ARCHANA BHARDWAJ

(GWAS : via Multiple approaches) GBIO0002

GBIO002 AB 2019

Important genetic terms

Given position in the genome (i.e. locus) has several associated alleles (A and G) which produce genotypes r_A/r_G



> Haplotypes

- Combination of alleles at different loci

GBIO002 AB 2019

GWAS main philosophy

- GWAS = Genome Wide Association Studies
- IDEA = GWAS involve scan for large number of genetic markers across the whole genome of many individuals to find specific genetic variations associated with the disease and/or other phenotype
- Find the genetic variation(s) that contribute(s) and explain(s) complex diseases

GWAS visually

- GWAS tries to uncover links between genetic basis of the disease
- Which set of SNPs explain the phenotype?

Genotype	Phenotype
ATGCAGTT	control
TTGCAGTT	control
CTGCAGTT	control
ATGCGGTT	case
TTGCGGTT	case
CTGCCGTT	case
SNP	





Genotype coding

For given bi-allelic marker/SNP/loci there could be total of

3 possible genotypes given alleles A and a

Genotype	Coding
AA	0
Aa	1
aa	2

Note: A is major allele and a is minor

GBIO002 AB 2019

Relationship between Genotypes and Phenotypes

- <u>Genotype</u>: Indicates the alleles that the organism has inherited regarding a particular trait.
- <u>Phenotype</u>: The actual visible trait of the organism.



Uses of GWAS

≻Identify genes that are responsible for traits of interest:

- Humans
- Animals
- Plants





➤Understanding biological mechanisms related to the trait of interest



Human Genome Statistics

Number of Chromosomes : 23 pairs
Genome Size : 3,079,843,747 Base pairs
No of Genes : 32,185

Gene: This is a sequence of nucleotides in the DNA that codes for a molecule (e.g., a protein)



Gene Structure





Let us identify signal (in from of SNPs) from GWAS DATA

GBIO002 AB 2019

PLINK SESSION

- Data Preparation
- Quality Control
- > Clustering



PLINK: Introduction

- PLINK is whole genome association analysis tool
- PLINK has a well-documented manual to explain all features
- PLINK is available for Linux, Mac OS and MS-DOS
- •gPLINK is the other version of PLINK that provides the graphical user interface
- Command line version is faster than graphical
 PLINK

PLINK: Download

To download PLINK:

<u>http://zzz.bwh.harvard.edu/plink/download.shtm</u>
<u>l#download</u>

- Uncompress the PLINK-1.07-dos.zip
- Click on the folder. There are two files
 - test.map contains the marker information
 - test.ped contains genotype data and sample information

PLINK: File Formats(1/2)

MAP Format

Each line of the MAP file describes a single marker and must contain exactly 4 columns:

- chromosome (1-22, X, Y or 0 if unplaced)
- rs# or snp identifier
- Genetic distance (morgans)
- Base-pair position (bp units)

Column1 Column2 Column3 Column4

1 snp1 0 1 1 snp2 0 2 A centimorgan, also known and written as a genetic map unit (gmu), is, at heart, a unit of probability. One cM is equal to the distance of two genes that gives a recombination frequency of one percent.

PLINK: File Formats(2/2)

•PED Format. This file is a white-space (space or tab) delimited file: the first six columns are mandatory:

Family ID, Individual ID, Paternal ID, Maternal ID, Sex (1=male; 2=female; other=unknown), phenotype

Column5 Column4 Column3 🔕 Column2 Column1 lumn6 1 0 1 1 Α G TG 1 1 1 0 0 ΔC 1 GG 1 CC 0 0 2 1 C 1 0 0 Δ 2 G 1 CC 0 0 1 GRIOODAR 1 0

Binary format:-> BED, BIM, and FAM

Transposed text format : and TFAM

We will be using BED, BIM and FAM in coming session

GBIO002 AB 2019

PLINK: File Formats (Example)

*.ped

1.00					
		n	n	3	n
	÷		.,	a	μ.

FID	IID	PID	MID	Sex	Ρ	rs1	rs2	rs3
1	1	0	0	2	1	СТ	AG	AA
2	2	0	0	1	0	cc	AA	AC
3	3	0	0	1	1	CC	AA	AC

Chr	SNP	GD	BPP
1	rs1	0	870000
1	rs2	0	880000
1	rs3	0	890000

*.fam

FID	IID	PID	MID	Sex	Ρ
1	1	0	0	2	1
2	2	0	0	1	0
3	3	0	0	1	1

*.bed Contains binary version of the SNP info of the *.ped file. (not in a format readable for humans)

*.bim

Chr	SNP	GD	BPP	Allele 1	Allele 2	
1	rs1	0	870000	С	т	
1	rs2	0	880000	А	G	
1	rs3	0	890000	A	с	

Covariate file

FID	IID	C1	C2	C3
1	1	0.00812835	0.00606235	-0.000871105
2	2	-0.0600943	0.0318994	-0.0827743
3	3	-0.0431903	0.00133068	-0.000276131

Legend					
FID	Family ID	rs{x}	Alleles per subject per SNP		
IID	Individual ID	Chr	Chromosome		
PID	Paternal ID	SNP	SNP name		
MID	Maternal ID	GD	Genetic distance (morgans)		
Sex	Sex of subject	BPP	Base-pair position (bp units)		
Ρ	Phenotype	C{x}	Covariates (e.g., Multidimensiona Scaling (MDS) components)		

PLINK: Data input (map and ped) Type command : plink --file test --noweb

Provide you detail information of SNPs count, individuals count and check for the missingness and frequency test

```
Options in effect:
        --file test
        --noweb
  (of 2) markers to be included from [ test.map ]
6 individuals read from [ test.ped ]
 individuals with nonmissing phenotypes
Assuming a disease phenotype (1=unaff, 2=aff, 0=miss)
Missing phenotype value is also -9
3 cases, 3 controls and 0 missing
6 males, 0 females, and 0 of unspecified sex
Before frequency and genotyping pruning, there are 2 SNPs
6 founders and 0 non-founders found
Total genotyping rate in remaining individuals is 1
0 SNPs failed missingness test ( GENO > 1 )
0 SNPs failed frequency test ( MAF < 0 )
After frequency and genotyping pruning, there are 2 SNPs
After filtering, 3 cases, 3 controls and 0 missing
After filtering, 6 males, 0 females, and 0 of unspecified sex
```

PLINK: Data input (bed format)

Convert map and ped into binary format

plink.exe --file test --make-bed --out test --noweb

Read binary file in plink

plink.exe –bfile test

Reading map (extended format) from [test.bim] 2 markers to be included from [test.bim] Reading pedigree information from [test.fam] 5 individuals read from [test.fam] 6 individuals with nonmissing phenotypes Assuming a disease phenotype (1=unaff, 2=aff, 0=miss) Missing phenotype value is also -9 3 cases, 3 controls and 0 missing 6 males, 0 females, and 0 of unspecified sex Reading genotype bitfile from [test.bed] Detected that binary PED file is v1.00 SNP-major mode Before frequency and genotyping pruning, there are 2 SNPs 6 founders and 0 non-founders found Total genotyping rate in remaining individuals is 1 0 SNPs failed missingness test (GENO > 1) 0 SNPs failed frequency test (MAF < 0) After frequency and genotyping pruning, there are 2 SNPs After filtering, 3 cases, 3 controls and 0 missing After filtering, 6 males, 0 females, and 0 of unspecified sex

Analysis finished: Sun Oct 20 18:49:10 2019

Comparision

2 (of 2) markers to be included from [test.map] 6 individuals read from [test.ped] 6 individuals with nonmissing phenotypes Assuming a disease phenotype (1=unaff, 2=aff, 0=miss) Missing phenotype value is also -9 3 cases, 3 controls and 0 missing 6 males, 0 females, and 0 of unspecified sex Before frequency and genotyping pruning, there are 2 SNPs 6 founders and 0 non-founders found Total genotyping rate in remaining individuals is 1 0 SNPs failed missingness test (GENO > 1) 0 SNPs failed frequency test (MAF < 0) After frequency and genotyping pruning, there are 2 SNPs After filtering, 3 cases, 3 controls and 0 missing After filtering, 6 males, 0 females, and 0 of unspecified sex

Reading map (extended format) from [test.bim] 2 markers to be included from [test.bim] Reading pedigree information from [test.fam] 6 individuals read from [test.fam] 6 individuals with nonmissing phenotypes Assuming a disease phenotype (1=unaff, 2=aff, 0=miss) Missing phenotype value is also -9 cases, 3 controls and 0 missing 6 males, 0 females, and 0 of unspecified sex Reading genotype bitfile from [test.bed] Detected that binary PED file is v1.00 SNP-major mode Before frequency and genotyping pruning, there are 2 SNPs 6 founders and 0 non-founders found Total genotyping rate in remaining individuals is 1 0 SNPs failed missingness test (GENO > 1) 0 SNPs failed frequency test (MAF < 0) After frequency and genotyping pruning, there are 2 SNPs After filtering, 3 cases, 3 controls and 0 missing After filtering, 6 males, 0 females, and 0 of unspecified sex

Analysis finished: Sun Oct 20 18:49:10 2019

TEXT output

BED OUTPUT

SIMILAR Conclusion

Example data

Download the example data from the course website (PLINK FOLDER)

- TSI_JPT_chr20_case_control.bed
- TSI_JPT_chr20_case_control.bim
- TSI_JPT_chr20_case_control.fam
- TSI_JPT_chr20_pheno_header.txt
- TSI_JPT_chr20_pheno.txt

By looking into file extension, **BED FORMAT**



GBIO002 AB 2019



Detection of LD, population stratification (comes under Filteration step) Lets Perform Quality filteration

Quality control processes

Missing genotype

Hardy-Weinberg Equilibrium

Minor Allele frequency

Linkage disequilibrium pruning

Missing genotype



Missing Genotypes

 To generate a list genotyping/missingness rate statistics:

- > plink.exe --bfile TSI_JPT_chr20_case_control --missing -noweb
 - plink.imiss
 plink.lmiss
 OUTPUT in running directory

It provides the detail missingness by individual and by SNP (locus), respectively.

Lets read log file

Reading map (extended format) from [TSI JPT chr20 case control.bim] 36302 markers to be included from [TSI JPT chr20 case control.bim] Reading pedigree information from [TSI_JPT_chr20_case_control.fam] 174 individuals read from [TSI_JPT_chr20_case_control.fam] 174 individuals with nonmissing phenotypes Assuming a disease phenotype (1=unaff, 2=aff, 0=miss) Missing phenotype value is also -9 88 cases, 86 controls and 0 missing 88 males, 86 females, and 0 of unspecified sex Reading genotype bitfile from [TSI_JPT_chr20_case_control.bed] Detected that binary PED file is v1.00 SNP-major mode Before frequency and genotyping pruning, there are 36302 SNPs 174 founders and 0 non-founders found Writing individual missingness information to [plink.imiss] Writing locus missingness information to [plink.lmiss] Total genotyping rate in remaining individuals is 0.998399 Ø SNPs failed missingness test (GENO > 1) ← 0 SNPs failed frequency test (MAF < 0) After frequency and genotyping pruning, there are 36302 SNPs After filtering, 88 cases, 86 controls and 0 missing After filtering, 88 males, 86 females, and 0 of unspecified sex

Clustering based on Missing Genotypes

- Systematic batch effects that induce missingness in parts of the sample will induce correlation between the patterns of missing data that different individuals display
- One approach to detecting correlation in these patterns, that might possibly idenity such biases, is to cluster individuals based on their identityby-missingness (IBM).

plink.exe --bfile TSI_JPT_chr20_case_control --cluster-missing --noweb

•which creates the files:

- plink.matrix.missing
- plink.cluster3.missing

which have similar formats to the corresponding IBS clustering files.

Missing Rate Per Person

•The initial step in all data analysis is to exclude individuals with too much missing genotype

data. This option is set as follows:

> plink.exe --bfile TSI_JPT_chr20_case_control --mind 0.1 --noweb

which means exclude with more than 10% missing genotypes.

•A line in the terminal output will appear, indicating how many individuals were removed due to low genotyping. If any individuals were removed, a file called **plink.irem** will be created, listing the Family and Individual IDs of these removed individuals.

Missing Rate Per SNP

Subsequent analyses can be set to automatically exclude SNPs on the basis of missing genotype rate, with the --geno option: the default is to include all SNPS (i.e. --geno 1).
To include only SNPs with a 90% genotyping rate (10% missing) use

> plink.exe --bfile TSI_JPT_chr20_case_control --geno 0.1 --noweb

 As with the --maf option, these counts are calculated after removing individuals with high missing genotype rates.

Hardy-Weinberg Equilibrium (1/2)

 To generate a list of genotype counts and Hardy-Weinberg test statistics for each SNP, use the option:

plink.exe --bfile TSI_JPT_chr20_case_control --hardy --noweb which creates a file: *plink.hwe.* The file has the following format

SNP SNP identifier

TEST Code indicating sample

- A1 Minor allele code
- A2 Major allele code

GENO Genotype counts:11/12/22

O(HET) observed hetrozygosity

E(HET) Expected hetrozygosity

P H-W p-value
Hardy-Weinberg Equilibrium (2/2)

•To exclude markers that failure the Hardy-Weinberg test at a specified significance threshold, use the option:

• plink.exe --bfile TSI_JPT_chr20_case_control --hwe 0.001 --noweb •By default this filter uses an exact test. The standard asymptotic (1 df genotypic chi-squared test) can be requested with the --hwe2 option instead of --hwe. •The following output will appear in the console window and in **plink.log**, detailing how many SNPs failed the Hardy-Weinberg test, for the sample as a whole, and (when PLINK) has detected a disease phenotype) for cases and controls separately:

```
Reading map (extended format) from [ TSI_JPT_chr20_case_control.bim ]
36302 markers to be included from [ TSI JPT chr20 case control.bim ]

    Input data

Reading pedigree information from [ TSI JPT chr20 case control.fam ]
174 individuals read from [ TSI_JPT_chr20_case_control.fam ]
174 individuals with nonmissing phenotypes
Assuming a disease phenotype (1=unaff, 2=aff, 0=miss)
Missing phenotype value is also -9
88 cases, 86 controls and 0 missing
88 males, 86 females, and 0 of unspecified sex
Reading genotype bitfile from [ TSI_JPT_chr20_case_control.bed ]
Detected that binary PED file is v1.00 SNP-major mode
Before frequency and genotyping pruning, there are 36302 SNPs
174 founders and 0 non-founders found
18 markers to be excluded based on HWE test ( p <= 0.001 )
        20 markers failed HWE test in cases
        18 markers failed HWE test in controls
Total genotyping rate in remaining individuals is 0.998399
0 SNPs failed missingness test ( GENO > 1 )
0 SNPs failed frequency test (MAF < 0)
After frequency and genotyping pruning, there are 36284 SNPs
After filtering, 88 cases, 86 controls and 0 missing
After filtering, 88 males, 86 females, and 0 of unspecified sex
```

Allele Frequency

how often an form of a gene shows up in a population over several generations

the number of copies of a particular allele divided by the number of copies of all alleles at the genetic place in a population. 6

GG

88

20

gg

ß

0000

B

Gg

8 8



Allele Frequency

•To generate a list of minor allele frequencies (MAF) for each SNP, based on all founders in the sample:

- plink.exe --bfile TSI_JPT_chr20_case_control --freq --noweb
- This will create a file: plink.frq with five columns:
 - CHR Chromosome
 - SNP SNP identifier
 - A1 Allele 1 code (minor allele)
 - A2 Allele 2 code (major allele)
 - MAF Minor allele frequency

NCHROBS Non-missing allele count

Minor Allele Frequency

•Once individuals with too much missing genotype data have been excluded, subsequent analyses can be set to automatically exclude SNPs on the basis of MAF (minor allele frequency):

plink --bfile TSI_JPT_chr20_case_control --maf 0.05 --noweb
 It means only include SNPs with MAF >= 0.05.
 The default value is 0.01. This quantity is based only on founders (i.e. individuals for whom the paternal and maternal individual codes and both 0).

option is appropriately counts alleles for X chromosome SNPs.

Reading map (extended format) from [TSI_JPT_chr20_case control.bim] 36302 markers to be included from [TSI_JPT_chr20_case_control.bim] lnput data Reading pedigree information from [TSI_JPT_chr20_case_control.fam] 174 individuals read from [TSI_JPT_chr20_case_control.fam] 174 individuals with nonmissing phenotypes Assuming a disease phenotype (1=unaff, 2=aff, 0=miss) Missing phenotype value is also -9 88 cases, 86 controls and 0 missing 88 males, 86 females, and 0 of unspecified sex Reading genotype bitfile from [TSI_JPT_chr20_case_control.bed] Detected that binary PED file is v1.00 SNP-major mode Before frequency and genotyping pruning, there are 36302 SNPs 174 founders and 0 non-founders found <u>Total genotyping rate in remaining individuals is 0.998399</u> Ø SNPs failed missingness test (GENO > 1) 7826 SNPs failed frequency test (MAF < 0.05) After frequency and genotyping pruning, there are 28476 SNPs After filtering, 88 cases, 86 controls and 0 missing After filtering, 88 males, 86 females, and 0 of unspecified sex

Count SNPs under MAF < 0.01 ?

```
Reading map (extended format) from [ TSI JPT chr20 case control.bim ] —
36302 markers to be included from [ TSI JPT chr20 case control.bim ]
Reading pedigree information from [ TSI JPT chr20 case control.fam ]
174 individuals read from [ TSI JPT chr20 case control.fam ]
174 individuals with nonmissing phenotypes
Assuming a disease phenotype (1=unaff, 2=aff, 0=miss)
Missing phenotype value is also -9
88 cases, 86 controls and 0 missing
88 males, 86 females, and 0 of unspecified sex
Reading genotype bitfile from [ TSI_JPT_chr20_case_control.bed ]
Detected that binary PED file is v1.00 SNP-major mode
Before frequency and genotyping pruning, there are 36302 SNPs
174 founders and 0 non-founders found
Total genotyping rate in remaining individuals is 0.998399
0 SNPs failed missingness test ( GENO > 1 )
4625 SNPs failed frequency test ( MAF < 0.01 )
After frequency and genotyping pruning, there are 31677 SNPs
After filtering, 88 cases, 86 controls and 0 missing
After filtering, 88 males, 86 females, and 0 of unspecified sex
```

Input data



Linkage disequilibrium (LD): This is a measure of non- random association between alleles at different loci at the same chromosome in a given population. SNPs are in LD when the frequency of association of their alleles is higher than expected under random assortment. LD concerns patterns of correlations between SNPs.

Linkage disequilibrium pruning (1/2)

•Sometimes it is useful to generate a pruned subset of SNPs that are in approximate linkage equilibrium with each other. This can be achieved via two commands:

--indep which prunes based on the variance inflation factor (VIF), which recursively removes SNPs within a sliding window;

-- indep-pairwise which is similar, except it is based only on pairwise genotypic correlation.

•The VIF pruning routine is performed:

Linkage disequilibrium pruning (2/2)

•Each is a simlpe list of SNP IDs; both these files can subsequently be specified as the argument for a -extract or --exclude command.

•The parameters for --indep are: window size in SNPs (e.g. 50), the number of SNPs to shift the window at each step (e.g. 5), the VIF threshold. The VIF is 1/(1-R^2) where R^2 is the multiple correlation coefficient for a SNP being regressed on all other SNPs simultaneously.

•That is, this considers the correlations between SNPs but also between linear combinations of SNPs.

plink.exe --bfile TSI_JPT_chr20_case_control --indep 50 5 2 --noweb

Reading map (extended format) from [TSI JPT chr20 case control.bim] 36302 markers to be included from [TSI JPT chr20 case control.bim] Reading pedigree information from [TSI JPT chr20 case control.fam] 174 individuals read from [TSI JPT chr20 case control.fam] 174 individuals with nonmissing phenotypes Assuming a disease phenotype (1=unaff, 2=aff, 0=miss) Missing phenotype value is also -9 88 cases, 86 controls and 0 missing 88 males, 86 females, and 0 of unspecified sex Reading genotype bitfile from [TSI JPT chr20 case control.bed] Detected that binary PED file is v1.00 SNP-major mode Before frequency and genotyping pruning, there are 36302 SNPs 174 founders and 0 non-founders found Total genotyping rate in remaining individuals is 0.998399 0 SNPs failed missingness test (GENO > 1) 0 SNPs failed frequency test (MAF < 0) After frequency and genotyping pruning, there are 36302 SNPs After filtering, 88 cases, 86 controls and 0 missing After filtering, 88 males, 86 females, and 0 of unspecified sex Performing LD-based pruning... Writing pruned-in SNPs to [plink.prune.in] Writing pruned-out SNPs to [plink.prune.out] Scanning from chromosome 20 to 20

Scan region on chromosome 20 from [rs4814683] to [rs6011394] For chromosome 20, 29618 SNPs pruned out, 6684 remaining Input data

How many snp in LD with window size "100", "200" ?

clustering

plink.exe --bfile TSI_JPT_chr20_case_control --cluster --noweb

which generates four output files:

plink.cluster0

plink.cluster1

plink.cluster2

plink.cluster3

that contain similar information but in different formats. The

The *.cluster0 file contains some information on the clustering process. This file can be safely ignored by most users.

The *.cluster1 file contains information on the final solution, listed by cluster.

The *.cluster2 file contains the same information but listed one line per individual

The *.cluster3 file is in the same format as cluster2 (one line per individual) but contains all solutions (i.e. every step of the clustering from moving from N clusters each of 1 individual (leftmost column after family and individual ID) to 1 cluster (labelled 0) containing all N individuals (the final, rightmost column)

Plink.cluster1

SOL-0 NA18946_NA18946 NA18955_NA18955 NA19087_NA19087 NA19002_NA19002

There is only one cluster.

What if we have more than one cluster?



We will perform this analysis in other R package

Association Analysis

Case/control

Multiple-testing correction

Basic case/control association test

To perform a standard case/control association analysis, use the option: plink.exe --bfile TSI_JPT_chr20_case_control --assoc --noweb which generates a file

plink.assoc

which contains the fields:

- CHR Chromosome
- SNP SNP ID
- BP Physical position (base-pair)
- A1 Minor allele name (based on whole sample)
- F_A Frequency of this allele in cases
- F_U Frequency of this allele in controls
- A2 Major allele name
- CHISQ Basic allelic test chi-square (1df)
- P Asymptotic p-value for this test
- OR Estimated odds ratio (for A1, i.e. A2 is reference)

plink.exe --bfile TSI_JPT_chr20_case_control --assoc --noweb

Reading map (extended format) from [TSI JPT chr20 case control.bim] 36302 markers to be included from [TSI_JPT_chr20_case_control.bim] Input data Reading pedigree information from [TSI JPT chr20 case control.fam] 174 individuals read from [TSI JPT chr20 case control.fam] 174 individuals with nonmissing phenotypes Assuming a disease phenotype (1=unaff, 2=aff, 0=miss) Missing phenotype value is also -9 88 cases, 86 controls and 0 missing 88 males, 86 females, and 0 of unspecified sex Reading genotype bitfile from [TSI JPT chr20 case control.bed] Detected that binary PED file is v1.00 SNP-major mode Before frequency and genotyping pruning, there are 36302 SNPs 174 founders and 0 non-founders found Total genotyping rate in remaining individuals is 0.998399 0 SNPs failed missingness test (GENO > 1) 0 SNPs failed frequency test (MAF < 0) After frequency and genotyping pruning, there are 36302 SNPs After filtering, 88 cases, 86 controls and 0 missing After filtering, 88 males, 86 females, and 0 of unspecified sex Writing main association results to [plink.assoc]

Adjustment for multiple testing

To generate a file of adjusted significance values that correct for all tests performed and other metrics, use the option:

plink.exe --*bfile TSI_JPT_chr20_case_control* --*assoc* --*adjust* --*noweb* which generates the file

plink.adjust

which contains the fields

CHR	Chromosome number
SNP	SNP identifer
UNADJ	Unadjusted p-value
GC	Genomic-control corrected p-values
BONF	Bonferroni single-step adjusted p-values
HOLM	Holm (1979) step-down adjusted p-values
SIDAK_SS	Sidak single-step adjusted p-values
SIDAK_SD	Sidak step-down adjusted p-values
FDR_BH	Benjamini & Hochberg (1995) step-up FDR control
FDR_BY	Benjamini & Yekutieli (2001) step-up FDR control
the fill of the second second line	a ter (fine and the second term the second term term term term term term term term

This file is sorted by significance value rather than genomic location, the most significant results being at the top.

- 1. Open plink.assoc.adjusted
- 2. Sort column « BONF »
- 3. Count SNPs has p value 0.05, 0.01?

Let us visualize GWAS result

LETS INSTALL R Pakcage

Open R window
 Install.packages(qqman)
 Load in library
 library(qqman)

> gwas <- data.frame(read.table(file="plink.assoc",header=TRUE))</pre>

manhattan(gwas, main = "Manhattan Plot", ylim = c(0, 10),col="blue")





It has data from Chr 20

Manhattan (multi chromosome view)



To work on cluster, we have Storage issue.

We will be working on small dataset.

It will have SNPs information for 1 chromosome.

R Session (multiple packages)

- > Data Preparation
- Preprocessing (similar to PLINK)
- > ANALYSIS : (similar to PLINK)
 - **Principal Component Analysis**
 - **Genome-Wide Association**
 - □ Functional insights into candidate markers

Introduction to GenABEL (1/2)

This library allows to do complete GWAS workflow

- GWAS data and corresponding attributes (SNPs, phenotype, sex, etc.) are stored in data object gwas.data-class
- The object attributes could be accessed with @

- phenotype data: gwaa_object@phdata
- number of people in study: gwaa_object@gtdata@nids

Introduction to GenABEL(2/2)

- number of SNPs: gwaa_object@gtdata@nsnps
- SNP names: gwaa_object@gtdata@snpnames
- Chromosome labels:gwaa_object@gtdata@chromosome
- SNPs map positions: gwaa_object@gtdata@map

GenABEL is an R library developed to facilitate Genome-Wide Association analysis of binary and quantitative traits.

Features of GenABEL :

- specific facilities for storage and manipulation of large data
- QC
- Maximum Likelihood estimation of linear, logistic and Cox regression on Genome-wide scale
- Specific functions to analyze and display the results

CONNECT TO SERVER (For Windows)

1. GO to Link https://putty.org/ Click to download





PuTTY is an SSH and telnet client, developed originally by Simon Tatham for the Windows platform. PuTTY is open source software that is available with source code and is developed and supported by a group of volunteers.

You can download PuTTY here.

Below suggestions are independent of the authors of PuTTY. They are not to be seen as endorsements by the PuTTY project.



Bitvise SSH Client

Bitvise SSH Client is an SSH and SFTP client for Windows. It is developed and supported professionally by Bitvise. The SSH Client is robust, easy to install, easy to use, and supports all features supported by PuTTY, as well as the following:

- graphical SFTP file transfer;
- single-click Remote Desktop tunneling;
- auto-reconnecting capability;
- dynamic port forwarding through an integrated proxy;
- an FTP-to-SFTP protocol bridge.

Bitvise SSH Client is free to use. You can download it here.



Bitvise SSH Server

Bitvise SSH Server is an SSH, SFTP and SCP server for Windows. It is robust, easy to install, easy to use, and works well with a variety of SSH clients, including Bitvise SSH Client, OpenSSH, and PUTTY. The SSH Server is developed and supported professionally by Bitvise.

You can download Bitvise SSH Server here.

CONNECT TO SERVER

- Open putty
- User id : username@ms801.montefiore.ulg.ac.be
- Enter Password

ms801 to ms818

	ms801.montefiore.ulg.ac.be - PuTTY -	>
	Using username "bhardwaj". bhardwaj@ms801.montefiore.ulg.ac.be's password:	
	Welcome to ubuntu 16.04.3 LIS (GNU/Linux 4.4.0-142-generic x86_64)	
p	* Documentation: https://help.ubuntu.com * Management: https://landscape.canonical.com	
	* Support: https://ubuntu.com/advantage	
1	325 packages can be updated. 8 updates are security updates.	
P	Last login: Mon Oct 21 12:56:04 2019 from 10.9.108.36 bhardwaj@ms801:~\$	

CONNECT TO SERVER (For Unix)

Open terminal

Enter command

ssh username@ms801.montefiore.ulg.ac.be

You will be connected to the server

GO TO « R »

💣 ms801.montefiore.ulg.ac.be - PuTTY	_	×
-bash: syntax error near unexpected token `doParallel' bhardwaj@ms801:~\$ library(SNPRelate) -bash: syntax error near unexpected token `SNPRelate' bhardwaj@ms801:~\$ R		^
R version 3.2.3 (2015-12-10) "Wooden Christmas-Tree" Copyright (C) 2015 The R Foundation for Statistical Computing Platform: x86_64-pc-linux-gnu (64-bit)		
R is free software and comes with ABSOLUTELY NO WARRANTY. You are welcome to redistribute it under certain conditions. Type 'license()' or 'licence()' for distribution details.		
Natural language support but running in an English locale R is a collaborative project with many contributors. Type 'contributors()' for more information and		
<pre>'citation()' on how to cite R or R packages in publications. Type 'demo()' for some demos, 'help()' for on-line help, or 'help.start()' for an HTML browser interface to help. Type 'g()' to guit R.</pre>		
		~

GRI0002 &R 2019

TRANSFER FILES TO SERVER

https://winscp.net/eng/download.php

Folder : CMI

CLID					
SELF					
<u>H</u> ost name:				Po	o <u>r</u> t number:
ms802.montefi	ore.ulg.ac.	be			22
<u>U</u> ser name:			Password:		
bhardwaj					
Edit				Ad	vanced

READY TO WORK !!

REQUIRED PACKAGES

- ➢ library(GenABEL)←
- library(dplyr)
- > library(snpStats)
- library(doParallel)
- library(SNPRelate)

If not installed, use following command

```
if (!requireNamespace("BiocManager", quietly = TRUE))
```

```
install.packages("BiocManager")
```

```
BiocManager::install("snpStats")
BiocManager::install("dplyr")
BiocManager::install("doParallel")
BiocManager::install("SNPRelate")
```

Already installed on server . No need to install it.
STUDY DESIGN

- Data Download : Directory under name of CMI
- Background : It has data from three population.
- TASK : What are significant SNPs in context to specific trait.

Files names :

gwas_indian.fam gwas_indian.bim gwas_indian.bed gwas_chinese.fam gwas_chinese.bim gwas_chinese.bed gwas_malay.fam gwas_malay.bim gwas_malay.bed GWASfunction.R conversionTable.RData 122Chinese_282lipids.txt 120Indian_282lipids.txt 117Malay_282lipids.txt

Population A

Population B

Population C

Lipids « omics » data

Read all files in R (Population data and lipids information)

Lets look into Data

PLINK-converted .bed, .fam and .bim Illumina files from each of the three ethnic groups.

Exercise

- Read data in PLINK
- How many SNPs present in each File ?
- How many indivduals belong to three datasets? Give count of males and females?

From PLINK to R package

> library(snpStats)

 We will use the function read.plink from the package snpStats and work on the resulting objects throughout the rest of the session.

pathC <- paste("/home/bhardwaj/22-OCT/gwas_chinese", c(".bed", ".bim", ".fam"), sep = "")
SNP_C <- read.plink(pathM[1], pathM[2], pathM[3])</pre>

pathI <- paste("/home/bhardwaj/22-OCT/gwas_indian", c(".bed", ".bim", ".fam"), sep = "")
SNP_I <- read.plink(pathI[1], pathI[2], pathI[3])</pre>

pathM <- paste("/home/bhardwaj/22-OCT/gwas_malay", c(".bed", ".bim", ".fam"), sep = "")
SNP_M <- read.plink(pathC[1], pathC[2], pathC[3])</pre>

Check no of input columns

```
if( ncol(SNP_C$genotypes) != ncol(SNP_I$genotypes)) {
    stop("Different number of columns in input files detected. This is not allowed.")
}
if( ncol(SNP_I$genotypes) != ncol(SNP_M$genotypes)) {
    stop("Different number of columns in input files detected. This is not allowed.")
}
```

 Append the \$genotypes element from each object row-wise (we've checked columns are equal in number)

SNP <- rbind(SNP_M\$genotypes, SNP_I\$genotypes, SNP_C\$genotypes)</pre>

Data merge (1/2)

Merge the three SNP datasets

SNP <- rbind(SNP_M\$genotypes, SNP_I\$genotypes, SNP_C\$genotypes)

Genotype will be merged

Lets merge all 3 maps are based on the same ordered set of SNPs) map <- SNP_M\$map</p>

> head(map)						
,	chromosome	snp.name	cM	position	allele.1	allele.2
rs4477212	1	rs4477212	NA	82154	<na></na>	A
kgp15717912	1	kgp15717912	NA	534247	<na></na>	C
kgp7727307	1	kgp7727307	NA	569624	<na></na>	C
kgp15297216	1	kgp15297216	NA	723918	<na></na>	G
rs3094315	1	rs3094315	NA	752566	G	A
rs3131972	1	rs3131972	NA	752721	A	G
>						

Data merge (2/2)

Get colnames from map data variable

colnames(map) <- c("chr", "SNP", "gen.dist", "position", "A1", "A2")</pre>

Lets merge all 3 fam are based on the same ordered set of SNPs)

SNP_M\$fam<- rbind(SNP_M\$fam, SNP_I\$fam, SNP_C\$fam)</pre>

	head (S	NP_M\$fam)					
		pedigree	member	father	mother	sex	affected
М1	11060903	M11060903	M11060903	NA	NA	1	NA
М1	11072004	M11072004	M11072004	NA	NA	2	NA
Μ1	11071302	M11071302	M11071302	NA	NA	1	NA
М1	11072317	M11072317	M11072317	NA	NA	1	NA
М1	11071307	M11071307	M11071307	NA	NA	2	NA
Μ1	11073011	M11073011	M11073011	NA	NA	2	NA
>							

Lets Read another Dataset : Lipids information of three population

Next we import and merge the three lipid data sets (stored as .txt)

lipidsMalay <- read.delim("/home/bhardwaj/22-OCT/117Malay_282lipids.txt", row.names = 1)</pre>

lipidsIndian <- read.delim("/home/bhardwaj/22-OCT/120Indian_282lipids.txt", row.names = 1)</pre>

lipidsChinese <- read.delim("/home/bhardwaj/22-OCT/122Chinese_282lipids.txt", row.names =
1)</pre>

all(Reduce(intersect, list(colnames(lipidsMalay), colnames(lipidsIndian), colnames(lipidsChinese))) == colnames(lipidsMalay))

> dim(lipidsChinese)
[1] 122 283

merge the three lipid data sets (stored as .txt)

lip <- rbind(lipidsMalay, lipidsIndian, lipidsChinese)</pre>

Country Information of samples

```
country <- sapply(list(SNP_M, SNP_I, SNP_C), function(k){
    nrow(k$genotypes)
})</pre>
```

> country [1] 110 105 108

sample.id Country						
1 M11060903	Μ					
2 M11072004	Μ					
3 M11071302	Μ					
4 M11072317	Μ					
5 M11071307	Μ					
6 M11073011	Μ					

Matched data in GWAS and Lipids: Country Information of samples

matchingSamples <- intersect(rownames(lip), rownames(SNP))
SNP <- SNP[matchingSamples,]
lip <- lip[matchingSamples,]
origin <- origin[match(matchingSamples, origin\$sample.id),]</pre>

Save session :

genData <- list(SNP = SNP, MAP = map, LIP = lip)
save(genData, origin, file = "PhenoGenoMap.RData")</pre>

Get GDS structure

Merge the three SNP datasets

SNP_M\$genotypes <- rbind(SNP_M\$genotypes, SNP_I\$genotypes, SNP_C\$genotypes)
colnames(map) <- c("chr", "SNP", "gen.dist", "position", "A1", "A2")
SNP_M\$fam<- rbind(SNP_M\$fam, SNP_I\$fam, SNP_C\$fam)</pre>

Rename SNPs present in the conversion table into rs IDs

load("conversionTable.RData")
mappedSNPs <- intersect(SNP_M\$map\$SNP, names(conversionTable))
newIDs <- conversionTable[match(SNP_M\$map\$SNP[SNP_M\$map\$SNP %in% mappedSNPs],
names(conversionTable))]
SNP_M\$map\$SNP[rownames(SNP_M\$map) %in% mappedSNPs] <- newIDs</pre>

write.plink("convertGDS", snps = SNP_M\$genotypes)

Writing FAM file to convertGDS.fam Writing extended MAP file to convertGDS.bim Writing BED file to convertGDS.bed (SNP-major mode) NULL

What are quality check we can do

- Based on MAF ? (As performed in plink)
- Based on HWE ? (As performed in plink)

Lets do it one by one and analyse the data.

Pre-Processing or QC

- In a nutshell, the pre-processing of the data consists in discarding SNPs with call rate < 1 or MAF < 0.1 discarding samples with call rate < 100%, IBD kinship coefficient > 0.1 or inbreeding coefficient |F| > 0.1
- Call rate is the proportion of SNPs (or samples) that were genotyped. For example, a call rate of 0.95 for a particular SNP (sample) means 5% of the values are missing. (as performed in PLINK)
- Minor-allele frequency (MAF) denotes the proportion of the least common allele for each SNP.
- Of course, it is harder to detect associations with rare variants and this is why we select against low MAF values. Most GWA studies typically report MAF thresholds of 0.05. (as performed in PLINK)

Load Packages

library(snpStats)
library(doParallel)
library(SNPRelate)
library(GenABEL)
library(dplyr)
source("GWASfunction.R")
load("PhenoGenoMap.RData")

Minor Allele Frequency

Lets use SNP call rate of 100%, MAF of 0.1 (very stringent) (as performed in PLINK)

maf <- 0.1 *# selected parameter*

callRate <- 1 # selected parameter

SNPstats <- col.summary(genData\$SNP)

SNPstats Detail Information

"Calls" "Call.rate" "Certain.calls" "RAF" "MAF" "P.AA" "P.AB" "P.BB" "z.HWE" "rs3094315" 319 1 1 0.860501567398119 0.139498432601881 0.0188087774294671 0.241379310344828 0.739811912225705 0.0968676724046835 "rs3131972" 319 1 1 0.766457680250784 0.233542319749216 0.0313479623824451 0.404388714733542 0.564263322884012 2.31429230658791 "kgp15275285" 319 1 1 0.855799373040752 0.144200626959248 0.0438871473354232 0.200626959247649 0.755485893416928 -3.34227919125702 "kgp5225889" 319 1 1 0.877742946708464 0.122257053291536 0.0188087774294671 0.206896551724138 0.774294670846395 -0.642784289569623 "rs11240777" 319 1 1 0.69435736677116 0.30564263322884 0.0689655172413793 0.473354231974922 0.457680250783699 2.05783717257905 "kgp8975187" 319 1 1 0.713166144200627 0.286833855799373 0.0626959247648903 0.448275862068966 0.489028213166144 1.70937237845253 "kgp5074587" 319 1 1 0.811912225705329 0.188087774294671 0.0376175548589342 0.300940438871473 0.661442006269592 -0.262046989104515

Lets Filter based on MAF

maf_call <- with(SNPstats, MAF > maf & Call.rate == callRate)

TRUE : Passed MAF FALSE : Failed MAF

maf call FALSE FALSE TRUE FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE TRUE TRUE TRUE FALSE TRUE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE TRUE FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE TRUE FALSE TRUE TRUE FALSE FALSE FALSE FALSE FALSE TRUE FALSE TRUE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE TRUE FALSE TRUE FALSE FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE FALSE TRUE FALSE FALSE FALSE FALSE TRUE FALSE FALSE TRUE FALSE TRUE TRUE TRUE

genData\$SNP <- genData\$SNP[,maf_call]
genData\$MAP <- genData\$MAP[maf_call,]
SNPstats <- SNPstats[maf_call,] # Read output
head(SNPstats)</pre>

Hardy-Weinberg equilibrium (HWE 1/2)

- Next, we need to consider samples that exhibit excessive heterozygosity technically speaking, deviations from the Hardy-Weinberg equilibrium (HWE),
- Overall, large and small |F| might indicate poor sample quality or inbreeding, respectively. We will exclude samples with |F|> 0.1.
 - Lets run Sample call rate & heterozygosity

callMat <- !is.na(genData\$SNP) Sampstats <- row.summary(genData\$SNP)

> head (Sa	mpstats)			
n an	Call.rate	Certain.calls	Heterozygosity	hetF
M11051101	1	1	0.3308424	0.09523034
M11050404	1	1	0.3165761	0.13424505
M11051105	1	1	0.3750000	-0.02552947
M11051901	1	1	0.4028533	-0.10170104
M11052507	1	1	0.3519022	0.03763720
M11050514	1	1	0.3491848	0.04506858
>				

Hardy-Weinberg equilibrium (HWE 2/2)

hetExp <- callMat %*% (2 * SNPstats\$MAF * (1 - SNPstats\$MAF)) # Hardy-Weinberg heterozygosity (expected)

hetObs <- with(Sampstats, Heterozygosity * (ncol(genData\$SNP)) * Call.rate) **# Hardy-Weinberg heterozygosity** (observed value)

Sampstats\$hetF <- 1-(hetObs/hetExp) (# hetF prediction)</pre>

Use sample call rate of 100%, het threshold of 0.1 (very stringent)

- het <- 0.1 # Set cutoff for inbreeding coefficient;</p>
- het_call <- with(Sampstats, abs(hetF) < het & Call.rate == 1)</pre>

> head(het_call) [,1] M11051101 TRUE M11050404 FALSE M11051105 TRUE M11051901 FALSE M11052507 TRUE M11050514 TRUE

- genData\$SNP <- genData\$SNP[het_call,]
 </pre>
- genData\$LIP <- genData\$LIP[het_call,]
 </pre>

Identity by Descent IBD

- Finally, we will investigate relatedness among samples using the kinship coefficient based on identity by descent (IBD).
- These functions from the package SNPRelate require GDS files. For this reason we first need to aggregate the .bed, .fam and .bim files from the three populations into convertGDS.
- The function snpgdsBED2GDS2 creates the GDS necessary for this part of the analysis.
- To determine the kinship coefficient between pairs of samples we will use a subset of uncorrelated SNPs in order to have unbiased estimates.

 For this purpose, we will use linkage disequilibrium (LD) as a measure of correlation between markers. LD ranges from 0 to 1, the higher its value the more likely two SNPs co-segregate and therefore correlate.

 Here, we will utilize the subset of SNPs with LD < 0.2 to determine the IBD kinship coefficient.

```
ld <- .2
kin <- .1
snpgdsBED2GDS(bed.fn = "convertGDS.bed", bim.fn = "convertGDS.bim",
fam.fn = "convertGDS.fam", out.gdsfn = "myGDS", cvt.chr = "char")
genofile <- snpgdsOpen("myGDS", readonly = F)
gds.ids <- read.gdsn(index.gdsn(genofile, "sample.id"))
gds.ids <- sub("-1", "", gds.ids)
add.gdsn(genofile, "sample.id", gds.ids, replace = T)
geno.sample.ids <- rownames(genData$SNP)</pre>
```

Lets filter for LD

Lets filter for LD (1/3)

snpSUB <- snpgdsLDpruning(genofile, ld.threshold = ld, sample.id = geno.sample.ids, snp.id = colnames(genData\$SNP))

> snpSUB <- snp	gdsLDpruning(genofile, ld.threshold = ld,
+	sample.id = geno.sample.ids,
+	snp.id = colnames(genData\$SNP))
SNP pruning bas	ed on LD:
Excluding 3528 S	SNPs on non-autosomes
Excluding 0 SNP	(monomorphic: TRUE, < MAF: NaN, or > missing rate: NaN)
Working space: 2	221 samples, 1472 SNPs
Using 1 (CPU	J) core
Sliding wind	low: 500000 basepairs, Inf SNPs
LD thresh	old: 0.2
Chromosome 0:	3.34%, 167/5000
167 SNPs are sel	ected in total.

Lets filter for LD (2/3)

```
snpset.ibd <- unlist(snpSUB, use.names = F)
# And now filter for MoM
ibd <- snpgdsIBDMoM(genofile, kinship = T,
sample.id = geno.sample.ids,
snp.id = snpset.ibd,
num.thread = 1)</pre>
```

```
> ibd <- snpgdsIBDMoM(genofile, kinship = T,</p>
            sample.id = geno.sample.ids,
+
            snp.id = snpset.ibd,
+
            num.thread = 1)
+
IBD analysis (PLINK method of moment) on SNP genotypes:
Excluding 4833 SNPs on non-autosomes
Excluding 0 SNP (monomorphic: TRUE, < MAF: NaN, or > missing rate: NaN)
Working space: 221 samples, 167 SNPs
    Using 1 (CPU) core
PLINK IBD: the sum of all working genotypes (0, 1 and 2) = 20713
PLINK IBD:
            Mon Oct 21 16:40:51 2019
                                          0%
PLINK IBD:
            Mon Oct 21 16:40:51 2019
                                          100%
```

Lets filter for LD (3/3)

ibdcoef <- snpgdsIBDSelection(ibd)
ibdcoef <- ibdcoef[ibdcoef\$kinship >= kin,]

Lets filter samples

```
related.samples <- NULL
while (nrow(ibdcoef) > 0) {
    # count the number of occurrences of each and take the top one
    sample.counts <- sort(table(c(ibdcoef$ID1, ibdcoef$ID2)), decreasing = T)
    rm.sample <- names(sample.counts)[1]
    cat("Removing sample", rm.sample, "too closely related to",
        sample.counts[1], "other samples.\n")</pre>
```

remove from ibdcoef and add to list ibdcoef <- ibdcoef[ibdcoef\$ID1 != rm.sample & ibdcoef\$ID2 != rm.sample,] related.samples <- c(as.character(rm.sample), related.samples)</pre>

After pre-processing, we are left with few samples . Note that your sample size might differ slightly as the LD pruning procedure is stochastic.

Filtered Data

genData\$SNP <- genData\$SNP[!(rownames(genData\$SNP) %in% related.samples),]
genData\$LIP <- genData\$LIP[!(rownames(genData\$LIP) %in% related.samples),]</pre>

Next :

What is the population structure in studied samples of india, chinese and malay population?

Principal Component Analysis

- Now that we are done with the pre-processing, it might be a good idea to examine the largest sources of variation in the genotype data and look out for outliers or clustering patterns, using <u>Principal Component Analysis (PCA)</u>.
- Because we are working with S4 objects, we will be using the PCA function from SNPRelate, snpgdsPCA. Let's plot the first two principal components (PCs).

Lets run the PCA

pca <- snpgdsPCA(genofile, sample.id = geno.sample.ids, snp.id = snpset.ibd, num.thread = 1)

pctab <- data.frame(sample.id = pca\$sample.id, PC1 = pca\$eigenvect[,1], PC2 =
pca\$eigenvect[,2], stringsAsFactors = F)</pre>

origin <- origin[match(pca\$sample.id, origin\$sample.id),]</pre>

pcaCol <- rep(rgb(0,0,0,.3), length(pca\$sample.id)) # Set black for chinese
pcaCol[origin\$Country == "I"] <- rgb(1,0,0,.3) # red for indian
pcaCol[origin\$Country == "M"] <- rgb(0,.7,0,.3) # green for malay</pre>

png("PCApopulation.png", width = 500, height = 500)
plot(pctab\$PC1, pctab\$PC2, xlab = "PC1", ylab = "PC2", col = pcaCol, pch = 16)
abline(h = 0, v = 0, lty = 2, col = "grey")
legend("top", legend = c("Chinese", "Indian", "Malay"), col = 1:3, pch = 16, bty =
"n")
dev.off()

OUTPUT



- As expected, the filtered SNP markers clearly delineate the indian population clearly different from other two populations.
- The results also suggest that Chinese and Malay are closer to each other than to Indian (this observation would be much better addressed with e.g. hierarchical clustering).

Genome-Wide Association

- Note that the glm function is used to determine the significance of association between each SNP and the trait of interest. (as performed in PLINK)
- glm is much more versatile than lm since it conducts Gaussian, Poisson, binomial and multinomial regression / classification, depending on how your trait of interest is distributed (all lipids in the phenotype file are Gaussian).
- This GWA function will not create a variable, but rather write a .txt summary table listing the coefficient estimate \beta , t and the corresponding P-value for each SNP, alongside with the corresponding genomic coordinates.

Choose trait for association analysis,

```
target <- "Cholesterol"
phenodata <- data.frame("id" = rownames(genData$LIP),
"phenotype" = scale(genData$LIP[,target]), stringsAsFactors = F)</pre>
```

Conduct GWAS

start <- Sys.time()
GWAA(genodata = genData\$SNP, phenodata = phenodata, filename = paste(target,
".txt", sep = ""))
Sys.time() - start # benchmark</pre>

Ams802.montefiore.ulg.ac.be - PuTTY	-		×	
1 > target <- "Cholesterol"			^	
<pre>> phenodata <- data.frame("id" = rownames(genData\$LIP), + "phenotype" = scale(genData\$LIP[.target]), stringsAsFactors = F</pre>				
<pre>> start <- Sys.time()</pre>				
> GWAA(genodata = genData\$SNP, phenodata = phenodata, filename = p	aste (target,		
.txt", sep = ""))				
1472 SNPs included in analysis.				n
221 samples included in analysis.				~~~
socket cluster with 2 nodes on host 'localhost'				
GWAS SNPs 1-148 (10% finished)				e
GWAS SNPs 149-296 (20% finished)				
GWAS SNPS 297-444 (30% finished)				
GWAS SNPS 445-592 (40% Finished)				
GWAS SNPs 741-888 (60% finished)				
GWAS SNPs 889-1036 (70% finished)				
GWAS SNPs 1037-1184 (80% finished)				
GWAS SNPs 1185-1332 (90% finished)				
GWAS SNPs 1333-1472 (100% finished)				
[1] "Done."				14
> Sys.time() - start # benchmark				r
Time difference of 4.262571 secs				

- Once finished, we can visualize the results using the so-called Manhattan plots.
- All we need is to load the .txt summary table written in the previous step, add a column with -\log_{10} (P) and plot these significance estimates against the genomic coordinates of all SNPs.

Lets draw Manhattan plot

GWASout <- read.table(paste(target, ".txt", sep = ""), header = T, colClasses = c("character", rep("numeric",4))) GWASout\$type <- rep("typed", nrow(GWASout)) GWASout\$Neg_logP <- -log10(GWASout\$p.value) GWASout <- merge(GWASout, genData\$MAP[,c("SNP", "chr", "position")]) GWASout <- GWASout[order(GWASout\$Neg_logP, decreasing = T),]</pre>

png(paste(target, ".png", sep = ""), height = 500, width = 1000)
GWAS_Manhattan(GWASout)
dev.off()



Chromosome

We see that a total of one SNPs is indicating association (none passes the 'hard' threshold). rs9509213
quantile-quantile (Q-Q) plot

- Before proceeding with this one hits, it is helpful to constrast the distribution of the resulting P-values against that expected by chance, as to ensure there is no confounding systemic bias.
- QQ plot using GenABEL estlambda function
 - png(paste(target, "_QQplot.png", sep = ""), width = 500, height = 500)
 - Iambda <- estlambda(GWASout\$t.value**2, plot = T, method = "median")</p>
 - > dev.off()

The resulting Q-Q plot clearly depicts a trend line (\lambda = 0.99, red) overlapping with x = y (black) and a slight deviation in the right tail, so we can be confident about our results.



Lambda should range from 0 to 1 (ideal condition).

The resulting Q-Q plot clearly depicts a trend line (\lambda = 1, red), overlapping with x = y (black) and a slight deviation in the right tail.

so we can be more confident about our results.

What if lambda value < 1 ?

How QQ plot looks like ?





Reason for this lesser lamda :

In this study, we worked on the small sample size. If we include entire dataset, lamda could be increased

Functional insights into candidate markers

- Next will try to find the functional relevance of these one candidate SNPs by searching for genes in their vicinity, using the <u>USCS Genome Browser</u> (enter Genome Browser, insert the SNP ID in the text box, enter and zoom out).
- I found that rs9509213 lands right on CRYL1 (crystalline lambda 1, intron sequence), rendering it an interesting candidate for follow-up studies.



USCS Genome Browser



Oct. 11, 2019 - New "group auto-scale" option

Oct. 9, 2019 - Expanded CRISPR track released for human (hg...

More news...

Subscribe

The UCSC Genome Browser is developed and maintained by the Genome Bioinformatics Group, a cross-departmental team within the UCSC Genomics Institute.

International Human Genome Project consortium completed the first working draft of the human genome assembly, forever ensuring free public access to the genome and the information it contains. A few weeks later, on July 7, 2000, the newly assembled genome was released on the web at http://genome.ucsc.edu, along with the initial prototype of a graphical viewing tool, the UCSC Genome Browser. In the ensuing years, the website has grown to include a broad collection of vertebrate and model organism assemblies and annotations, along with a large suite of tools for viewing, analyzing and downloading data.

Genomics





Genomic comparision with other species

Exercice : Identify gene name of following SNPs

rs11083846 rs11636802 rs13397985 rs13401811 rs1679013 rs17483466 rs210142

Advance Functional interpretation

Open browser <u>https://amp.pharm.mssm.edu/Enrichr/</u>	and paste gene
Enrichr Enrich	Login Register 7.644.583 lists analyzed
Analyze What's New? Libraries Find a Gene About Help	229,071 terms 123 libraries
Input data	
Choose an input file to upload. Either in BED format or a list of genes. For a quantitative set, add a comma and the level of membership of that gene. The membership level is a number between 0.0 and 1.0 to represent a weight for each gene, where the weight of 0.0 will completely discard the gene from the enrichment analysis and the weight of 1.0 is the maximum. Try an example BED file. Choose File No file chosen	optionally followed by ip. Try two examples:
	0 gene(s) entered
Enter a brief description for the list in case you want to share it. (Optional)	Submit
Contribute	Submit
Please acknowledge Enrichr in your publications by citing the following references: Chen EY, Tan CM, Kou Y, Duan Q, Wang Z, Meirelles GV, Clark NR, Ma'ayan A. Enrichr: interactive and collaborative H	TML5

https://amp.pharm.mssm.edu/Enrichr/



Login | Register

Transcription Pathways Ontologies Disease/Drugs Cell Types Misc Legacy Crowd

Description Sample gene list (5 genes)



- Click on pathways
- Click on ontologies
 - Cellular , Molecular and Biological functions??

Pathways Transcription Ontologies Diseases/Drugs

Cell Types Misc Legacy

Crowd

Description No description available (1252 genes)



GO Biological Process 2018



Table Clu



Click the bars to sort. Now sorted by **p-value ranking**.

SVG PNG JPG

positive regulation of transcription, DNA-templated (GO:0045893)

regulation of transcription from RNA polymerase II promoter (GO:0006357)

positive regulation of gene expression (GO:0010628)

positive regulation of transcription from RNA polymerase II promoter (GO:0045944)

regulation of signal transduction by p53 class mediator (GO:1901796)

cellular response to DNA damage stimulus (GO:0006974)

regulation of transcription, DNA-templated (GO:0006355)

DNA repair (GO:0006281)

regulation of apoptotic process (GO:0042981)

negative regulation of transcription, DNA-templated (GO:0045892)

Let us try Unknown dataset

- 1. Download gene _ist .txt from course website.
- 2. Predict Gene ontology
- **3. Predict different pathways**
- 4. Any idea, genes playing specific biological functional ??

Additional Reference

https://onlinelibrary.wiley.com/doi/full/10.1002/mpr.1608



FROM TODAY SESSION

- Performed different steps of GWAS and identified signals
- Performed annotation of identified signals (different ways)
- Genomic information
- PATHWAYS
- Gene ontology

NEXT SESSION

POST GWAS : ADVANCE FUNCTIONAL ANNOTATION