

GenSpace-Pgx_day2

July 25, 2020

1 GenSpace Workshop for Personalized Medicine

1.0.1 Date: July 26th, 2020

1.0.2 Objective: How to calculate Polygenic Risk Scores (PRS) to understand SSRI response

1.0.3 Facilitator: Kumar Veerapen, PhD and Caitlin Cooney, CGC

Creation of this hands-on workshop material is thanks to the inspiration from materials provided by Laramie Duncan, PhD, and Hanyang Shen, MS, Stanford University.

Materials from this notebook is complementary from the lecture slides found [here](#)

1.1 Learning objectives

- 1) To learn basic commands in Jupyter notebooks
- 2) To understand the calculation of PRS
- 3) To determine association of PRS to a genetic trait of interest
- 4) To define statistical limitations of PRS
- 5) To figure out what is considered a best fit PRS model

1.2 What is Polygenic Risk Score (PRS)?

An area in genetics that has large traction is understanding polygenicity of genetic traits. In recent years genome-wide association studies (GWAS) have discovered thousands of genetic variants associated both with diseases (depression, heart disease) and complex traits (height). Complex traits are traits derived from multiple genes, and exhibit a large range of attributes.

Polygenic risk scores evaluate these variants to help provide a personalized and more accurate risk/prediction assessment such as depression and height. The term PRS is synonymous with Genomic Risk Scores (GRS), Risk Profile Scores (RPS), or simply genetic scores.

Despite being a trendy method in understanding genetic traits, these scores are far from perfect. The PRS that have been published thus far usually tend to explain less than [20% of a trait variance in complex traits](#).

There are strengths and weaknesses to utilize PRS which are covered in the lecture session of this hands-on workshop.

1.3 What are we doing?

This hand-on workshop will allow you to explore a conventional approach to computing PRS using a ubiquitous genetics tool called [PLINK](#).

Problem statement: Scientifically, we know that approx [53-65% of patients with depression respond well to antidepressants](#).

Hypothesis: We hypothesize that this variability is potentially caused by genetic risk.

Scientific question/objective: Are individuals with higher depression PRS more responsive to selective serotonin reuptake inhibitors (SSRI) [a type of antidepressant]?

1.4 Step by step overview

1.4.1 First...

We will obtain a reference data that contains effect sizes needed for calculating PRS.

Effect size : a statistical concept that measures the strength of the relationship between two variables on a numeric scale.

Example : If we look at a population's height, on average men are taller than women. This difference between the height of men and the height of women is known as the effect size. If the effect size is high, there is a higher height difference between men and women. If the effect size is small, there is less height difference.

The effect sizes (or summary statistics) that we will be using were generated from a reference data published by [David Howard and co in Nature Neuroscience \(2019\)](#). Their study attempted to understand the genetic risk of major depressive disorder in a very large collection of depressed patients. The reference data from this paper will contain effect sizes that we will need to calculate PRS in our test data which is a modified version of the 1000 Genomes Project.

[The 1000 Genomes Project ran between 2008 and 2015, creating the largest public catalogue of human variation and genotype data](#)

In this project, variation means to look for polymorphisms (genetic variants with frequencies of at least 1% in the populations studied)

1.4.2 Second...

We will calculate the PRS to answer our scientific question (refer to previous cell).

Using the effect sizes from our reference data (Howard et al 2019), we will then calculate PRS from samples obtained from a modified version of the 1000 Genomes data on a binary trait called selective serotonin reuptake inhibitor (SSRI) response. An SSRI is a typical group of antidepressants prescribed to patients with depressive symptoms.

The trait that you are analysing was simulated where

1 = the individual did not show an improvement in their depressive symptoms (measured by depression scores)

2 = the individual showed an improvement in their depressive symptoms (measured by depression scores)

1.4.3 Third...

We will use the PRS calculation to find an association to the trait that we are testing – SSRI response.

Finally, we will then use a statistical programming tool called R in order to 1) analyse the PRS generated; and 2) visualize our PRS using plots.

In order to provide a little understanding of the code that we are running, we have included some details in the comment section of your workshop.

Let's start with finding out how to use Jupyter notebook in order to run your analysis.

1.5 Using Jupyter

1.5.1 Running cells

Evaluate cells using SHIFT + ENTER. Select the next cell and run it

```
[1]: print('Hello, world')
```

```
[1] "Hello, world"
```

1.5.2 Modes

Jupyter has two modes, a **navigation mode** and an **editor mode**.

Navigation mode:

- BLUE cell borders
- UP / DOWN move between cells
- ENTER while a cell is selected will move to **editing mode**.
- Many letters are keyboard shortcuts! This is a common trap.

Editor mode:

- GREEN cell borders
- UP / DOWN/ move within cells before moving between cells.
- ESC will return to **navigation mode**.
- SHIFT + ENTER will evaluate a cell and return to **navigation mode**.

1.5.3 Cell types

There are several types of cells in Jupyter notebooks. The two you will see here are **Markdown** (text) and **Code**.

```
[2]: # This is a code cell  
my_variable = 5
```

This is a markdown cell, so even if something looks like code (as below), it won't get executed!

```
my_variable += 1
```

1.5.4 Tips and tricks

Keyboard shortcuts:

- SHIFT + ENTER to evaluate a cell
- ESC to return to navigation mode
- y to turn a markdown cell into code
- m to turn a code cell into markdown
- a to add a new cell **above** the currently selected cell
- b to add a new cell **below** the currently selected cell
- d, d (repeated) to delete the currently selected cell
- TAB to activate code completion

To try this out, create a new cell below this one using **b**, and print `my_variable` by starting with `print(my` and pressing TAB!

1.6 Set up R environment with packages that you will need

[R is a statistical programming tool](#) (that's free) is a tool that many analysts use in order to analyse their datasets. This is especially true in genetics where we are cheap with buying software and analytic tools. With R, you can set up your environment to get additional tools to run some of your analyses.

The common practice in R is to load all of your needed tools (called packages) at the beginning of your script, i.e. the next line.

`fmsb` will be used in order to calculate the variance explained by the PRS that you will be computing. The first line here, downloads the package into your environment.

Let's load the package into your environment now

```
[3]: library(fmsb)
```

When you want to look up a function in R, you can use a `?` sign

```
[ ]:
```

2 Let's get started!

2.1 Step 0: Exploration of your files and directories

Where are you?

```
[4]: getwd()
```

`~/Users/kumar/Dropbox (Partners HealthCare)/Teaching_GenSpace_Pharmacogenetics/genspace_prs/resources`

What files are in your directory (a.k.a. folder)?

```
[5]: list.files()
```

1. 'GenSpace-Pgx_day2.ipynb' 2. 'HELP' 3. 'MDD_2019_logOR_score'
 4. 'MDD_2019_logORpVal' 5. 'MDD_2019_pvalue_score' 6. 'OUTCOME' 7. 'Pherandom.reduced_1000_Genomes.bed' 8. 'Pherandom.reduced_1000_Genomes.bim' 9. 'Pherandom.reduced_1000_Genomes.fam' 10. 'plink' 11. 'plink.log' 12. 'q.ranges.GWASsig_to_1' 13. 'QC' 14. 'Ranalysis_PRS.R'

These are the files that we will be using for the workshop!

What files are what?

2.1.1 PLINK files

For [PLINK](#), we have 3 files that will be used which are `.bed`, `.bim`, and `.fam` files.

These are known as binary PLINK files where

.bed file File name: `Pherandom.reduced_1000_Genomes.bed`

contains the genetic information for each individual in 0 and 1s. This file would contain individual level data such as deidentified sample ID, family ID (if applicable), and genotypes at multiple positions in the genome.

This file cannot be viewed because it is a binary file.

.bim file File name: `Pherandom.reduced_1000_Genomes.bim`

If the bed file contains the individual ID, the `.bim` file contain that generalized information. Think of the `.bim` files as Google map, and the `.bed` file as a person on the map.

This file contains locations of the genetic information that was captured such as which chromosome the variants are on. It has 6 columns and the most important are the 1st, 2nd, 4th, 5th, and 6th:

- (1) Chromosome number
- (2) SNP name / ID (oftentimes and 'rs' number)
- (3) Chromosome Position
- (4) Reference / Wildtype Allele
- (5) Variant Allele

Let's view it with

```
[6]: bim.file <- read.table("Pherandom.reduced_1000_Genomes.bim", header=F, sep="\t")
     head(bim.file)
```

V1	V2	V3	V4	V5	V6
1	rs11579015	0	1036959	C	T
1	rs12029885	0	1125348	C	T
1	rs113908945	0	1183858	T	G
1	rs7549601	0	1213224	A	T
1	rs3766176	0	1486903	C	T
1	rs4648729	0	1808769	C	T

.fam file File name: `Pherandom.reduced_1000_Genomes.fam`

contains the individual ID and family name (for this workshop, ALL individuals have been deidentified).

The 1st, 2nd, and 6th columns are the important ones for this workshop.

- (1) Family ID = FID
- (2) Individual ID = IID
- (3) Trait identification.

Note that 1 means control and 2 means case, per PLINK conventions. The trait values here are randomly generated where :

1 would indicate that the individual did not show an improvement in their depressive symptoms.

2 would indicate that the individual showed an improvement in their depressive symptoms

Let's view this file

```
[7]: fam.file <- read.table("Pherandom.reduced_1000_Genomes.fam", header=F, sep="")
     head(fam.file)
```

V1	V2	V3	V4	V5	V6
HG00096	HG00096	0	0	0	2
HG00097	HG00097	0	0	0	2
HG00099	HG00099	0	0	0	2
HG00100	HG00100	0	0	0	2
HG00101	HG00101	0	0	0	1
HG00102	HG00102	0	0	0	2

2.1.2 Files used for scoring

The reference data that we are using for computing PRS are also here in this environment. A reminder the reference data used in this workshop is from [David Howard and co on genetic associations in major depressive disorder published in 2019](#).

There are 3 files that we will need which are the

- 1) "Score File",
- 2) "p-value file", and
- 3) the "TRANCHE file"

The "Score file" File name: `MDD_2019_logOR_score`.

This file contains the effect sizes from David Howard's paper needed for our PRS calculation. There are 3 columns in this file which are SNP ID, Risk allele, and logOR (the effect size that we keep talking about).

Let's view it

```
[8]: score.file <- read.table("MDD_2019_logOR_score", header=T, sep="")
head(score.file)
```

SNP	allele	OR
rs11579015	T	0.02620370
rs12029885	T	0.00970278
rs113908945	T	0.03649590
rs7549601	A	-0.00879859
rs3766176	T	-0.02839950
rs4648729	T	-0.00810274

The “p-value file” File name : MDD_2019_pvalue_score.

This file contains the p-values from David Howard’s paper that were used to tell if an effect size was considered significant or not.

As you are aware, we will need effect sizes in order to calculate PRS. These effect sizes were considered significant if they have a p-value that are less than 5×10^{-8} (0.000000005) from David Howard’s paper. This value is a typical value used in genome-wide association studies (GWAS). That p-value would imply that the variant identified was 5×10^{-6} (0.0000005)% a false positive. The usage of p-values is ubiquitous in all of statistics so that analysts do not go down rabbit holes of signals which were bad apples to begin with.

This file contains 2 columns: SNP ID and p-value

Let’s view the file (Notice that there are p-values which are traditionally, insignificant)

```
[9]: pval.file <- read.table("MDD_2019_pvalue_score", header=T, sep="")
head(pval.file)
```

SNP	pval
rs11579015	0.12570
rs12029885	0.56880
rs113908945	0.03714
rs7549601	0.63200
rs3766176	0.02987
rs4648729	0.45270

The “TRANCHE file” File name: q.ranges.GWASsig_to_1.

This file contains tranches / bin / groups of p-value thresholds that will be used to group effect sizes for PRS calculations. The most strict of all p-values in the tranche file is the GWAS significant tranche (PSi). This would set our program (PLINK) to calculate PRS using effect sizes that were all having a p-value that are lower than 5×10^{-8} or 0.000000005 (PSi). Tranches listed in the file range from the strictest p-value threshold to the most lenient which includes all the variants that were tested in David Howard’s paper (p-value threshold of lower than 1.0 (all)).

Let’s view the file:

```
[10]: tranche.file <- read.table("q.ranges.GWASsig_to_1", header=F, sep="")
head(tranche.file)
```

V1	V2	V3
Psi	0	5e-08
Pe6	0	1e-06
pe4	0	1e-04
pe3	0	1e-03
pe2	0	1e-02
P05	0	5e-02

Note : In this analysis, we will be using **pe2** (0.001) for our downstream applications.

2.2 Step 1 : Calculating Principal Components

(Timing Check) *Note* : This step will be done if we have time in the workshop. Also this step has been done for you and saved in your directory, just in case we are short on time.

Before we can start our PRS computations, we would need to capture the variance contributed by the different ancestries in our 1000 Genomes data. With multiple different ancestries, this may confound (a.k.a. mess up) our downstream analysis because humans are different from each other. By calculating principal components, we can capture some of the intrinsic differences observed in the 1000 Genome samples.

NOTE : there are a few other steps that we would need to do to clean the data up that I have done for you which include pruning the data for variants which are informative such as common variants etc.

The tool that we can use to compute principal components is PLINK version 1.90.

To run this:

```
[11]: system("./plink --bfile QC/Pherandom.reduced_1000_Genome.QC --pca 20 --out_
      ↳OUTCOME/PCA_20_for_1000_Genomes", intern=T)
```

```
1. 'PLINK v1.90b6.2 64-bit (12 Jun 2018) www.cog-genomics.org/plink/1.9/' 2. '(C)
2005-2018 Shaun Purcell, Christopher Chang GNU General Public License v3'
3. 'Logging to OUTCOME/PCA_20_for_1000_Genomes.log' 4. 'Options in ef-
fect:' 5. ' -bfile QC/Pherandom.reduced_1000_Genome.QC' 6. ' -out OUT-
COME/PCA_20_for_1000_Genomes' 7. ' -pca 20' 8. " 9. '16384 MB RAM detected; re-
serving 8192 MB for main workspace.' 10. '88086 variants loaded from .bim file.' 11. '2504
people (0 males, 0 females, 2504 ambiguous) loaded from .fam.' 12. 'Ambiguous sex IDs
written to OUTCOME/PCA_20_for_1000_Genomes.nosex .' 13. '2504 phenotype values
loaded from .fam.' 14. 'Using up to 4 threads (change this with -threads).' 15. 'Before main
variant filters, 2504 founders and 0 nonfounders present.' 16. 'Calculating allele frequencies...
0%\b\b1%\b\b2%\b\b3%\b\b4%\b\b5%\b\b6%\b\b7%\b\b8%\b\b9%\b\b10%\b\b\b11%\b\b\b12%\b\b\b
done.' 17. '88086 variants and 2504 people pass filters and QC.' 18. 'Note: No phenotypes present.'
19. '\r60 markers complete.\r120 markers complete.\r180 markers complete.\r240 markers
complete.\r300 markers complete.\r360 markers complete.\r420 markers complete.\r480 markers
complete.\r540 markers complete.\r600 markers complete.\r660 markers complete.\r720 markers
complete.\r780 markers complete.\r840 markers complete.\r900 markers complete.\r960 markers
complete.\r1020 markers complete.\r1080 markers complete.\r1140 markers complete.\r1200
markers complete.\r1260 markers complete.\r1320 markers complete.\r1380 markers com-
plete.\r1440 markers complete.\r1500 markers complete.\r1560 markers complete.\r1620 markers
```


complete.\r1680 markers complete.\r1740 markers complete.\r1800 markers complete.\r1860
markers complete.\r1920 markers complete.\r1980 markers complete.\r2040 markers com-
plete.\r2100 markers complete.\r2160 markers complete.\r2220 markers complete.\r2280 markers
complete.\r2340 markers complete.\r2400 markers complete.\r2460 markers complete.\r2520
markers complete.\r2580 markers complete.\r2640 markers complete.\r2700 markers com-
plete.\r2760 markers complete.\r2820 markers complete.\r2880 markers complete.\r2940 markers
complete.\r3000 markers complete.\r3060 markers complete.\r3120 markers complete.\r3180
markers complete.\r3240 markers complete.\r3300 markers complete.\r3360 markers com-
plete.\r3420 markers complete.\r3480 markers complete.\r3540 markers complete.\r3600 markers
complete.\r3660 markers complete.\r3720 markers complete.\r3780 markers complete.\r3840
markers complete.\r3900 markers complete.\r3960 markers complete.\r4020 markers com-
plete.\r4080 markers complete.\r4140 markers complete.\r4200 markers complete.\r4260 markers
complete.\r4320 markers complete.\r4380 markers complete.\r4440 markers complete.\r4500
markers complete.\r4560 markers complete.\r4620 markers complete.\r4680 markers com-
plete.\r4740 markers complete.\r4800 markers complete.\r4860 markers complete.\r4920 markers
complete.\r4980 markers complete.\r5040 markers complete.\r5100 markers complete.\r5160
markers complete.\r5220 markers complete.\r5280 markers complete.\r5340 markers com-
plete.\r5400 markers complete.\r5460 markers complete.\r5520 markers complete.\r5580 markers
complete.\r5640 markers complete.\r5700 markers complete.\r5760 markers complete.\r5820
markers complete.\r5880 markers complete.\r5940 markers complete.\r6000 markers com-
plete.\r6060 markers complete.\r6120 markers complete.\r6180 markers complete.\r6240 markers
complete.\r6300 markers complete.\r6360 markers complete.\r6420 markers complete.\r6480
markers complete.\r6540 markers complete.\r6600 markers complete.\r6660 markers com-
plete.\r6720 markers complete.\r6780 markers complete.\r6840 markers complete.\r6900 markers
complete.\r6960 markers complete.\r7020 markers complete.\r7080 markers complete.\r7140
markers complete.\r7200 markers complete.\r7260 markers complete.\r7320 markers com-
plete.\r7380 markers complete.\r7440 markers complete.\r7500 markers complete.\r7560 markers
complete.\r7620 markers complete.\r7680 markers complete.\r7740 markers complete.\r7800
markers complete.\r7860 markers complete.\r7920 markers complete.\r7980 markers com-
plete.\r8040 markers complete.\r8100 markers complete.\r8160 markers complete.\r8220 markers
complete.\r8280 markers complete.\r8340 markers complete.\r8400 markers complete.\r8460
markers complete.\r8520 markers complete.\r8580 markers complete.\r8640 markers com-
plete.\r8700 markers complete.\r8760 markers complete.\r8820 markers complete.\r8880 markers
complete.\r8940 markers complete.\r9000 markers complete.\r9060 markers complete.\r9120 mark-
ers complete.\r9180 markers complete.\r9240 markers complete.\r9300 markers complete.\r9360
markers complete.\r9420 markers complete.\r9480 markers complete.\r9540 markers com-
plete.\r9600 markers complete.\r9660 markers complete.\r9720 markers complete.\r9780 markers
complete.\r9840 markers complete.\r9900 markers complete.\r9960 markers complete.\r10020
markers complete.\r10080 markers complete.\r10140 markers complete.\r10200 markers com-
plete.\r10260 markers complete.\r10320 markers complete.\r10380 markers complete.\r10440
markers complete.\r10500 markers complete.\r10560 markers complete.\r10620 markers com-
plete.\r10680 markers complete.\r10740 markers complete.\r10800 markers complete.\r10860
markers complete.\r10920 markers complete.\r10980 markers complete.\r11040 markers com-
plete.\r11100 markers complete.\r11160 markers complete.\r11220 markers complete.\r11280
markers complete.\r11340 markers complete.\r11400 markers complete.\r11460 markers com-
plete.\r11520 markers complete.\r11580 markers complete.\r11640 markers complete.\r11700
markers complete.\r11760 markers complete.\r11820 markers complete.\r11880 markers com-
plete.\r11940 markers complete.\r12000 markers complete.\r12060 markers complete.\r12120

markers complete.\r82740 markers complete.\r82800 markers complete.\r82860 markers complete.\r82920 markers complete.\r82980 markers complete.\r83040 markers complete.\r83100 markers complete.\r83160 markers complete.\r83220 markers complete.\r83280 markers complete.\r83340 markers complete.\r83400 markers complete.\r83460 markers complete.\r83520 markers complete.\r83580 markers complete.\r83640 markers complete.\r83700 markers complete.\r83760 markers complete.\r83820 markers complete.\r83880 markers complete.\r83940 markers complete.\r84000 markers complete.\r84060 markers complete.\r84120 markers complete.\r84180 markers complete.\r84240 markers complete.\r84300 markers complete.\r84360 markers complete.\r84420 markers complete.\r84480 markers complete.\r84540 markers complete.\r84600 markers complete.\r84660 markers complete.\r84720 markers complete.\r84780 markers complete.\r84840 markers complete.\r84900 markers complete.\r84960 markers complete.\r85020 markers complete.\r85080 markers complete.\r85140 markers complete.\r85200 markers complete.\r85260 markers complete.\r85320 markers complete.\r85380 markers complete.\r85440 markers complete.\r85500 markers complete.\r85560 markers complete.\r85620 markers complete.\r85680 markers complete.\r85740 markers complete.\r85800 markers complete.\r85860 markers complete.\r85920 markers complete.\r85980 markers complete.\r86040 markers complete.\r86100 markers complete.\r86160 markers complete.\r86220 markers complete.\r86280 markers complete.\r86340 markers complete.\r86400 markers complete.\r86460 markers complete.\r86520 markers complete.\r86580 markers complete.\r86640 markers complete.\r86700 markers complete.\r86760 markers complete.\r86820 markers complete.\r86880 markers complete.\r86940 markers complete.\r87000 markers complete.\r87060 markers complete.\r87120 markers complete.\r87180 markers complete.\r87240 markers complete.\r87300 markers complete.\r87360 markers complete.\r87420 markers complete.\r87480 markers complete.\r87540 markers complete.\r87600 markers complete.\r87660 markers complete.\r87720 markers complete.\r87780 markers complete.\r87840 markers complete.\r87900 markers complete.\r87960 markers complete.\r88020 markers complete.\r88080 markers complete.\r88086 markers complete.\rRelationship matrix calculation complete.' 20. '[extracting eigenvalues and eigenvectors]\r-pca: Results saved to OUTCOME/PCA_20_for_1000_Genomes.eigenval and' 21. 'OUTCOME/PCA_20_for_1000_Genomes.eigenvec '

--bfile is to call in the PLINK binary files that we will need. This contains the 1000 Genomes samples

--pca generate principal components and state number of components wanted, i.e. 20

--out is to print out the principal components

Where is our output PCs?

```
[12]: list.files(path="OUTCOME", pattern="PCA_20_for_1000_Genomes")
```

1. 'PCA_20_for_1000_Genomes.eigenval'
2. 'PCA_20_for_1000_Genomes.eigenvec'
3. 'PCA_20_for_1000_Genomes.log'
4. 'PCA_20_for_1000_Genomes.nosex'

2.3 Step 2: Calculating Polygenic Risk Scores

We have understood what files contain what pieces of information/data that we would need for our final analysis.

We now get to compute PRS for our data. In order to do this, we can run the PLINK command as below, where description of each flag follows the code block

```
[13]: system("./plink --bfile Pherandom.reduced_1000_Genome --score_
↳MDD_2019_logORpVal 1 2 3 header no-mean-imputation --q-score-range q.ranges.
↳GWASsig_to_1 MDD_2019_pvalue_score --extract QC/Pherandom.
↳reduced_1000_Genome.QC.clumped.valid.snp --allow-no-sex --out OUTCOME/MDD",
↳intern=T)
```

Warning message in system("./plink --bfile Pherandom.reduced_1000_Genome --score MDD_2019_logORpVal 1 2 3 header no-mean-imputation --q-score-range q.ranges.GWASsig_to_1 MDD_2019_pvalue_score --extract QC/Pherandom.reduced_1000_Genome.QC.clumped.valid.snp --allow-no-sex --out OUTCOME/MDD", :

"running command './plink --bfile Pherandom.reduced_1000_Genome --score MDD_2019_logORpVal 1 2 3 header no-mean-imputation --q-score-range q.ranges.GWASsig_to_1 MDD_2019_pvalue_score --extract QC/Pherandom.reduced_1000_Genome.QC.clumped.valid.snp --allow-no-sex --out OUTCOME/MDD' had status 2"

1. 'PLINK v1.90b6.2 64-bit (12 Jun 2018) www.cog-genomics.org/plink/1.9/'
2. '(C) 2005-2018 Shaun Purcell, Christopher Chang GNU General Public License v3'
3. 'Logging to OUTCOME/MDD.log.'
4. 'Options in effect:'
5. ' -allow-no-sex'
6. ' -bfile Pherandom.reduced_1000_Genome'
7. ' -extract QC/Pherandom.reduced_1000_Genome.QC.clumped.valid.snp'
8. ' -out OUTCOME/MDD'
9. ' -q-score-range q.ranges.GWASsig_to_1 MDD_2019_pvalue_score'
10. ' -score MDD_2019_logORpVal 1 2 3 header no-mean-imputation'
11. "
12. '16384 MB RAM detected; reserving 8192 MB for main workspace.'

--bfile is to call in the PLINK binary files that we will need. This contains the 1000 Genomes samples

--score is to call in the effect size "Score File" from David Howard's 2019 MDD paper.

We also parse in a few arguments with this flag which are

1 2 3 : where 1 is for the 1st column containing SNP ID, 2nd column is for the effective allele, and 3 is for the effect size estimate

header : is that the file has a header line

no-mean-imputation: we do not want our results for an individual's PRS to depend on the other individuals in the file

--q-score-range is to call in the p-value tranches/bins/groups "TRANCHE file" for grouping the p-values and p-values "p-value file" that we parsed in from David Howard's 2019 MDD paper

--extract is to compensate for genetic structure observed in humans called linkage disequilibrium

--allow-no-sex allows for plink to parse the phenotypes of data where sex is missing

--out is to print out the PRS computation

Where is our output?

```
[14]: list.files(path="OUTCOME", pattern="MDD.P")
```

1. 'MDD.P05.profile' 2. 'MDD.P10.profile' 3. 'MDD.P20.profile' 4. 'MDD.P30.profile'
5. 'MDD.P40.profile' 6. 'MDD.P50.profile' 7. 'MDD.P75.profile' 8. 'MDD.Pe6.profile'
9. 'MDD.Psi.profile'

```
[15]: list.files(path="OUTCOME", pattern="MDD.p")
```

1. 'MDD.pe2.profile' 2. 'MDD.pe3.profile' 3. 'MDD.pe4.profile'

2.4 Step 3: Analysis of MDD PRS to SSRI Response status

First, let's read the polygenic scores (pe2 or p-value tranche of 0.001, i.e. anything lower than 0.001 is included) that we have previously computed and save it into an object named `dataset`

```
[16]: dataset <- read.table("OUTCOME/MDD.pe2.profile",header=T)
```

```
[17]: head(dataset)
```

FID	IID	PHENO	CNT	CNT2	SCORE
HG00096	HG00096	2	20824	9951	-0.000898552
HG00097	HG00097	2	20824	10033	-0.000864247
HG00099	HG00099	2	20824	9958	-0.000560535
HG00100	HG00100	2	20824	9826	-0.000698385
HG00101	HG00101	1	20824	10096	-0.000790333
HG00102	HG00102	2	20824	9977	-0.000402355

What's the range of the polygenic risk scores?

```
[18]: summary(dataset$SCORE)
```

Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
-1.111e-03	-3.326e-04	-8.591e-05	-5.469e-05	2.383e-04	8.793e-04

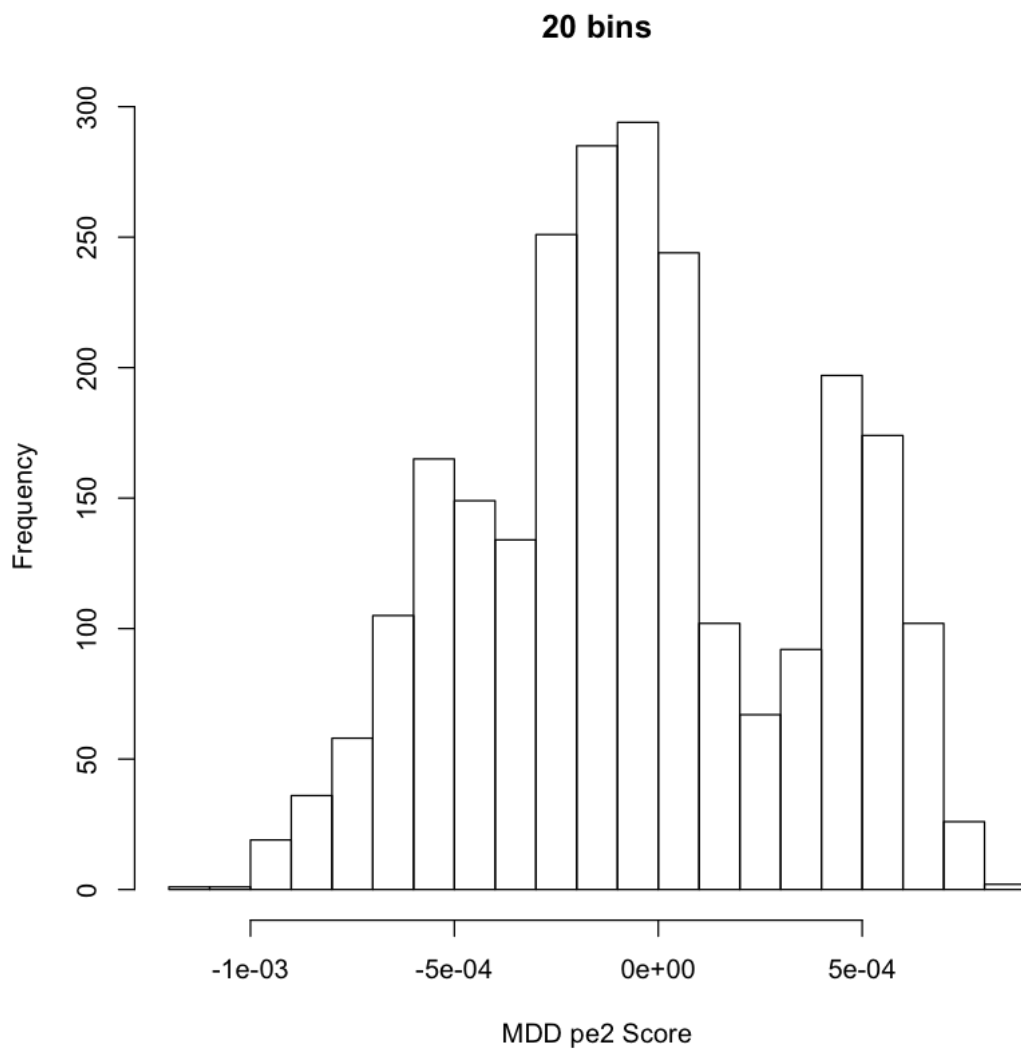
2.5 Step 3.1: Data exploration

Now, we will plot histograms of the PRS, and we will specify how many bins should be included (20 and 100).

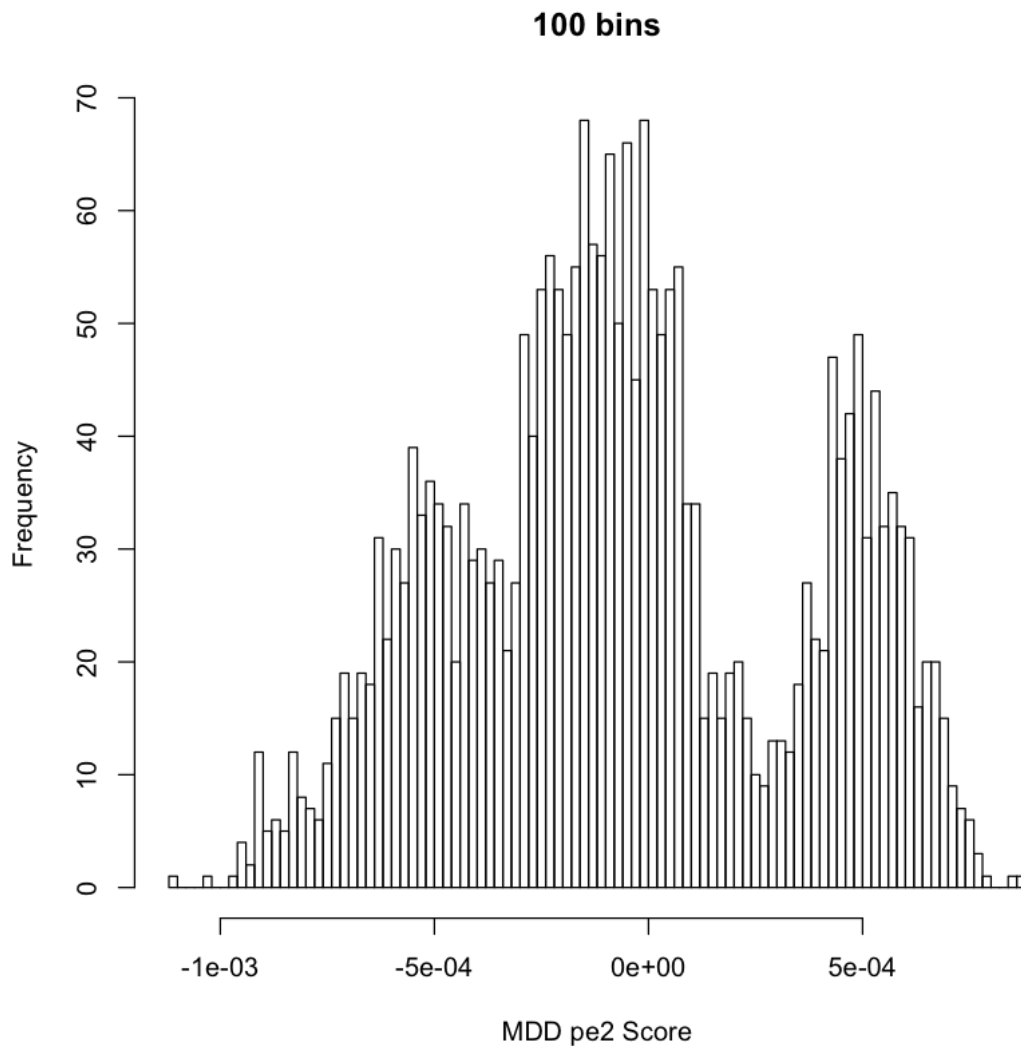
The additional commands below specify the title of the plot ('main=') and the label for the x-axis ('xlab='). It's always a good idea to plot your data throughout the analysis process. Histograms are one of my favourite methods.

using `R:base` to plot out how the histograms look like

```
[19]: hist(dataset$SCORE, main="20 bins",xlab="MDD pe2 Score",breaks=20)
```



```
[20]: hist(dataset$SCORE,main="100 bins",xlab="MDD pe2 Score",breaks=100)
```



Recall that we can use density plots to compare case/control status for polygenic risk load.

```
[21]: ## calculating the density function
noResponse <- density(dataset[dataset$PHENO==1,]$SCORE)
Response <- density(dataset[dataset$PHENO==2,]$SCORE)
```

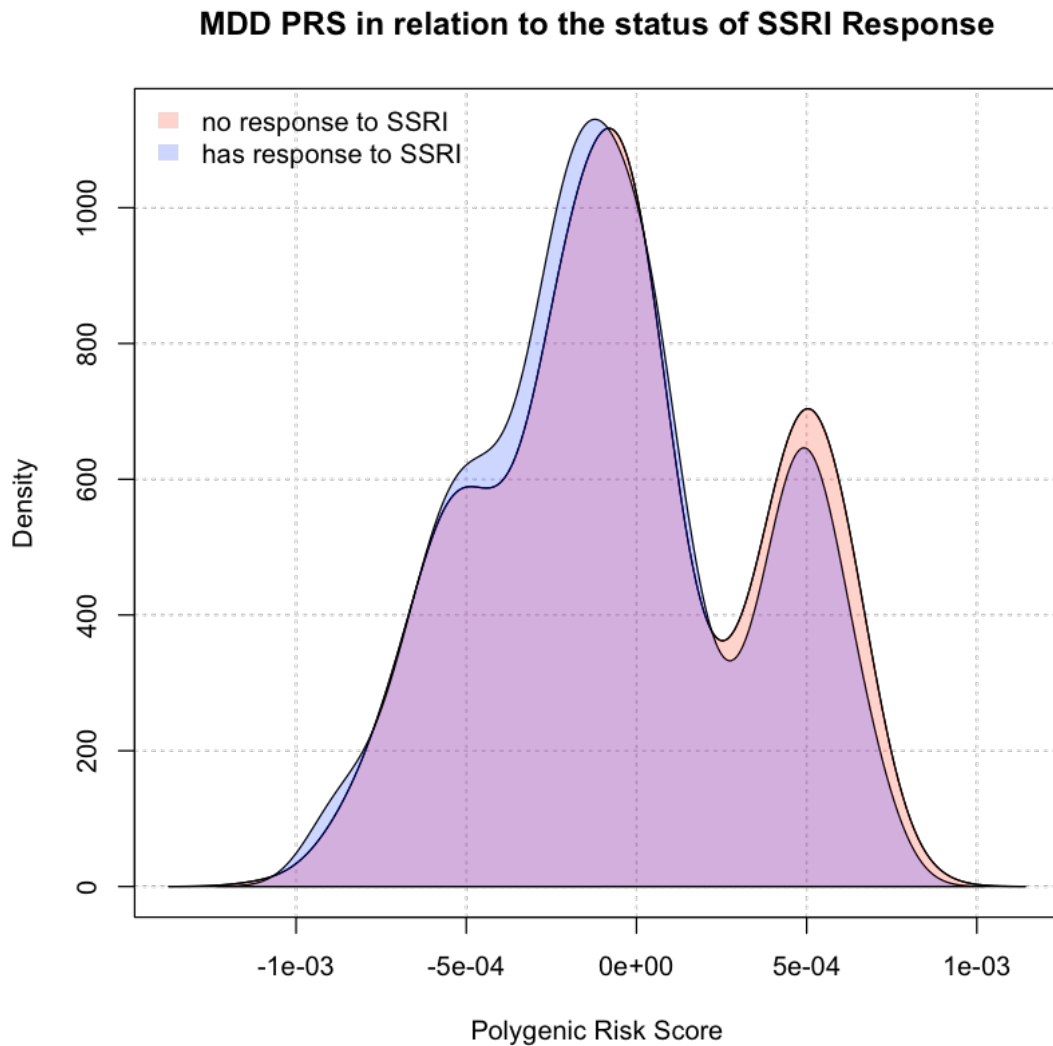
```
[22]: ## calculate the range for the plot
xlim <- range(noResponse$x, Response$x)
ylim <- range(0, noResponse$y, Response$y)
```

```
[23]: #pick the colours
noResponseCol <- rgb(1,0,0,0.2)
ResponseCol <- rgb(0,0,1,0.2)
```

```
[24]: ## plot the no response plot and and set up most of the plot parameters
plot(noResponse, xlim = xlim, ylim = ylim, xlab = 'Polygenic Risk Score',
     main = 'MDD PRS in relation to the status of SSRI Response',
     panel.first = grid())

#put our density plots in
polygon(noResponse, density = -1, col = noResponseCol)
polygon(Response, density = -1, col = ResponseCol)

## add a legend in the corner
legend('topleft', c('no response to SSRI', 'has response to SSRI'),
      fill = c(noResponseCol, ResponseCol), bty = 'n',
      border = NA)
```



2.5.1 What do you think?

Is there a difference/shift in one of the groups? If yes, do you think it would imply a potential polygenic risk load?

2.5.2 Analysis of PRS? Yes/No?

Now that we see some slight difference, is there a statistical way to analyse this?

Why, yes! With something known as a regression analysis. A regression analysis allows us to test for whether a trait is associated with the PRS that we computed.

We can then report a p-value to see if the PRS is statistically significantly associated with the trait. In this case, “**is MDD PRS associated with SSRI response;”.

2.6 Step 3.2 Capturing ancestry differences in human

(Timing Check) Humans are different from each other! We need to control for this so that we can improve accuracy of modeling the trait that we plan on modeling.

To capture this, we use something known as principal components that we have already precomputed (or we did it in our workshop today).

```
[25]: # Read the PC file into a file called 'pca'
pca <- read.table("OUTCOME/PCA_20_for_1000_Genomes.eigenvec")

# Provide the column names with this command.
names(pca) <- c("FID", "IID", paste0("PC", 1:20))

# Merge the polygenic score file (dataset) and the pca file (pca). Merging will
  ↳ be based on the "FID" and "IID" columns.
dataset <- merge(dataset, pca, by=c("FID", "IID"))
```

2.7 Step 3.3 PRS Analysis

Now that we have loaded the principal components, we can finally model our data (to find the association) for “**is MDD PRS associated with SSRI response;”.

We will be using something called a logistic regression analysis. Logistic because the trait that we are testing is a yes/no (binary/dichotomous). A regression analysis, in short would allow us to test whether there is an association of PRS to SSRI response.

We will run two regression analyses.

The first will include the polygenic score as the outcome, adjusting for population ancestry (using 10 PCs), and the variable that we are interested in which is PRS (SCORE):

```
[26]: Fullmodel <- glm(data = dataset, PHENO-1 ~ SCORE + PC1 + PC2 + PC3 + PC4 + PC5 +
  PC6 + PC7 + PC8 + PC9 + PC10, family = "binomial")
```

The second model only includes 10 PCs and *excludes* (reduced) the PRS (SCORE).

Why? We will need to use this for a later analysis in order to see what heritability (trait differences) are we estimating from the PRS

Later we will calculate how much phenotypic variance is explained by the full model (PRS + PCs) and how much phenotypic variance is explained by the reduced model (PCs only).

The difference in variance explained between the two models is the amount of variance explained by the polygenic score term. Other than association of a tested trait, this is the primary question that is usually asked in a polygenic scoring study.

```
[27]: Reducedmodel <- glm(data = dataset, PHENO-1 ~ PC1 + PC2 + PC3 + PC4 + PC5 +  
                          PC6 + PC7 + PC8 + PC9 + PC10, family = "binomial")
```

Let's examine the models

```
[28]: summary(Fullmodel)
```

Call:

```
glm(formula = PHENO - 1 ~ SCORE + PC1 + PC2 + PC3 + PC4 + PC5 +  
     PC6 + PC7 + PC8 + PC9 + PC10, family = "binomial", data = dataset)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-1.249	-1.179	-1.060	1.175	1.307

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-0.04311	0.04307	-1.001	0.317
SCORE	-111.72604	289.59694	-0.386	0.700
PC1	-3.30588	4.91131	-0.673	0.501
PC2	-0.68915	3.43034	-0.201	0.841
PC3	0.48982	2.22353	0.220	0.826
PC4	-0.06881	2.19792	-0.031	0.975
PC5	-0.38738	2.00784	-0.193	0.847
PC6	-0.89865	2.00377	-0.448	0.654
PC7	-1.61608	2.00859	-0.805	0.421
PC8	1.27905	2.01003	0.636	0.525
PC9	0.60880	2.00185	0.304	0.761
PC10	-1.60740	2.00765	-0.801	0.423

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 3470.4 on 2503 degrees of freedom
Residual deviance: 3461.6 on 2492 degrees of freedom
AIC: 3485.6

Number of Fisher Scoring iterations: 3


```
[29]: summary(Reducedmodel)
```

Call:

```
glm(formula = PHENO ~ 1 + PC1 + PC2 + PC3 + PC4 + PC5 + PC6 +  
     PC7 + PC8 + PC9 + PC10, family = "binomial", data = dataset)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-1.236	-1.179	-1.063	1.175	1.301

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-0.03700	0.04004	-0.924	0.3555
PC1	-5.03497	2.00970	-2.505	0.0122 *
PC2	0.38581	2.00101	0.193	0.8471
PC3	0.86421	2.00060	0.432	0.6658
PC4	0.28212	2.00081	0.141	0.8879
PC5	-0.45073	2.00105	-0.225	0.8218
PC6	-0.93826	2.00109	-0.469	0.6392
PC7	-1.63442	2.00798	-0.814	0.4157
PC8	1.26208	2.00949	0.628	0.5300
PC9	0.60320	2.00171	0.301	0.7632
PC10	-1.66414	2.00219	-0.831	0.4059

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 3470.4 on 2503 degrees of freedom
Residual deviance: 3461.8 on 2493 degrees of freedom
AIC: 3483.8

Number of Fisher Scoring iterations: 3

2.8 Step 3.4 Phenotype/ Trait Variance that is explained by PRS

Let's now calculate the phenotype / trait variance (that's sometimes also known as heritability) explained by the polygenic scoring 'SCORE' term.

We will be using a function from a package that we initially loaded. The package is called `fmsb` and the function that we will be using is called `NagelkerkeR2`.

Nagelkerke's R² is a metric that we use to explain how much of a trait tested is measured by PRS. This ranges from 0-1.

We would usually hope and pray and wish it's as close to 1 but in reality, as we covered at the top of this workshop – barely 20% is ever explained!

This command now calculates Nagelkerke's R2. The subsequent command extracts just the R2 value.

```
[30]: Fullmodel.rsquare <- NagelkerkeR2(Fullmodel)[[2]]
      Fullmodel.rsquare
```

0.00467737481470882

Similarly, let's compute Nagelkerke's R2 for the reduced model

```
[31]: Reducedmodel.rsquare <- NagelkerkeR2(Reducedmodel)[[2]]
      Reducedmodel.rsquare
```

0.00459837968804138

Each of the values above explain variances in each model. Now, the rationale behind using reduced vs full model comes in.

We will take the difference between the two models to see how much is the variance being explained by MDD PRS to SSRI Response.

Let's now compute the R2 difference between the models.

```
[32]: diff.rsquare <- Fullmodel.rsquare - Reducedmodel.rsquare
```

2.9 Step 3.5 Results (FINALLY)

Let's now pull all of the results into one file so that we understand what our data is telling us.

First, we will obtain the full model's (including PRS) effect size and p-value for the polygenic score term. Note, again, that this is not significant.

```
[33]: Estimate <- coef(summary(Fullmodel))[2,c(1,4)]
```

Next, we will place these values into a **results** file that contains the difference in R2 **diff.rsquare**

In the subsequent **names(results)** line, we can name column headers manually in order to improve clarity.

```
[34]: results <- as.data.frame(cbind(Estimate[1],Estimate[2],diff.rsquare))
      names(results) <- c("OR","p-value","NagelkerkeR2")
```

So, what is the summary of our results?

drum roll please

```
[35]: results
```

	OR	p-value	NagelkerkeR2
Estimate	-111.726	0.699646	7.899513e-05

3 Breakout session

If we have the time, we will try out some practical assignments that we do in genetics.

As we eluded to earlier, the variance explained by PRS towards a trait is one of our primary goals. Based on the model that you have ran with the **pe2** tranche, we noticed an R^2 of 7.899513×10^{-5} .

Do you think that this will *increase* when we loosen up the threshold i.e. use **pe05**, **pe10**?

This would allow **more SNPs** to be used as reference from the reference data into our test data.

Do you think that this will *decrease* if we make this threshold more stringent i.e. use **pe4**, **pe6**?

This would allow **less SNPs** to be used as reference from the reference data into our test data.

Let's test this out!

3.1 Breakout room 1: Testing the INCREASE in p-value tranche threshold

Use **pe10**

File to use: OUTCOME/MDD.P10.profile

3.2 Breakout room 2: Testing the DECREASE in p-value tranche threshold

Use **pe4**

File to use: OUTCOME/MDD.pe4.profile

3.3 What do you have to do?

- 1) Introduce yourselves!
- 2) Identify a note-taker (and a back up, just in case). This person will also share their screen with the group for code reviewing.
- 3) Identify a reporter who will share your group's responses with the larger group.

Kumar will pop in and out of your rooms to check in; please use the "Ask for Help" button to bring Kumar into your group as and when needed

3.3.1 Steps-by-step guide

- 1) Copy and paste the analysis steps for **Step 3** (Steps 3.1 through 3.5) into separate cells below this cell. **OR** you can just edit the above cells. It's really up to you!
- 2) Choose the appropriate files to run your analysis. All you would have to do is substitute with the appropriate p-value tranche threshold that you were assigned. The command that you would edit is `dataset <- read.table("OUTCOME/MDD.pe2.profile",header=T)`
- 3) Within your group, discuss the
 - 3.1) [effect size](#) and p-value contributed by the MDD PRS to SSRI response
 - 3.2) Is the effect size the same?

3.3) More importantly, what is the Nagelkerke's R^2 and how does it compare to pe^2 ? What do you think this means?

3.4) Can genetics be any more complex than this? *groan*

We will broadcast a 2-minute warning when it's time to start wrapping up.

4) Present your group's work by sharing the screen of your reporter. *Each group will have 3-5 minutes to share their discussion points.*

[]: