

1 **BICF Nanocourse: Genome Analysis**  
2 **Workshop for: Exome-/genome-sequencing in population based studies**

3  
4 Julia Kozlitina  
5 [Julia.Kozlitina@UTSouthwestern.edu](mailto:Julia.Kozlitina@UTSouthwestern.edu)  
6 May 2, 2019  
7

8 **Today we are going to:**  
9

- 10 - Convert VCF file into format suitable for association analysis  
11 - Perform QC of population sequencing data  
12 - Perform association analysis  
13     o Single-variant test for common variants  
14     o Gene-based test for rare variants  
15 - Summarize and assess the quality of the results  
16

17 **This tutorial will use the following software:**  
18

- 19 PLINK            (<https://www.cog-genomics.org/plink2>)  
20                    Command-line genetic analysis toolset  
21  
22 Haploview        (<https://www.broadinstitute.org/haploview/haploview>)  
23                    Graphical tool for viewing PLINK results and SNP analysis  
24  
25 Locuszoom        ([locuszoom.org/](https://locuszoom.org/))  
26                    Graphical tool for visualizing regional association results  
27  
28 EFACTS           (<https://genome.sph.umich.edu/wiki/EFACTS>)  
29                    versatile software pipeline to perform various statistical tests for identifying  
30                    genome-wide association from sequence data  
31

32 #####

33 **1. Getting started**

34 #####

- 35  
36 Log into BioHPC and log into the compute node  
37 - Set up a WebGUI (<https://portal.biohpc.swmed.edu/terminal/webgui>) session on BioHPC  
38 - Launch via "connect with VNC client", open using TurboVNC  
39     (<https://sourceforge.net/projects/turbovnc/>).  
40 - You can also launch the session by "connect via web" but copying and pasting may not  
41     work under this mode.  
42 - Open a terminal window - you should be in your home directory.  
43     /home2/trainXX  
44

45 Now prepare the environment and data

46

47 1. Copy session4 material into your directory and work from there

48

```
49 cp -r /archive/nanocourse/genome_analysis/shared/session4 .  
50 cd session4
```

51

52 2. Load the necessary modules

53

```
54 module load R/3.4.1-gccmkl  
55 module load plink/1.9  
56 module load locuszoom/1.4  
57 module load epacts/3.3.2
```

58

59 Check that PLINK is working by typing:

60

```
61 plink
```

62

63 This will provide a description of PLINK, basic syntax example and a list of some commands.

64

65

```
66 #####
```

67

## 2. Datasets

```
68 #####
```

69

70 The data used in this exercise are from 661 African and 503 European ancestry individuals from  
71 the 1000 Genomes project (<http://www.internationalgenome.org>). From the whole-genome  
72 sequencing data, a subset of ~171,000 bi-allelic SNPs, mostly in exonic regions, was extracted. The  
73 genotypes along with a simulated disease status, quantitative phenotype and some covariates are  
74 contained in the following files.

75

76	1kg_Exome.vcf.gz	genotype data for 1164 individuals
77	1kg_data.covar	additional covariates to be used in analysis
78	1kg_data.BIN.pheno	case-control status
79	1kg_data.QT.pheno	quantitative phenotype
80	1kg_sample_info.txt	sample information

81

82

```
83 #####
```

84

## 3. Explore the data and convert VCF to PLINK BED/BIM/FAM Format

```
85 #####
```

86

87 3.1 Go to the data directory by typing at the command prompt:

88

```
cd ~/session4/data
```

89

90

91 3.2 Check the VCF file  
92 `gunzip -c lkg_Exome.vcf.gz | head -265 | cut -f 1-10`  
93

94 3.3 Convert the VCF file to PLINK format for QC and analysis:  
95 `plink --vcf lkg_Exome.vcf.gz --keep-allele-order --double-id --make-bed --out`  
96 `lkg_data`  
97

98 The `--keep-allele-order` option keeps the REF and ALT alleles as defined in the VCF file.  
99 PLINK by default forces the more common allele to be REF allele (A2), and the less common  
100 allele to be ALT allele (A1), regardless of which is REF and ALT in the VCF. To learn more, type:  
101 `plink --help --keep-allele-order`  
102

103 The `--make-bed` command above will produce the following output files:  
104

105	<code>lkg_data.bed</code>	genotype data in binary format
106	<code>lkg_data.bim</code>	chromosomal map file for SNPs included in <code>.bed</code> file
107	<code>lkg_data.fam</code>	family pedigree information
108	<code>lkg_data.log</code>	log file containing all the commands and options

109

110 (a) BED is a binary file that contains the genotype information, similar to a standard PED file,  
111 but in machine-readable format (it takes much less storage space (10%), and allows for  
112 faster processing in PLINK). If we could read it, it would contain the genotype data with 1  
113 line per individual and 1 column for each SNP:

```
114 A A A A C C C C  
115 A A A A C C C C  
116 A A A A C C C C  
117 A A A A T C C C  
118 ...  
119
```

120  
121 (b) BIM file contains information on the SNPs included in the `.bed` file. The first 6 columns are  
122 CHR, SNP, cM, Position, Allele 1 (minor), Allele 2 (major). To view the first few lines of the  
123 BIM file, type:

124 `head lkg_data.bim`  
125

126  
127 which should produce the following output:

```
128  
129 1 rs75333668 0 762320 T C  
130 1 rs201186828 0 865545 A G  
131 1 rs148711625 0 865584 A G  
132 1 rs146327803 0 865625 A G  
133 ...  
134
```

135 To see how many variants are in the genotype file:

136 `wc -l lkg_data.bim`  
137

138 (c) FAM file contains the pedigree information, the same as the first 6 columns of a standard  
139 PED file. It has 6 columns: family ID, individual ID, paternal ID, maternal ID, sex (1 = Male, 2  
140 = Female, 0 = unknown), and phenotype (1=unaffected control, 2=affected case, 0 or -9 =  
141 missing).

```
142 head 1kg_data.fam
143
144 HG00096    HG00096    0    0    0    -9
145 HG00097    HG00097    0    0    0    -9
146 HG00099    HG00099    0    0    0    -9
147 HG00100    HG00100    0    0    0    -9
148 HG00101    HG00101    0    0    0    -9
149 ...
```

151  
152 Notice that sex variable is set to unknown for all individuals (since this information was not  
153 provided in the VCF). We can update this information using the following command:

```
154 plink --bfile 1kg_data --keep-allele-order --update-sex
155 1kg_sample_info.txt 3 --make-bed --out 1kg_data_temp
```

```
156 head 1kg_data_temp.fam
157
158 HG00096    HG00096    0    0    1    -9
159 HG00097    HG00097    0    0    2    -9
160 HG00099    HG00099    0    0    2    -9
161 HG00100    HG00100    0    0    2    -9
162 HG00101    HG00101    0    0    1    -9
```

163  
164  
165 (d) **Phenotype file.** Instead of the phenotype in the 6<sup>th</sup> column of FAM file, it is possible to  
166 load a different phenotype to the binary file set from a white-space- or tab-delimited file,  
167 with at least three columns: FID, IID, Phenotype value, using the option `--pheno`  
168 (additional columns will be ignored unless `--pheno-name` is specified):

```
169 plink --bfile 1kg_data --pheno 1kg_data.BIN.pheno
```

170  
171  
172 To view the file, type:

```
173 head 1kg_data.BIN.pheno
174 head 1kg_data.QT.pheno
175
176 FID        IID        Pheno      lPheno
177 HG00096    HG00096    54.82      4
178 HG00097    HG00097    57.4       4.05
179 HG00099    HG00099    24.79      3.21
180 HG00100    HG00100    31.89      3.46
181 HG00101    HG00101    17.17      2.84
182 ..
```

183  
184  
185 (e) **Covariate file.** Covariate files are similar to phenotype files, and contain additional  
186 covariates that will be used in analysis. To load the covariates, use the option `--covar`.  
187

```
188
189 plink --bfile lkg_data --covar lkg_data.covar
190
191 head lkg_data.covar
192
193 FID IID Sex AGE PC1 PC2
194 HG00096 HG00096 1 55 -0.0136039 -0.0147257
195 HG00097 HG00097 2 63 -0.0131045 -0.0141718
196 HG00099 HG00099 2 52 -0.0136478 -0.0128483
197 HG00100 HG00100 2 52 -0.0130089 -0.0139981
198 HG00101 HG00101 1 37 -0.0130738 -0.0130549
```

```
201 #####
```

#### 4. Some pointers to working with PLINK

```
203 #####
```

- 205 - PLINK always generates a LOG file, which includes the details of the implemented
- 206 commands, and any warning messages. It is very useful for checking if the software is
- 207 successfully completing commands.
- 208
- 209 - Exact syntax and spelling is **very important**
- 210 e.g., "--bfile" is not the same as "-bfile"
- 211
- 212 - PLINK has excellent web documentation
- 213 PLINK 1.07: <http://pngu.mgh.harvard.edu/purcell/plink/>
- 214 PLINK 2.0: <https://www.cog-genomics.org/plink2>

```
216 #####
```

#### 5. Data QC

```
218 #####
```

220 **Note:** In this exercise, we assume that samples have already undergone standard quality control  
221 steps: all gender discordant samples, duplicates, discordant duplicate pairs, have been excluded.  
222 These steps can be implemented using the commands below. For more details, see Anderson et  
223 al., 2010 [PMID: 21085122].

```
224 plink --bfile lkg_data --check-sex --out lkg_data
225 plink --bfile lkg_data --het --out lkg_data
```

227

228 5.1 Exclude SNPs with a missing genotype call rate of >10% (--geno 0.1) and individuals with a  
229 missing genotype call rate of >10% across SNPs (--mind 0.1).

```
230 plink --bfile lkg_data_temp --keep-allele-order --geno 0.1 --mind 0.1 --
231 make-bed --out lkg_data_temp2
```

233

234 5.2 Exclude SNPs with minor allele count (mac) <5 (since single-variant tests have low power to  
235 detect the effects of extremely rare variants).

```
236  
237 plink --bfile 1kg_data_temp2 --keep-allele-order --mac 5 --make-bed --out  
238 1kg_data_temp3
```

239  
240 5.3 Compute HWE p-values (we have to do it separately for AFR and EUR):

```
241  
242 plink --bfile 1kg_data_temp3 --keep-allele-order --filter 1kg_sample_pop.txt AFR  
243 --hardy midp gz --out 1kg_AFR  
244
```

245 This will produce a file 1kg\_AFR.hwe.gz. Extract all SNPs with  $P(\text{HWE}) < 1e-6$ :

```
246  
247 gunzip -c 1kg_AFR.hwe.gz | awk '{if($9 <= 1e-6) print $0}' >  
248 SNPs_fail_HWE_AFR.txt  
249 head SNPs_fail_HWE_AFR.txt  
250  
251 awk '{print $2}' SNPs_fail_HWE_AFR.txt > SNPs_fail_AFR.txt  
252
```

253 Follow the same steps for EUR population:

```
254  
255 plink --bfile 1kg_data_temp3 --keep-allele-order --filter 1kg_sample_pop.txt EUR  
256 --hardy midp gz --out 1kg_EUR  
257  
258 gunzip -c 1kg_EUR.hwe.gz | awk '{if($9 <= 1e-6) print $0}' >  
259 SNPs_fail_HWE_EUR.txt  
260 awk '{print $2}' SNPs_fail_HWE_EUR.txt > SNPs_fail_EUR.txt  
261
```

262 Combine the two lists:

```
263  
264 cat SNPs_fail_AFR.txt SNPs_fail_EUR.txt > SNPs_fail_HWE.txt  
265
```

266 Filter out the failed SNPs:

```
267  
268 plink --bfile 1kg_data_temp3 --keep-allele-order --exclude SNPs_fail_HWE.txt --  
269 make-bed --out 1kg_data.pass  
270
```

271 Remove the temp files:

```
272  
273 rm *temp*  
274
```

275 5.4. (Optional) Summarize the allele frequencies:

```
276  
277 plink --bfile 1kg_data --keep-allele-order --filter 1kg_sample_pop.txt AFR --  
278 freq gz --out 1kg_AFR  
279 plink --bfile 1kg_data --keep-allele-order --filter 1kg_sample_pop.txt EUR --  
280 freq gz --out 1kg_EUR  
281  
282  
283
```

284 #####

285 **6. PCA**

286 #####

287

288 We will use pre-computed PCs included in the `1kg_data.covar` file for analysis. To calculate  
289 principal components, you could follow the steps below. For more details, see Anderson et al.  
290 (2010).

291

- 292 1. Extract high-quality independent variants (i.e., not in linkage disequilibrium) from the QC'd  
293 genotype file:

```
294 #plink --bfile 1kg_data.pass --geno 0.05 --maf 0.05 --indep-pairwise 50 5  
295 0.2 --out 1kg_data.pass  
296
```

297

298 This will produce a list of variants to include and to exclude:

```
299 1kg_data.prune.in  
300 1kg_data.prune.out  
301
```

- 302 2. Extract the pruned-in variants from your genotype file and, optionally, from a reference file

303 (e.g. 1000 genomes):

```
304  
305 #plink --bfile 1kg_data.pass --extract 1kg_pca_snps.txt --make-bed --out  
306 1kg_pca_temp  
307 #plink --bfile 1000G_data_full --extract 1kg_data.prune.in --make-bed --  
308 out 1000G_pca  
309 #plink --bfile 1000G_pca --bmerge 1kg_pca_temp --make-bed --out 1kg_  
310 merged  
311
```

- 312 3. Compute the PC's:

313

```
314 #plink --bfile 1kg_merged --pca --out 1kg_data_pca
```

315 This will produce two files containing PCA summary and the first 20 PC scores for each  
316 individual in the sample:

```
317 1kg_data.eigenval ## Summary  
318 1kg_data.eigenvec ## PC's  
319
```

320 #####

321 **7. Association analysis for binary trait (case/control status)**

322 #####

323

324 Create a directory to write plink output:

325

```
326 cd ~/session4  
327 mkdir plink_out  
328
```

329 Update the phenotype value:

330

```
331 plink --bfile ~/session4/data/1kg_data.pass --pheno
332 ~/session4/data/1kg_data.BIN.pheno --make-bed --out
333 ~/session4/data/1kg_data.pass
```

- 334
- 335 **1. Basic association test (allelic).** To perform a basic  $\chi^2$  test, which compares frequencies of
- 336 alleles in cases versus controls, type:

```
337
338 cd ~/session4/plink_out
339 plink --bfile ~/session4/data/1kg_data.pass --assoc --out data
```

340

341 This will create an output file 'data.assoc'. It has one row per SNP containing the chromosome

342 [CHR], the SNP identifier [SNP], the base-pair location [BP], the minor allele [A1], the frequency

343 of the minor allele in the affected/cases [F\_A] and unaffected/controls [F\_U], the major allele

344 [A2] and statistical data for an allelic association test including the  $\chi^2$  test statistic [CHISQ], the

345 asymptotic *P*-value [P] and the estimated OR for association between the minor allele and

346 disease [OR].

```
347
348 head data.assoc
```

```
349
350 CHR      SNP      BP      A1      F_A      F_U      A2      CHISQ      P      OR
351 1      rs75333668 762320  T      0.06117 0.05626  C      0.2529      0.615  1.093
352 1      rs148711625 865584  A      0.02039 0.0245  G      0.4488      0.5029 0.8288
353 1      rs146327803 865625  A      0.004078 0.001815 G      0.9918      0.3193 2.252
354 1      rs41285790 865628  A      0.001631 0.002722 G      0.3223      0.5702 0.5986
355 1      rs9988179 865694  T      0.01631 0.01089  C      1.259      0.2618 1.506
356 1      rs116730894 865700  T      0.003263 0.001815 C      0.4732      0.4915 1.8
357 1      rs149677938 874456  A      0.002447 0.001815 G      0.1082      0.7422 1.349
```

358

359 **Note:** this test assumes HWE, and may not work optimally when genotype frequencies deviate

360 from HWE in cases or controls. Use only as a descriptive summary.

- 361
- 362 **2. Association between genotype frequencies and disease status.** When there are no covariates
- 363 to consider, carry out a simple  $\chi^2$  test of association which compares genotype frequencies in
- 364 cases versus controls, by using the --model option:

```
365
366 plink --bfile ~/session4/data/1kg_data.pass --model --out data
```

367

368 This command will perform the test of association under several genetic models:

- 369 • Genotypic (2 df) test
- 370 • Cochran-Armitage trend test (additive model)
- 371 • Allelic test (1df)
- 372 • Dominant gene action (1df) test
- 373 • Recessive gene action (1df) test

374

375

376 This creates the output file 'data.model'. It contains five rows per SNP, one for each of the

377 association tests described in **table 2**. Each row contains the chromosome [CHR], the SNP

378 identifier [SNP], the minor allele [A1], the major allele [A2], the test performed [TEST: GENO

379 (genotypic association); TREND (Cochran-Armitage trend); ALLELIC (allelic association); DOM



380 (dominant model); and REC (recessive model)], the cell frequency counts for cases [AFF] and  
 381 controls [UNAFF], the  $\chi^2$  test statistic [CHISQ], the degrees of freedom for the test [DF] and the  
 382 asymptotic *P* value [*P*].

383  
 384 head data.model

385  
 386

CHR	SNP	A1	A2	TEST	AFF	UNAFF	CHISQ	DF	P
1	rs75333668	T	C	GENO	3/69/541	2/58/491	NA	NA	NA
1	rs75333668	T	C	TREND	75/1151	62/1040	0.2492	1	0.6176
1	rs75333668	T	C	ALLELIC	75/1151	62/1040	0.2529	1	0.615
1	rs75333668	T	C	DOM	72/541	60/491	NA	NA	NA
1	rs75333668	T	C	REC	3/610	2/549	NA	NA	NA
1	rs148711625	A	G	GENO	0/25/588	0/27/524	NA	NA	NA
1	rs148711625	A	G	TREND	25/1201	27/1075	0.4593	1	0.498
1	rs148711625	A	G	ALLELIC	25/1201	27/1075	0.4488	1	0.5029
1	rs148711625	A	G	DOM	25/588	27/524	NA	NA	NA

396  
 397 Note: Genotypic, dominant and recessive tests will not be conducted if any one of the cells in  
 398 the table of case control by genotype counts contains less than five observations. This is  
 399 because the  $\chi^2$  approximation may not be reliable when cell counts are small. To change the  
 400 behavior, use the '--cell' option. For example, to lower the threshold to 3, one would type

401  
 402 plink --bfile ~/session4/data/1kg\_data.pass --model --cell 3 --out data

403  
 404 **3.** Another option for small counts is to use Fisher's exact test. Type

405  
 406 plink --bfile ~/session4/data/1kg\_data.pass --model fisher --out fisher

407  
 408 This will create an output file 'fisher.model'.

409  
 410 head fisher.model

411  
 412

CHR	SNP	A1	A2	TEST	AFF	UNAFF	P
1	rs75333668	T	C	GENO	3/69/541	2/58/491	0.904
1	rs75333668	T	C	TREND	75/1151	62/1040	0.6176
1	rs75333668	T	C	ALLELIC	75/1151	62/1040	0.6595
1	rs75333668	T	C	DOM	72/541	60/491	0.7113
1	rs75333668	T	C	REC	3/610	2/549	1
1	rs148711625	A	G	GENO	0/25/588	0/27/524	0.5703
1	rs148711625	A	G	TREND	25/1201	27/1075	0.498
1	rs148711625	A	G	ALLELIC	25/1201	27/1075	0.5748
1	rs148711625	A	G	DOM	25/588	27/524	0.5703

422  
 423 Warning: still reports Cochran-Armitage test results under allelic test (Chi-square, 1df)

424  
 425 **4.** When there are covariates (usually sex, age, principal components of ancestry), perform  
 426 association tests using logistic regression:

427  
 428 plink --bfile ~/session4/data/1kg\_data.pass --logistic --hide-covar --covar  
 429 ~/session4/data/1kg\_data.covar --out data

430  
 431 By default, this command performs a test of association assuming a multiplicative model. To  
 specify a genotypic, dominant or recessive model in place of a multiplicative model, include

432 the model option `--genotypic`, `--dominant` or `--recessive`, respectively. To include sex as  
433 a covariate, include the option `--sex` (in our case, sex is included in the covariate file, so will  
434 be automatically used).

435  
436 `head data.assoc.logistic`

```
437  
438  
439  
440  
441  
442  
443  
444  
445  
446
```

CHR	SNP	BP	AI	TEST	NMISS	OR	STAT	P
1	rs75333668	762320	T	ADD	1164	0.9238	-0.4272	0.6692
1	rs148711625	865584	A	ADD	1164	0.6893	-1.284	0.1992
1	rs146327803	865625	A	ADD	1164	1.809	0.7036	0.4817
1	rs41285790	865628	A	ADD	1164	0.8182	-0.2181	0.8273
1	rs9988179	865694	T	ADD	1164	1.333	0.7688	0.442
1	rs116730894	865700	T	ADD	1164	1.631	0.5621	0.5741
1	rs149677938	874456	A	ADD	1164	1.124	0.1277	0.8984

447 To output top association results:

```
448  
449 awk '{if($9 <= 1e-4) print $0}' data.assoc.logistic >  
450 data.assoc.logistic.top.txt
```

```
451  
452 head data.assoc.logistic.top.txt
```

```
453  
454  
455  
456  
457
```

2	rs17188434	157096776	C	ADD	1164	3.251	4.002	6.29e-05
7	rs2108225	107453103	G	ADD	1164	0.6469	-4.51	6.475e-06

```
458 #####
```

## 459 8. Data visualization and interpretation

```
460 #####
```

461

462 (a) **Quantile-quantile plots.** To create a quantile-quantile plot of p-values, follow these steps.

463

- 464 i. Start R software (type R at the prompt).
- 465 ii. To create a q-q plot based on the results of chi-square tests (performed in 7.2 above), copy  
466 and paste the following commands at the prompt:

```
467  
468 data <- read.table("data.model", header=TRUE);  
469 obs <- -log10(sort(data[data$TEST == "TREND", ]$P));  
470 exp <- -log10(c(1:length(obs))/(length(obs) + 1));  
471 pdf("pvalue.chisq.qq.plot.pdf");  
472 plot(exp, obs, ylab="Observed(-logP)", xlab="Expected(-logP)", ylim=c(0,  
473 8), xlim=c(0,6));  
474 abline(a=0, b=1, col=1, lwd=1.5, lty=2);  
475 dev.off()
```

476

477 Open the file `"pvalue.chisq.qq.plot.pdf"`. What do you think about this plot?

478

- 479 iii. Now generate a similar plot based on the results of logistic regression analysis.

480

```
481  
482 data <- read.table("data.assoc.logistic", header=TRUE);  
483 obs <- -log10(sort(data[data$TEST == "ADD", ]$P));  
484 exp <- -log10(c(1:length(obs))/(length(obs) + 1));  
pdf("pvalue.logistic.qq.plot.pdf");
```

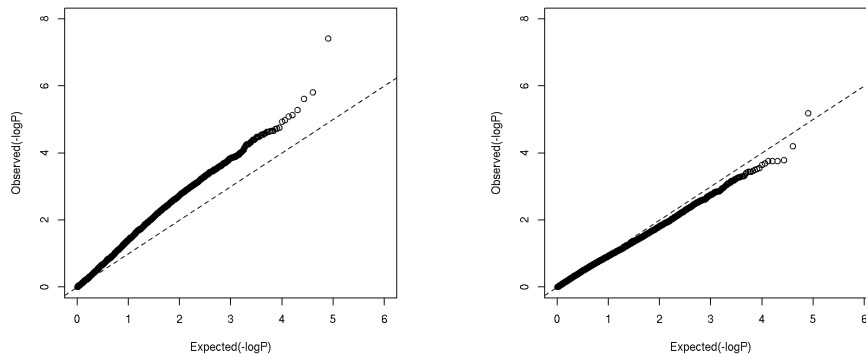
```

485 plot(exp, obs, ylab="Observed(-logP)", xlab="Expected(-logP)", ylim=c(0,
486 8), xlim=c(0,6));
487 abline(a=0, b=1, col=1, lwd=1.5, lty=2);
488 dev.off()
489 q()

```

490  
491 Open the file "pvalue.logistic.qq.plot.pdf". What do you think about this plot?

492  
493 **Figure 1:** Q-Q plots based on association analysis of a binary trait.



495  
496  
497  
498

499 **(b) Calculate the genomic control inflation factor  $\lambda$  for GWA studies.**

- 500 (i) To obtain the inflation factor, include the `--adjust` option in any of the PLINK commands  
501 described in Step 4. For example, the inflation factor based on logistic regression assuming a  
502 multiplicative model is obtained by typing

```

503
504 plink --bfile ~/session4/data/lkg_data.pass --logistic --hide-covar --covar
505 ~/session4/data/lkg_data.covar --adjust --out data

```

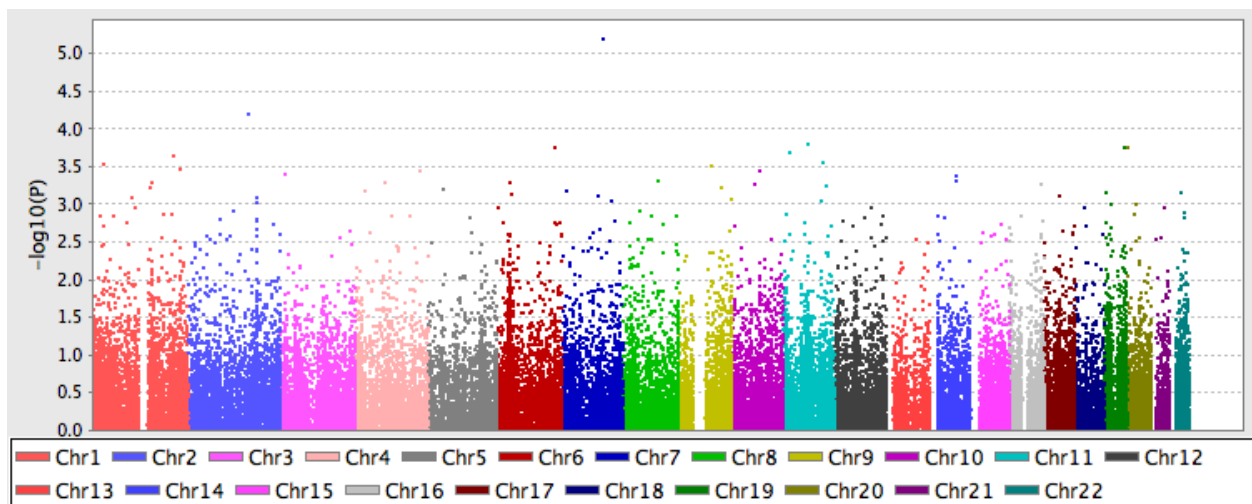
- 506 (ii) Open the PLINK log file 'data.log', which records the inflation factor. The inflation factor for  
507 our GWA study is 1.0077, indicating that no population stratification is detected in our GWA  
508 data (values <1.1 are considered "acceptable")
- 509
- 510 (iii) GC adjustment is based on the median p-value, and does not capture other features of the  
511 distribution (e.g., tail behavior), so can over- or under-correct. Use a diagnostic (to detect if  
512 there is evidence of population stratification) rather than to correct p-values.

513  
514

515 **(c) Manhattan plots.**

- 516
- 517 (i) Start Haploview (`java -jar ~/session4/bin/Haploview.jar`). In the 'Welcome to  
518 Haploview' window, select the 'PLINK Format' tab. Click the 'browse' button and select the  
519 SNP association output file created in Step 7. We select association results from the file

- 520 'data.assoc.logistic'. Select the corresponding MAP file, which will be the '.bim' file for the  
 521 binary file format. We select our GWA study file '1kg\_data.pass.bim'. Leave other options as  
 522 they are (ignore pairwise comparison of markers > 500 kb apart and exclude individuals with  
 523 > 50% missing genotypes). Click 'OK'.  
 524  
 525 (ii) Select the association results relevant to the test of interest by selecting 'TEST' in the  
 526 dropdown tab to the right of 'Filter:', '=' in the dropdown menu to the right of that and the  
 527 PLINK keyword corresponding to the test of interest in the window to the right of that. We  
 528 select PLINK keyword 'ADD' to visualize results for allelic tests of association in our GWA  
 529 study. Click the gray 'Filter' button. Click the gray 'Plot' button. Leave all options as they are  
 530 so that 'Chromosomes' is selected as the 'X-Axis'. Choose 'P' from the drop-down menu for  
 531 the 'Y-Axis' and '-log10' from the corresponding dropdown menu for 'Scale:'. Click 'OK' to  
 532 display the Manhattan plot.  
 533  
 534 (iii) To save the plot as a scalable vector graphics file, click the button 'Export to scalable vector  
 535 graphics:' and then click the 'Browse' button (immediately to the right) to select the  
 536 appropriate title and directory. **Or, after the plot is generated, right click with your mouse  
 537 and choose "Save as..." from the menu, to save the graph as a PNG file.**  
 538



539  
 540  
 541 **Figure 2:** Manhattan plot.  
 542

- 543 (iv) To create a Manhattan plot in R, start R software and copy the following commands:

```
544 source('~/.session4/bin/Manhattan.plot.R')
545 data <-read.table("data.assoc.logistic", header=TRUE);
546 cl <-c("red", "navyblue", "darkgreen", "gold", "deepskyblue4", "magenta4", "slategray")
547 png('Manhattan_plot.png', width = 8.5, height = 3.5, units = "in", res=300)
548 par(mar=c(4.1,4.1,1.6,1.1), cex.lab=1.4, cex.axis=1.3, mgp=c(2.75, .95, 0), las=1,
549 font=2)
550 m.plot(data$P, data$CHR, data$BP, cex=0.75, pch=16, cex.axis=1.3, cex.lab=1.5, col=cl,
551 mgp=c(2.75, .95, 0), pt.cex=0.9, main="", ylab=expression(paste(-log[10], ' P-value')));
552 abline(h=-log10(0.05/sum(!is.na(data$P))), lty=2, col='gray37', lwd=0.75)
```

554 dev.off()

555

556 #####

## 557 9. Quantitative traits

558 #####

559  
560 (a) **Basic quantitative trait association.** To load a quantitative phenotype, use the option `--pheno`.  
561 To obtain a basic association test between genotype and a quantitative trait, type:

562

563

```
563 plink --bfile ~/session4/data/1kg_data.pass --pheno
```

```
564 ~/session4/data/1kg_data.QT.pheno --assoc --out data
```

565 This will generate the file `'data.qassoc'`, with the following columns:

566	CHR	SNP	BP	NMISS	BETA	SE	R2	T	P
567	1	rs75333668	762320	1164	-3.015	1.785	0.00245	-1.689	0.09141
568	1	rs148711625	865584	1164	-2.871	2.899	0.0008431	-0.9902	0.3223
569	1	rs146327803	865625	1164	-2.423	7.75	8.412e-05	-0.3127	0.7546
570	1	rs41285790	865628	1164	-7.605	9.159	0.000593	-0.8303	0.4065
571	1	rs9988179	865694	1164	-2.285	3.664	0.0003345	-0.6236	0.533

572

573

574 (b) As with a binary trait, we typically want to include covariates (such as age, gender and  
575 ancestry). To do that, use linear regression (`--linear`) to test the association.

```
576 plink --bfile ~/session4/data/1kg_data.pass --linear --pheno
```

```
577 ~/session4/data/1kg_data.QT.pheno --pheno-name Pheno --hide-covar --covar
```

```
578 ~/session4/data/1kg_data.covar --out data
```

579 View the file `"data.assoc.linear"`.

580	CHR	SNP	BP	A1	TEST	NMISS	BETA	STAT	P
581	1	rs75333668	762320	T	ADD	1164	-1.031	-0.551	0.5818
582	1	rs148711625	865584	A	ADD	1164	-1.042	-0.3552	0.7225
583	1	rs146327803	865625	A	ADD	1164	-1.361	-0.1759	0.8604
584	1	rs41285790	865628	A	ADD	1164	-9.112	-0.995	0.3199
585	1	rs9988179	865694	T	ADD	1164	-0.5524	-0.1503	0.8805

586

587 (c) Generate a q-q plot of the results in R. Start R software.

```
588 data <- read.table("data.assoc.linear", header=TRUE);
```

```
589 obs <- -log10(sort(data[data$TEST == "ADD", ]$P));
```

```
590 exp <- -log10(c(1:length(obs))/(length(obs) + 1));
```

```
591 pdf("pvalue.linear.qq.plot.pdf");
```

```
592 plot(exp, obs, ylab="Observed(-logP)", xlab="Expected(-logP)", ylim=c(0,
```

```
593 max(obs)), xlim=c(0,6));
```

```
594 abline(a=0, b=1, col=1, lwd=1.5, lty=2);
```

```
595 dev.off()
```

596

597 What do you think?

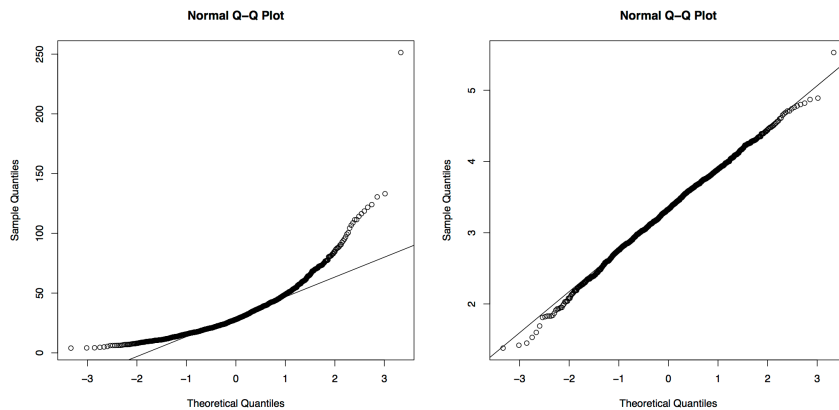
598 (d) What are the assumptions of linear regression analysis? What was the distribution of the  
599 quantitative trait? Generate a normal q-q plot.

```
600 pheno <-read.table('~/.session4/data/1kg_data.QT.pheno', h=T); dim(pheno)
601 pheno[1:2,]
602
603 pdf("Normal.qq.plot.pheno.pdf");
604 qqnorm(pheno$Pheno); qqline(pheno$Pheno)
605 dev.off()
606
607 pdf("Normal.qq.plot.logpheno.pdf");
608 qqnorm(pheno$lPheno); qqline(pheno$lPheno)
609 dev.off()
610 q()
```

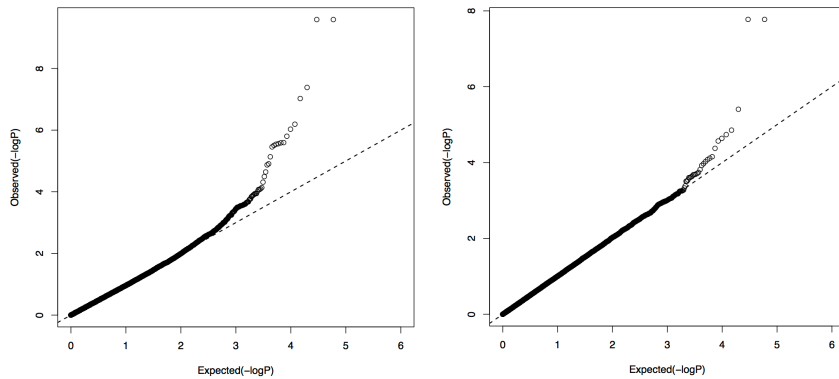
612 (e) Now re-run the association analysis using a log-transformed phenotype. Create a new q-q plot  
613 and compare the results.

```
614 plink --bfile ~/.session4/data/1kg_data.pass --linear --pheno
615 ~/.session4/data/1kg_data.QT.pheno --pheno-name lPheno --hide-covar --covar
616 ~/.session4/data/1kg_data.covar --out data2
```

```
617 R
618 data <- read.table("data2.assoc.linear", header=TRUE);
619 obs <- -log10(sort(data[data$TEST == "ADD", ]$P));
620 exp <- -log10(c(1:length(obs))/(length(obs) + 1));
621 pdf("pvalue.linear.logpheno.qq.plot.pdf");
622 plot(exp, obs, ylab="Observed(-logP)", xlab="Expected(-logP)", ylim=c(0,
623 max(obs)), xlim=c(0,6));
624 abline(a=0, b=1, col=1, lwd=1.5, lty=2);
625 dev.off()
626
```



627



628  
 629 **Figure 3:** Top panels: Normal q-q plot of raw phenotype data (top left) and log-transformed values  
 630 (top right). Lower panels: q-q plots of p-values based on the association analysis of raw phenotype  
 631 data (lower left) and log-transformed values (lower right).

632 (f) Generate a Manhattan plot and create a plot of regional association results for the top hit.

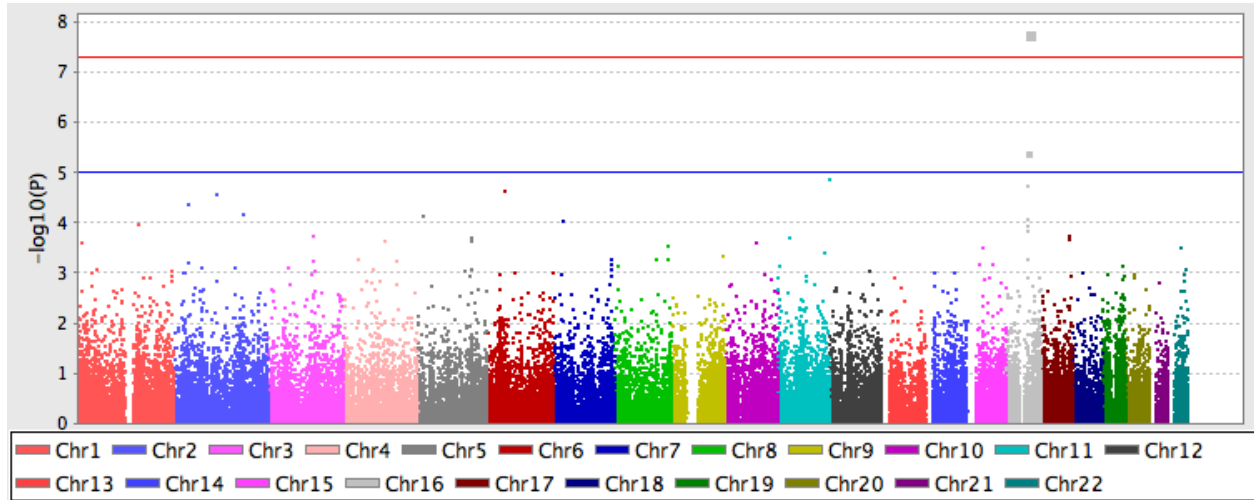
```
633 data[order(data$P), ][1:10,]
```

634	635	636	637	638	639	640	641	642	643	644	645	646	647	648	649	650	651
CHR	SNP	BP	A1	TEST	NMISS	BETA	STAT	P									
64190	16	rs1421085	53800954	C	ADD	1164	0.1754	5.683	1.678e-08								
64191	16	rs1558902	53803574	A	ADD	1164	0.1754	5.683	1.678e-08								
64198	16	rs9941349	53825488	T	ADD	1164	0.1295	4.637	3.931e-06								
71675	19	rs141060900	7691062	G	ADD	1164	-1.0920	-4.601	4.678e-06								
51756	11	rs6590705	133334522	A	ADD	1164	-0.1640	-4.363	1.396e-05								
64193	16	rs17817449	53813367	G	ADD	1164	0.1043	4.302	1.832e-05								

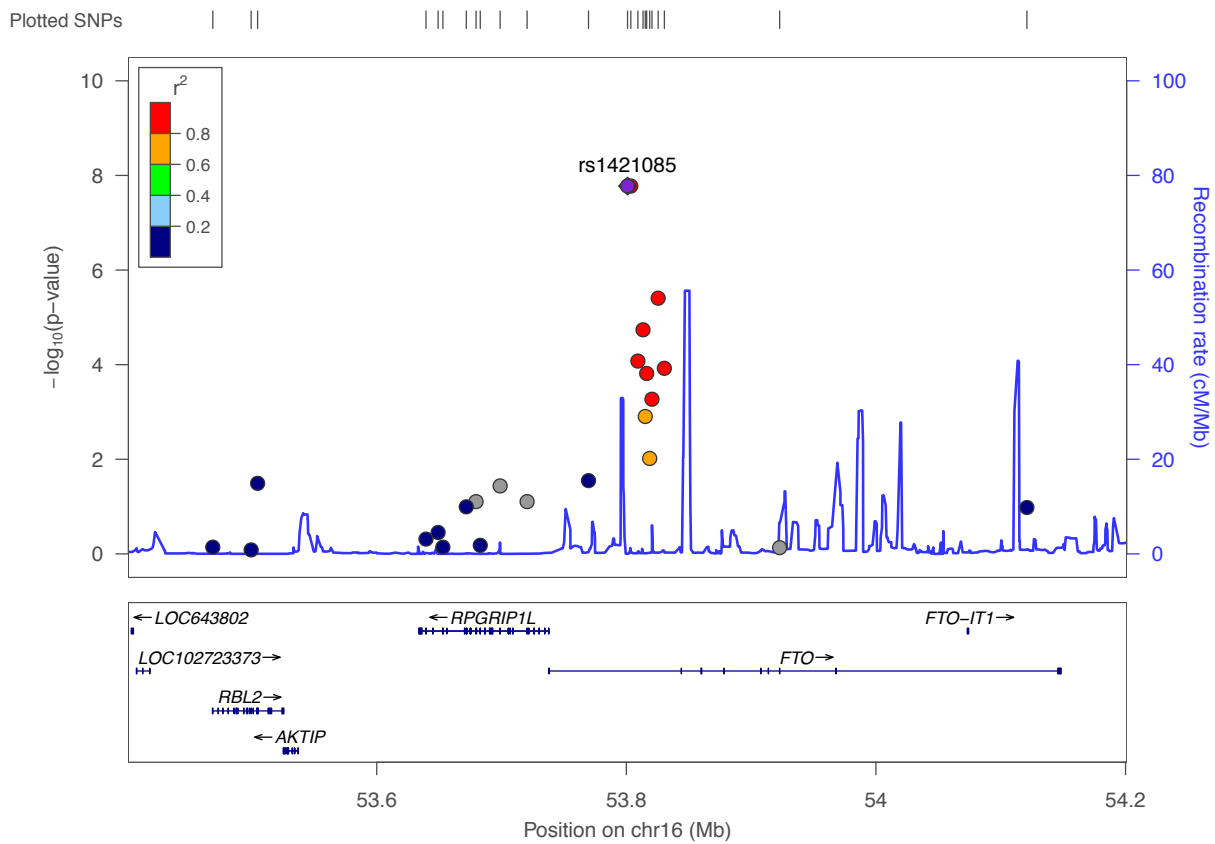
643 Go to Locuszoom website: <http://locuszoom.sph.umich.edu/locuszoom/>

- 644 - Click on "Plot Using your data"
- 645 - Choose file: "data2.assoc.linear"
- 646 - P-Value Column Name: P
- 647 - Marker Column Name: SNP
- 648 - Column Delimiter: WhiteSpace
- 649 - SNP Reference Name: rs1421085
- 650 - Choose Genome Build/LD Population (we will leave EUR)
- 651 - Click on "Plot the data" at the bottom of the page.

652  
 653 **Figure 4:** Manhattan plot of association results for a quantitative trait.



**Figure 5:** Plot of regional association results around the top SNP.



#####  
**10. Gene-based tests**



661 #####

662  
663 To perform gene based tests for rare variants, we will use [EPACTS](#). We will use a built-in example  
664 in EPACTS package, using genotype data from chr 20 on a subset of 1000 genomes project  
665 participants.

```
666  
667 cd ~/session4  
668 run_epacts -shell  
669 ./myrun_epacts.sh  
670 exit
```

671  
672 The script above will perform single-variant association analysis for a binary phenotype DISEASE  
673 and then a burden test for variants with MAF<0.05. The annotated results can be viewed in the  
674 directory `epacts_out`

675  
676  
677  
678

679 **References:**

680 Anderson et al. (2010) Data quality control in genetic case-control association studies. Nature  
681 Protocols, 5(9), 1564.

682 Clarke et al. (2011). Basic statistical analysis in genetic case-control studies. Nature Protocols, 6(2),  
683 121.