomics.com/>) to estimate variance components (from the model on slide 14 of lecture 3), where our data will be SNP markers from the whole genome or a chromosome. This software is a command line tool that has the most support on the Linux operating system.

The GCTA program uses an argument interface similar to PLINK. It is advised that you visit GCTA's website <<http://cnsgenomics.com/software/gcta/>> and familiarise yourself a little with the software (perhaps for five minutes). GCTA has many functions but one of its primary uses is variance component estimation via restricted maximum likelihood (REML).

Today we will take a look at some of its primary functions – building the SNP genetic relationship matrix (GRM), and variance component analysis. In this practical we will use GCTA to do genomic restricted maximum likelihood to estimate variance components (GREML). This will be done for the height and serum transferrin levels data seen in the previous practical. We will interrogate these results using R.

We begin with building the GRM for the data we used in the previous practical. If you are interested in the conceptual basis for the GRM built with GCTA please refer to Yang et al. (2010). Peak RAM usage for building the GRM is high and thus it is likely that your process may fail if your computer does not have enough resources (around 8GB of RAM). If your system does not have this amount of RAM, we will use a much smaller subset of the data to build the GRM as an exercise; the real GRM is already stored in your practical4/data folder and will be used in subsequent analyses. Many of the flags that you have seen in PLINK are present in GCTA. Remember that if you need help with the syntax for GCTA please take a look at the url <<http://cnsgenomics.com/software/gcta/>>. We will **not attempt to build the GRM** as the run time on a good computer is approximately ten minutes. If you would like to build the GRM in your spare time listings 1 and 2 should work on an 8GB RAM machine and a 2GB RAM machine respectively. **IMPORTANT - Remember to export your path again.**

```
1 $ gcta --bfile practical_4/data/QIMRX_cleaned --make-grm --autosome --out practical_4/results/QIMRX
```

Listing 1 Building a GRM with GCTA

```
1 $ gcta --bfile practical_4/data/QIMRX_cleaned_small --make-grm --autosome
2 $ --out practical_4/results/QIMRX_small
```

Listing 2 Building a smaller GRM with GCTA

GCTA prints output to the console whilst processing the GRM and saves this information to a .log file. Some of the key summary statistics of the GRM built above are seen below.

```
1 $ Summary of the GRM:
2 $ Mean of diagonals = 1.00083
3 $ Variance of diagonals = 9.96466e-05
4 $ Mean of off-diagonals = -0.000206112
5 $ Variance of off-diagonals = 4.47574e-05
```

Listing 3 Summary of GRM

This process generates two files QIMRX.grm.gz and QIMRX.grm.id (for older versions on Mac and Windows) or a binary version QIMRX.grm.bin with auxiliary file QIMRX.grm.N.bin (on Linux). Let's take a look at these files in R. Alternatively you can just use `head practical_4/data/QIMRX.grm.gz` from the command line.

```

1 > # Read in the gzipped GRM file
2 > grm <- read.table("practical_4/data/QIMRX.grm.gz")
3 > head(grm)
4   V1 V2   V3   V4
5   1  1 265805 0.982474300
6   2  1 265771 0.430762300
7   2  2 265791 0.997155800
8   3  1 261504 0.001788883
9   3  2 261492 0.001014439
10  3  3 261529 1.000038000

```

Listing 4 Take a look at the GRM

The gzipped GRM is stored in row form with each row having four elements. The first two columns correspond to the (i, j) position of the lower triangular matrix, the third column is the number of non missing SNPs for this row-column calculation, and the fourth contains an estimate of the genetic relatedness.

Estimating proportion of phenotypic variation due to additive genetic factors using GCTA

Let's use the GRM matrix to estimate the proportion of phenotypic variance explained by additive genome-wide SNPs for height and serum transferrin. Open the terminal or command prompt and execute the following command

```

1 $ gcta --grm practical_4/data/QIMRX --pheno practical_4/data/HT_T_X.pheno --mphen 1
2 $      --reml --out practical_4/results/QIMRX_1
3 $ gcta --grm practical_4/data/QIMRX --pheno practical_4/data/HT_T_X.pheno --mphen 2
4 $      --reml --out practical_4/results/QIMRX_2

```

Listing 5 Estimating variance components via GREML

We will use R to take a look at the output files that GCTA has calculated. Follow the listing below to read in the files and answer the following questions

```

1 > # Read in GREML result files
2 > hsq.1 <- read.table("practical_4/results/QIMRX_1.hsq",
3                     header = T, fill = T)
4 > hsq.2 <- read.table("practical_4/results/QIMRX_2.hsq",
5                     header = T, fill = T)
6 > head(hsq.1)
7   Source Variance SE
8   V(G)    0.795498 0.051003
9   V(e)    0.238877 0.041478
10  Vp      1.034375 0.028005
11  V(G)/Vp 0.769062 0.040840
12  logL    -1416.193
13  logL0   -1466.335
14  LRT 100.283
15  df      1
16  Pval    0
17  n      2836

```

Listing 6 GCTA .hsq file in R

If you prefer the command line you can do this in one line at the command line

```

1 $ # Look at the heritability file provided you are in SISG_AQG_2015 folder
2 $ cat practical_4/results/QIMRX_1.hsq
3   Source Variance SE
4   V(G)    0.795498 0.051003
5   V(e)    0.238877 0.041478
6   Vp      1.034375 0.028005
7   V(G)/Vp 0.769062 0.040840
8   logL    -1416.193
9   logL0   -1466.335
10  LRT 100.283
11  df      1
12  Pval    0
13  n      2836

```

Listing 7 Command line GCTA .hsq file

In the above listing 7 G represents genetic, e residual, and p phenotype.

Exercise 1

- What is the percentage of phenotypic variance that is explained by common SNPs for both traits?
- Are the heritability estimates significant?
- Are the heritability values what you expect?

We will now take a closer look at some of the properties of the GRM by reading using R

```

1 > # Name the columns of the GRM
2 > names(grm) <- c("IND_1", "IND_2", "SNP_NUM", "REL")
3 > dim(grm)
4      11817091      4
5 # Take out the diagonal elements
6 > grm.diag <- grm[which(grm$IND_1 == grm$IND_2), ]
7 > dim(grm.diag)
8      4861      4
9 > head(grm.diag)
10      IND_1 IND_2 SNP_NUM      REL
11      1      1  265805 0.9824743
12      2      2  265791 0.9971558
13      3      3  261529 1.0000380
14      4      4  265743 1.0015680
15      5      5  265796 1.0079730
16      6      6  265615 1.0066640
17 > # Take out the GRM off-diagonal elements
18 > grm.off.diag <- grm[which(grm$IND_1 != grm$IND_2), ]
19 > # Make a histogram of the diagonals
20 > hist(grm.diag[, 4], breaks = 2500, freq = F,
21 >      xlab = "GRM diagonals", xlim = c(0.95, 1.2), main = "")
22 > # Make a histogram of the GRM off-diagonal relatedness estimates
23 > par(mfrow = c(2, 1))
24 > hist(grm.off.diag[, 4], breaks = 2500, freq = F,
25 >      xlab = "GRM off-diagonals", xlim = c(-0.1, 0.1), main = "")
26 > hist(grm.off.diag[which(grm.off.diag[, 4] > 0.1), 4],
27 >      breaks = 200, freq = F,
28 >      xlab = "GRM off-diagonals", xlim = c(0.1, 1.1), main = "")

```

Listing 8 Looking at GRM diagonals and off-diagonals

In the above results the relatedness may have affected the estimate of heritability (Fig. 2 panel 2). We will remove the relatedness and see whether the results change. This is done with GCTA via the `--grm-cutoff 0.025` flag and the line in listing 11.

Exercise 2

- Repeat the REML analyses as in exercise 1 but with relatedness removed
- Compare the results with those in exercise 1 (using listing 9) and from Yang et al. (2010) (for trait 1 height)

```

1 > # Read in GREML result files without relatedness.
2 > grm.nr <- read.table("practical_4/data/QIMRX_nr.grm.gz")
3 > # Create the same figures as above
4 > # Note that these are only example file extensions
5 > # and may change depending on what you want to call the files
6 > hsq.1.nr <- read.table("practical_4/results/QIMRX_nr_1.hsq",
7 >      header = T, fill = T)
8 > hsq.2.nr <- read.table("practical_4/results/QIMRX_nr_2.hsq",
9 >      header = T, fill = T)

```

Listing 9 GCTA GREML .hsq result files with no relatives

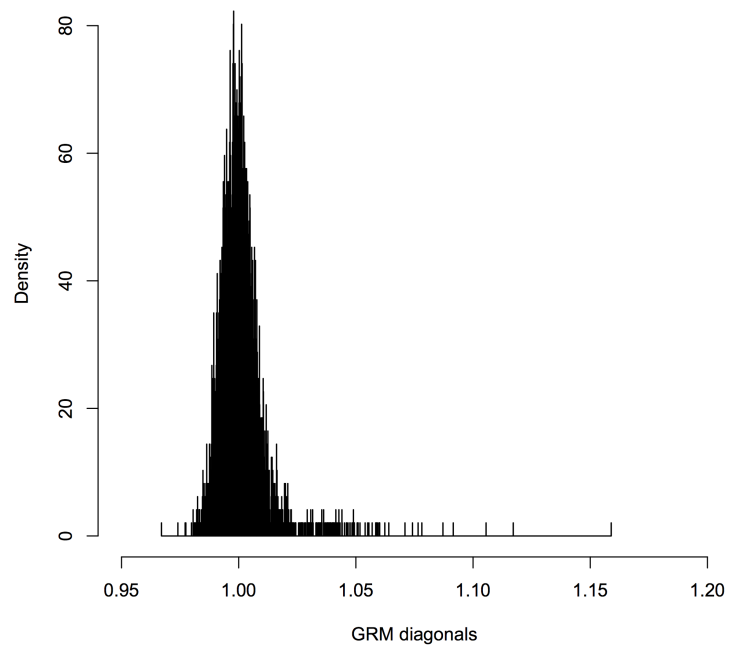


Figure 1 Plot of QIMR GRM diagonals.

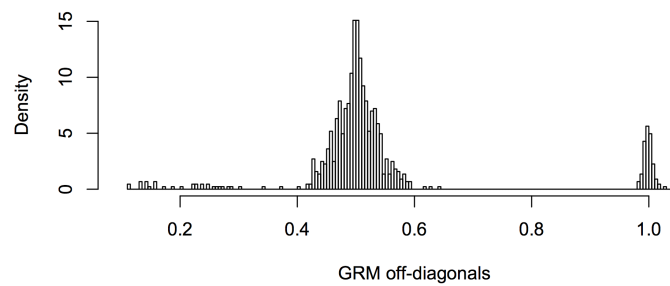
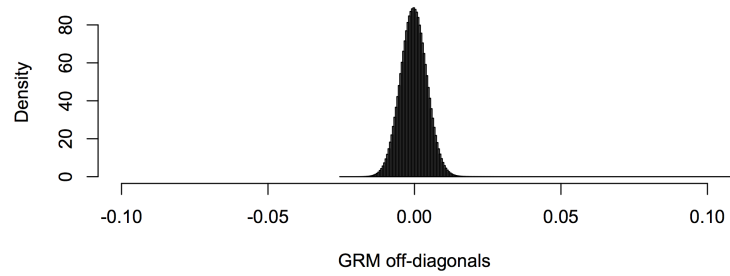


Figure 2 Plot of QIMR GRM off diagonals.

Partitioning the variance via minor allele frequency

We will now investigate partitioning variance components by creating two GRM matrices. One will be created with SNPs with lower MAFs and another with those that have higher MAFs. This will allow us to investigate (in a very imprecise way) the genetic architecture of the trait by trying to understand whether rare or common variants contribute more or less to the proportion of phenotypic variance explained by additive genetic variance (tagged by genome wide SNPs). The SNP files used are located in the `practical_4/data` directory; these files along with the `--extract` flag were used to build the two GRMs (listing 10). In order to filter on relatedness we will keep the individuals from the relatedness threshold GRM built above and the flag `--keep` (listing 10). The GRMs were built with the command and it is best if you try to build these in your own time as the process is computationally expensive.

```
1 $ gcta --bfile practical_4/data/QIMRX_cleaned --extract practical_4/data/bot_maf_snps.txt --autosome
2 $      --make-grm --keep QIMRX_nr.grm.id --out practical_4/data/QIMRX_bot_maf_snps
3 $ gcta --bfile practical_4/data/QIMRX_cleaned --extract practical_4/data/top_maf_snps.txt --autosome
4 $      --make-grm --keep QIMRX_nr.grm.id --out practical_4/data/QIMRX_top_maf_snps
```

Listing 10 Preparing the GRMs for variance partitioning

Given these two GRMs we can estimate the proportion of phenotypic variance that can be explained by additive genetic variants from low MAF SNPs versus higher MAF SNPs. Attempt to run this listing and answer the following questions

```
1 $ gcta --mgrm practical_4/data/QIMRX_multi.txt --pheno practical_4/data/HT_T_X.pheno
2 $      --mphenos 1 --reml --out practical_4/results/QIMRX_1_multi_nr
```

Listing 11 Partitioning variance components via GREML

Exercise 3

- Repeat for phenotype 2
- What do you observe for the different variance component estimates from the GRMs of low and high MAF SNPs?
- Is this what we expect?

References

- Michael Lynch, Bruce Walsh, et al. *Genetics and analysis of quantitative traits*, volume 1. Sinauer Sunderland, MA, 1998.
- Jian Yang, Beben Benyamin, Brian P McEvoy, Scott Gordon, Anjali K Henders, Dale R Nyholt, Pamela A Madden, Andrew C Heath, Nicholas G Martin, Grant W Montgomery, et al. Common snps explain a large proportion of the heritability for human height. *Nature genetics*, 42(7):565–569, 2010.
- Jian Yang, S Hong Lee, Michael E Goddard, and Peter M Visscher. Gcta: a tool for genome-wide complex trait analysis. *The American Journal of Human Genetics*, 88(1):76–82, 2011.