

Introduction to Genetic Epidemiology

Shea J Andrews, PhD

3/14/24

Table of contents

I	Introduction	4
	Schedule	6
	Week 1 04/22 - 04/26	6
	Monday 04/22	6
	Wednesday 04/24	6
	Friday 04/26	7
	Week 2 04/29 - 05/03	8
	Monday 04/29	8
	Wednesday 05/01	8
	Friday 05/03	9
	Week 3 05/06 - 05/10	9
	Monday 05/06	9
	Wednesday 05/08	10
	Friday 05/10	10
1	Conda Env	11
	HABS-HD	12
	Phenotypes	12
	Descriptive Table	14
	Genotyping	19
	HapMap III	19
II	Genome-wide Association Studies	22
	GWAS QC	23
	SNP QC	23
	Call Rate & Allele frequency	23
	Hardy Weinberg Equilibrium	25
	Sample QC	26
	Call Rate	26
	Sex Discordance	27
	Pruning	28
	Heterozygosity check	28

Cryptic Relatedness	29
Population Substructure	33
Scree Plot	34
Population substructure	34
Exclude Samples	35
GWAS	36
Import	36
Phenotype File	36
Covariate File	36
GWAS	37
Manhattan Plot	37
GWAS-SS	38
AD GWAS	38
MungeSumstats	39
Manhattan Plot	39
Genetic Ancestry	40
Principal Component Analysis	40
ADMIXTURE	40
 III Polygenic Risk Scores	 41
3 Polygenic Risk Scores	42
 IV Mendelian Randomization	 43
4 Mendelian Randomization	44
 Appendices	 44
References	45

Part I

Introduction

This is a Quarto book.

To learn more about Quarto books visit <https://quarto.org/docs/books>.

Schedule

Welcome to this BMI Mini-Course on an Methods in Genetic Epidemiology. The course structure is divided into three key methodologies: genome-wide association studies, polygenic risk scores, and Mendelian randomization. Each trainee will apply these methods to a trait of their choosing and, at the course's conclusion, present their findings to the class.

Week 1 04/22 - 04/26

Monday 04/22

Alzheimer's disease

In this mini-course, we will frequently reference Alzheimer's disease to illustrate various genetic epidemiology methods. This session aims to introduce Alzheimer's disease and explore significant findings related to its genetic architecture.

Readings

- Knopman, D. S. et al. Alzheimer disease. [Nat Rev Dis Primers 7, 33 \(2021\)](#).
- Alzheimer's Association. 2023 Alzheimer's disease facts and figures. [Alzheimer's Dementia \(2023\)](#).
- Kornblith, E. et al. Association of Race and Ethnicity With Incidence of Dementia Among Older Adults. [Jama 327, 1488–1495 \(2022\)](#).
- Andrews, S. J. et al. The complex genetic architecture of Alzheimer's disease: novel insights and future directions. [eBioMedicine 90, 104511 \(2023\)](#).
- Andrews, S. J., Fulton-Howard, B. & Goate, A. Interpretation of risk loci from genome-wide association studies of Alzheimer's disease. [Lancet Neurology 19, 326–335 \(2020\)](#).

Wensday 04/24

Genome-Wide Association Studies

Genome-Wide Association Studies (GWAS) are foundational to various genetic analysis methodologies. In this session, we will delve into what a GWAS entails, the process of conducting one, and then engage in a hands-on exercise to carry out our own GWAS.

Readings

- Uffelmann, E. et al. Genome-wide association studies. [Nat Rev Methods Primers 1, 59 \(2021\)](#).
- Abdellaoui, A., Yengo, L., Verweij, K. J. H. & Visscher, P. M. 15 years of GWAS discovery: Realizing the promise. [Am J Hum Genetics \(2023\)](#)
- Marees, A. T. et al. A tutorial on conducting genome-wide association studies: Quality control and statistical analysis. [Int J Method Psych 27, e1608 \(2018\)](#).
- MacArthur, J. A. L. et al. Workshop proceedings: GWAS summary statistics standards and sharing. [Cell Genom 1, 100004 \(2021\)](#).

Tools

- [PLINK](#)¹
- [BCFTools](#)²
- [MungeSumStats](#)³

Friday 04/26

Genetic Ancestry

Genetic ancestry explores the lineage and heritage inferred from our DNA, providing insights into population history and individual heritage. This session will introduce the concepts and methodologies used in determining genetic ancestry, emphasizing their importance in genetic epidemiology research.

Readings

- Lewis, A. C. F. et al. Getting genetic ancestry right for science and society. [Science 376, 250–252 \(2022\)](#).
- National Academies of Sciences, Engineering, and Medicine. 2023. Using Population Descriptors in Genetics and Genomics Research: A New Framework for an Evolving Field. [Washington, DC: The National Academies Press](#).

Tools

- [ADMIXTURE](#)⁴
- [RFmix](#)⁵

Week 2 04/29 - 05/03

Monday 04/29

Heritability & Genetic Correlations

Heritability quantifies the proportion of phenotype variance attributable to genetic factors, whereas genetic correlations assess the extent of shared genetic architecture between traits. In this session, we will concentrate on the tools utilized to estimate these metrics from GWAS summary statistics.

Readings

- Rhee, W. van, Peyrot, W. J., Schork, A. J., Lee, S. H. & Wray, N. R. Genetic correlations of polygenic disease traits: from theory to practice. [Nat Rev Genetics 20, 567–581 \(2019\)](#).
- Barry, C.-J. S. et al. How to estimate heritability: a guide for genetic epidemiologists. [Int J Epidemiol \(2022\)](#)

Tools

- [LDSC](#)⁶
- [HDL](#)⁷
- [GenomicSEM](#)⁸

Wednesday 05/01

Polygenic Risk Scores I

Polygenic risk scores (PRS) measure an individual's total genetic liability for a trait. This session will cover the process of constructing a PRS and assessing its performance in predicting the trait.

Readings

- Choi, S. W., Mak, T. S.-H. & O'Reilly, P. F. Tutorial: a guide to performing polygenic risk score analyses. [Nat Protoc 15, 2759–2772 \(2020\)](#).
- Wand, H. et al. Improving reporting standards for polygenic scores in risk prediction studies. [Nature 591, 211–219 \(2021\)](#).
- Lennon, N. J. et al. Selection, optimization and validation of ten chronic disease polygenic risk scores for clinical implementation in diverse US populations. [Nat. Med. 1–8 \(2024\)](#)

Tools

- [PRSice2](#)⁹
- [PRSet](#)¹⁰

Friday 05/03

Polygenic Risk Scores II

The accuracy of polygenic risk scores (PRS) diminishes as the genetic distance from the training population increases. This session will explore cross-ancestry PRS methods designed to enhance PRS accuracy across diverse populations.

Readings

- Kachuri, L. et al. Principles and methods for transferring polygenic risk scores across global populations. [Nat. Rev. Genet. 1–18 \(2023\) doi:10.1038/s41576-023-00637-2](#).
- Ding, Y. et al. Polygenic scoring accuracy varies across the genetic ancestry continuum. [Nature 618, 774–781 \(2023\)](#).

Tools

- [PRS-CSx](#)¹¹

Week 3 05/06 - 05/10

Monday 05/06

Mendelian Randomization I

Mendelian Randomization (MR) is a method employed to identify causal risk factors for diseases. This session will cover the fundamentals of MR and demonstrate how to execute a two-sample MR analysis.

Readings

- Sanderson, E. et al. Mendelian randomization. [Nat Rev Methods Primers 2, 6 \(2022\)](#).
- Davies, N. M., Holmes, M. V. & Smith, G. D. Reading Mendelian randomisation studies: a guide, glossary, and checklist for clinicians. [BMJ 362, k601 \(2017\)](#).
- Hemani, G. et al. The MR-Base platform supports systematic causal inference across the human phenome. [Elife 7, e34408 \(2018\)](#).

Tools

- [TwoSampleMR](#)¹²

Wednesday 05/08

Mendelian Randomization II

A crucial aspect of Mendelian Randomization (MR) studies is assessing whether the causal associations derived from MR analyses remain valid despite potential violations of MR's underlying assumptions. This session will focus on diagnostic and sensitivity analyses in MR, along with guidance on effectively reporting MR findings.

Readings

- Skrivankova, V. W. et al. Strengthening the Reporting of Observational Studies in Epidemiology Using Mendelian Randomization. [JAMA 326, 1614–1621 \(2021\)](#).
- Skrivankova, V. W. et al. Strengthening the reporting of observational studies in epidemiology using mendelian randomisation (STROBE-MR): explanation and elaboration. [BMJ 375, n2233 \(2021\)](#).

Friday 05/10

Trainee Presentations

Upon concluding this mini-course, trainees will showcase the results of their analyses.

1 Conda Env

```
#| eval: false
name: genetic_epi
channels:
  - conda-forge
  - bioconda
dependencies:
  - plink=1.90b6.21
```

```
#| eval: false
conda env create -f env/genetic_epi.yml
```

HABS-HD

Phenotypes

```
library(tidyverse)
setwd('~/.gitcode/IntroGeneticEpi/')
SAVE_VISUALIZATIONS_PATH <- "~/Downloads"

aa_v1_path = 'resources/HABSHD/v5/HD 1 African American 50+ Request 355.csv'
ma_v1_path = 'resources/HABSHD/v5/HD 1 Mexican American 50+ Request 355.csv'
nhw_v1_path = 'resources/HABSHD/v5/HD 1 Non-Hispanic White 50+ Request 355.csv'

hd_cols = spec(read_csv(nhw_v1_path, guess_max = 10000))
aa_v1.raw = read_csv(aa_v1_path, col_types = hd_cols, na = c("", "NA", "9999", "-9999", "
  janitor::clean_names()
ma_v1.raw = read_csv(ma_v1_path, col_types = hd_cols, na = c("", "NA", "9999", "-9999", "
  janitor::clean_names()
nhw_v1.raw = read_csv(nhw_v1_path, col_types = hd_cols, na = c("", "NA", "9999", "-9999", "
  janitor::clean_names()

habshd.raw <- bind_rows(
  aa_v1.raw, ma_v1.raw, nhw_v1.raw
) %>%
mutate(
  id_race_white = as.factor(id_race_white),
  id_race_black = as.factor(id_race_black),
  id_race_indian_alaska = as.factor(id_race_indian_alaska),
  id_race_asian = as.factor(id_race_asian),
  id_race_japanese = as.factor(id_race_japanese),
  id_race_korean = as.factor(id_race_korean),
  id_race_vietnamese = as.factor(id_race_vietnamese),
  id_race_native_hawaiian = as.factor(id_race_native_hawaiian),
  id_race_guam_chamorro = as.factor(id_race_guam_chamorro),
  id_race_samoan = as.factor(id_race_samoan),
```

```

id_race_other_pacific = as.factor(id_race_other_pacific),
id_race_other = as.factor(id_race_other),
id_hispanic = as.factor(id_hispanic),
id_hispanic_other = as.factor(id_hispanic_other),
race = case_when(
  id_hispanic != 1 ~ "Hispanic",
  id_race_white == 1 & id_hispanic != 2 ~ "NHW",
  id_race_black == 1 ~ "Black",
  TRUE ~ "Other")
)

habshd <- habshd.raw %>%
  mutate(
    abeta40 = ifelse(is.na(r3_qtx_plasma_abeta42), r5_qtx_plasma_abeta40, r3_qtx_plasma_abeta40),
    abeta42 = ifelse(is.na(r3_qtx_plasma_abeta42), r5_qtx_plasma_abeta42, r3_qtx_plasma_abeta42),
    ptau181 = ifelse(is.na(r3_qtx_plasma_ptau181), r5_qtx_plasma_ptau181, r3_qtx_plasma_ptau181),
    total_tau = ifelse(is.na(r3_qtx_plasma_total_tau), r5_qtx_plasma_total_tau, r3_qtx_plasma_total_tau),
    nfl = ifelse(is.na(r3_qtx_plasma_nf_l), r5_qtx_plasma_nf_l, r3_qtx_plasma_nf_l)
  ) %>%
  select(
    med_id, age, id_gender, interview_language, adi_state_rank, race,
    id_education, smoke_ever, cdx_cog, cdx_depression, cdx_hypertension,
    cdx_diabetes, cdx_dyslipidemia, cdr_sum,
    om_bp1_dia, om_bp1_sys,
    om_height, om_weight, om_bmi, om_ab_circumference,
    bw_chol_total, bw_ldl_chol, bw_hdl_chol, bw_hba1c, gds_total,
    abeta40, abeta42, ptau181, total_tau, nfl,
    apoe4_snp
  )

write_csv(habshd, "work/habshd_pheno.csv")

#descriptive table
description <- c("Medical ID","", "1 = Female <br> 0 = Male",
  "Language in which <br> interview was administered",
  "Area Deprivation Index", "Black, Hispanic, NHW",
  "Years of Education", "Ever Smoked<br> (1:Yes, 0:No)",
  "Cognitive Disorder:<br> 0: Cognitively Unimpaired <br>1:Mild Cognitive Impairment<br>2: Dementia", "Depression<br> (1:Yes, 0:No)",
  "Hypertension <br> (1:Yes,0:No)", "Diabetes <br> (1:Yes,0:No)",
  "High Cholesterol (1:Yes,0:No)",
  "Clinical Dementia Rating (CDR):<br> Sum of Boxes",

```

```

      "Diastolic BP", "Systolic BP", "Height (in)", "Weight (lbs)",
      "BMI", "Abdominal circumference (in)", "Total Cholesterol",
      "LDL Cholesterol<br> (bad)", "HDL Cholesterol<br> (good)",
      "Hemoglobin", "Geriatric Depression Scale (GDS)",
      "abeta40", "abeta42", "ptau181", "total_tau", "nfl",
      "APOE Genotype")

table_desc <- data.frame(cbind(names(habshd), description))
table_desc %>% kbl(caption = '', col.names = c("Variable", "Description"),
                  escape = FALSE) %>%
  kable_classic(full_width = FALSE, html_font = "Ariel") %>%
  kable_styling(font_size = 16, position = "center") %>%
  column_spec(1:2, border_left = F, border_right = F) %>%
  pack_rows("Demographics", 1, 7) %>%
  pack_rows("Clinical", 8, 25) %>%
  pack_rows("Imaging", 26, 30) %>%
  pack_rows("Genomics", 31, 31)

```

Descriptive Table

Variable	Description
Demographics	
med_id	Medical ID number
	Not an MRN
age	Age (yrs)
id_gender	1:Female
	0: Male
interview_language	Language in which interview was administered
	1:English 2:Spanish
adi_state_rank	Area Deprivation Index
	Levels: 1,2,...,10
	1: least disadvantaged 10: most disadvantaged
race	Black, NHW, Hispanic
id_education	Years of Education
Clinical	
smoke_ever	Ever smoked? 0:No 1:Yes

cdx_cog	Cognitive Disorder 0: Cognitively Unimpaired 1:Mild Cognitive Impairment 2: Dementia
cdx_depression	Depression 0:No 1:Yes
cdx_hypertension	Hypertension 0:No 1:Yes
cdx_diabetes	Diabetes 0:No 1:Yes
cdx_dyslipidemia.	High Cholesterol 0:No 1:Yes
cdr_sum	Clinical Dementia Rating (CDR) Sum of Boxes
gds_total	Geriatric Depression Scale (GDS) sum of GDS 1 to GDS 30
om_bp1_dia	Diastolic BP
om_bp1_sys	Systolic BP
om_height	Height (in)
om_weight	Weight (lbs)
om_bmi	Body Mass Index (BMI)
om_ab_circumference	Abdominal Circumference (in)
bw_chol_total	Total Cholesterol (mg/dL)
bw_ld_lchol	LDL Cholesterol (mg/dL) (bad)
bw_hdl_chol	HDL Cholesterol (mg/dL) (good)
bw_hba1c	Hemoglobin A1C% of total Hgb
Imaging	
abeta40	$A\beta_{40}$
abeta42	$A\beta_{42}$
ptau181	Phospho-Tau (pg/mL) Average CV: 0.07065 Average LLOD: 0.016 Average HLOD:349
total_tau	Total Tau
nfl	Neurofilament Light (pg/mL) Average CV: 0.038 Average LLOD: 0.038 Average HLOD:1800
Genetics	
apoe4_snp	APOE Genotype E2E3, E2E4, E3E3, E3E4, E4E4

```

habshd[which(habshd$adi_state_rank=="GQ"),] <- NA
habshd[which(habshd$adi_state_rank=="PH"),] <- NA

```

```

habshd[which(habshd$adi_state_rank=="Invalid Address"),] <- NA
habshd[which(habshd$smoke_ever==2),] <- NA
habshd$adi_state_rank <- as.integer(habshd$adi_state_rank)

theme_gtsummary_compact()
demographics_table <- habshd %>% select(age,id_gender,interview_language,
                                       adi_state_rank, race, id_education) %>%
  tbl_summary(., by = race,
    statistic = list(
      all_continuous() ~ "{mean}<br> ({sd})",
      all_categorical()~ "{p}%",
      digits = all_continuous()~2,
      label = c(age~"Age", id_gender ~ "Gender",
        interview_language ~ "Interview Language",
        adi_state_rank~ "ADI State Rank",
        id_education~"Education"),
      missing_text = "(Missing)" %>%
      modify_header(label = "**Demographic <br> Variables**",
        all_stat_cols() ~ "**{level}**<br> N = {n}") %>%
      as_gt() %>%
      tab_options(column_labels.border.top.color = "black",
        column_labels.border.bottom.color = "black",
        table_body.border.bottom.color = "black",
        table_body.hlines.color = "white",
        table.font.size = 12,
        container.width = 500,
        container.height =500) %>%
      fmt_markdown(columns = everything())

gtsave(demographics_table,filename = file.path(SAVE_VISUALIZATIONS_PATH, "demographics_sum

#convert to factors
cdx_cols <- names(habshd %>% select(starts_with("cdx_")))
habshd[cdx_cols] <-lapply(habshd[cdx_cols], factor)

clinical_table1 <- habshd %>% select(smoke_ever,cdx_cog,cdx_depression,
                                       cdx_hypertension,cdx_diabetes,
                                       cdx_dyslipidemia,cd_r_sum, om_bp1_dia,
                                       race) %>%
  tbl_summary(., by = race,
    statistic = list(

```



```

        all_continuous() ~ "{mean} ({sd})",
        all_categorical()~ "{p}%",
    digits = all_continuous()~2,
    label = c(smoke_ever ~ "Smoke ",
              cdx_cog ~ "Cognitive Disorder",
              cdx_depression ~ "Depression",
              cdx_hypertension ~ "Hypertension",
              cdx_diabetes ~ "Diabetes",
              cdx_dyslipidemia ~ "Displeidemia",
              cdr_sum ~ "CDR Total Score",
              om_bp1_dia ~ "Diastolic BP"),

    missing_text = "(Missing)" %>%
    modify_header(label = "**Clinical Variables**",
                  all_stat_cols() ~ "**{level}**<br> N = {n}") %>%
    as_gt() %>%
    tab_options(
      column_labels.border.top.color = "black",
      column_labels.border.bottom.color = "black",
      table_body.border.bottom.color = "black",
      table_body.hlines.color = "white",
      table.font.size = 12,
      container.height = 700,
      container.width = 700) %>%
    fmt_markdown(columns = everything())

gtsave(clinical_table1,filename = file.path(SAVE_VISUALIZATIONS_PATH, "clinical_table1.png)

clinical_table2 <- habshd %>% select(om_bp1_sys,om_height,om_weight,
                                   om_bmi,om_ab_circumference,bw_chol_total,
                                   bw_ld_lchol,bw_hdl_chol,race,bw_hba1c,
                                   gds_total,race) %>%
tbl_summary(., by = race,
            statistic = list(
              all_continuous() ~ "{mean} ({sd})",
              all_categorical()~ "{p}%",
            digits = all_continuous()~2,
            label = c(om_bp1_sys~"Systoliuc BP",
                      om_height ~ "Height (in)",
                      om_weight~ "Weight(lbs)",
                      om_bmi ~"BMI",
                      om_ab_circumference~ "Abdominal <br>

```

```

        Circumference (in)",
        bw_chol_total ~ "Total Cholesterol",
        bw_ld_lchol~ "LDL <br> Cholesterol",
        bw_hdl_chol~ "HDL <br> Cholesterol",
        bw_hba1c~ "Hemoglobin",
        gds_total ~ "GDS Total"),
    missing_text = "(Missing)" %>%
  modify_header(label = "**Clinical <br> Variables**",
    all_stat_cols() ~ "**{level}**<br> N = {n}") %>%
  as_gt() %>%
  tab_options(
    column_labels.border.top.color = "black",
    column_labels.border.bottom.color = "black",
    table_body.border.bottom.color = "black",
    table_body.hlines.color = "white",
    table.font.size = 12,
    container.height = 700,
    container.width = 700) %>%
  fmt_markdown(columns = everything())
gtsave(clinical_table2,filename = file.path(SAVE_VISUALIZATIONS_PATH, "clinical_table2.png"))

habshd$apoe4_snp = as.factor(habshd$apoe4_snp)
imaging_genetics_table <- habshd %>% select(abeta40, abeta42,ptau181, total_tau,
  nfl, apoe4_snp,race) %>%
  tbl_summary(., by = race,
    statistic = list(
      all_continuous() ~ "{mean} ({sd})",
      all_categorical()~ "{p}%",
    digits = all_continuous()~2,
    label = c(abeta40~"AB40",
      abeta42~"AB42",
      ptau181 ~ "pTau",
      total_tau ~ "Total Tau",
      nfl~ "Plasma NFL",
      apoe4_snp ~ "APOE4 SNP"),
    missing_text = "(Missing)" %>%
  modify_header(label = "**Imaging & Genetic <br> Variables**",
    all_stat_cols() ~ "**{level}**<br> N = {n}") %>%
  as_gt() %>%
  tab_options(

```

```

        column_labels.border.top.color = "black",
        column_labels.border.bottom.color = "black",
        table_body.border.bottom.color = "black",
        table_body.hlines.color = "white",
        table.font.size = 12,
        container.width = 700,
        container.height = 700) %>%
    fmt_markdown(columns = everything())

gtsave(imaging_genetics_table, filename = file.path(SAVE_VISUALIZATIONS_PATH, "imaging_gene

! [Table1] (~ /Downloads/demographics_summary.png)

! [Table2] (~ /Downloads/clinical_table1.png)
! [Table3] (~ /Downloads/clinical_table2.png)
! [Table4] (~ /Downloads/imaging_genetic_table.png)

```

Genotyping

TBD

HapMap III

Download the hapmap_3.3.hg38.vcf.gz file from the [Broad's google bucket](#)

```

bcftools view -i 'AF > 0 && TYPE="snp" && N_ALT=1' resources/genetic_epi/resources_broad_h
bcftools view -H > work/hapmap3_snps.txt

hm3.raw <- read_table("work/hapmap3_snps.txt", col_names = F)

hm3 <- hm3.raw %>%
  mutate(
    cpri = glue::glue("{X1}:{X2}:{X4}:{X5}"),
    X1 = as.numeric(str_replace(X1, 'chr', ''))
  ) %>%
  filter(!is.na(X1)) %>%
  rename(chr = X1, pos = X2, rsid = X3, ref = X4, alt = X5) %>%
  select(-X6)

```

```

out <- hm3 %>%
  distinct(cpra, .keep_all = T) %>%
  distinct(rsid, .keep_all = T)

out %>%
  select(cpra) %>%
  write_tsv(., 'work/hm3_extract.txt', col_names = F)

out %>%
  select(cpra, rsid) %>%
  write_tsv(., 'work/hm3_crpa_rsids.txt', col_names = F)

plink \
  --bfile resources/HABSHD/genotypes/all \
  --keep-allele-order \
  --extract work/hm3_extract.txt \
  --make-bed \
  --out work/habshd_hm3

plink \
  --bfile work/habshd_hm3 \
  --keep-allele-order \
  --update-name work/hm3_crpa_rsids.txt \
  --make-bed \
  --out work/habshd_rsids

```

2

Part II

Genome-wide Association Studies

GWAS QC

TBD

SNP QC

SNP level QC consists of removing markers with excessive missingness or low allele frequency. This QC increases the power to identify true associations with disease risk by removing sub-optimal markers that can increase false positives.

Call Rate & Allele frequency

95% was used as the SNP call rate threshold (usually 95% or higher), and 1% was used as the MAF threshold (usually 1% or higher).

Filtering SNPs on MAF and call rate can be done in PLINK 1.9 by typing the following (or similar) at the shell prompt. This uses 95% and 1% for the call-rate and MAF, respectively:

```
# Generate frequency reports
plink \
  --bfile work/habshd_rsid \
  --keep-allele-order \
  --freq \
  --out work/habshd_snpqc

plink \
  --bfile work/habshd_rsid \
  --keep-allele-order \
  --freqx \
  --out work/habshd_snpqc

# Filter on call rate and maf
plink \
  --bfile work/habshd_rsid \
  --keep-allele-order \
```

```

--geno 0.05 --maf 0.01 \
--make-bed --out work/habshd_snpqc

## ==== SNP Level Filtering ====
# ---- readin plink .frq ---- ##
message("reading plink frq file")
freq.raw <- read_table('work/habshd_snpqc.frq', col_names = T,
  col_types = cols(
    CHR = col_double(),
    SNP = col_character(),
    A1 = col_character(),
    A2 = col_character(),
    MAF = col_double(),
    NCHROBS = col_double()
  ))

# ---- readin plink .frqx ---- ##
message("reading plink frqx file")
freqx.raw <- read_tsv('work/habshd_snpqc.frqx', col_names = T,
  col_types = cols(
    CHR = col_double(),
    SNP = col_character(),
    A1 = col_character(),
    A2 = col_character(),
    `C(HOM A1)` = col_double(),
    `C(HET)` = col_double(),
    `C(HOM A2)` = col_double(),
    `C(HAP A1)` = col_double(),
    `C(HAP A2)` = col_double(),
    `C(MISSING)` = col_double()
  ))

# ---- SNP level statisitcs ----
snps <- freq.raw %>%
  full_join(freqx.raw, by = c("CHR", "SNP", "A1", "A2")) %>%
  rename(AA = `C(HOM A1)`, AB = `C(HET)`, BB = `C(HOM A2)`, missing = `C(MISSING)`) %>%
  mutate(Call.rate = 1 - (missing / c(AA + AB + BB + missing))) %>%
  mutate(Call = Call.rate >= 1 - 0.05) %>%
  mutate(Call.maf = MAF < 0.01)

```

Figure @ref(fig:MAFxcallrate) shows the SNP call rate versus minor allele frequency across all typed SNPs in the study. The dashed lines denote the MAF and call rate QC thresholds. xxx

SNPs were removed due to low call rate and xxx SNPs were removed due to low minor allele frequency.

```
MAFxcallrate.p <- ggplot(data = snps, aes(x = MAF, y = Call.rate)) +  
  geom_point(alpha = 0.3, size = 0.5) +  
  geom_hline(yintercept = 1 - 0.05, linetype = 2, colour = 'red') +  
  geom_vline(xintercept = 0.01, linetype = 2, colour = 'red') +  
  scale_x_log10(breaks = scales::trans_breaks("log10", function(x) round(10^x, 3))) +  
  labs(y = 'Proportion of called genotypes', x = 'Minor Allele Frequency (log)') +  
  theme_bw() + annotation_logticks()  
  
ggsave('results/plots/MAFxcallrate.png', plot = MAFxcallrate.p, height = 4, width = 6, uni
```

Hardy Weinberg Equilibrium

Violations of Hardy Weinberg Equilibrium can indicate either the presence of population sub-structure, or the occurrence of genotyping error. It is common practice to assume that violations are indicative of genotyping error and remove SNPs in which the HWE test statistic has a corresponding p-value of less than 1×10^{-6} . A threshold of xxx is used here.

For case-control data, HWE is generally not tested in cases to not exclude real selection against a phenotype, so it is best to include case-control status in the PLINK files.

Filtering SNPs on Hardy Weinberg Equilibrium for autosomes only can be done in PLINK by typing the following at the shell prompt:

```
plink \  
  --bfile work/habshd_snpqc \  
  --keep-allele-order \  
  --autosome \  
  --hardy \  
  --hwe 0.000001 \  
  --make-bed --out work/habshd_hwe  
  
# ---- readin plink .hwe ---- ##  
message("reading plink hwe file")  
hwe.raw <- read_table2('work/habshd_hwe.hwe', col_types = cols(  
  CHR = col_integer(),  
  SNP = col_character(),  
  TEST = col_character(),  
  A1 = col_character(),
```

```

    A2 = col_character(),
    GENO = col_character(),
    `O(HET)` = col_double(),
    `E(HET)` = col_double(),
    P = col_double()
  ))

snps <- snps %>%
  full_join(hwe.raw, by = c("CHR", "SNP", "A1", "A2")) %>%
  mutate(hwe = P > 0.000001) %>%
  as_tibble()

suppressPackageStartupMessages(library(ggtern))

hweplot <- snps %>%
  filter(!is.na(P)) %>%
  mutate(alpha = ifelse(hwe, 0.2, 0.8),
         hwe = ifelse(hwe, "Pass", "Fail")) %>%
  ggtern::ggtern(aes(x = AA, y = AB, z = BB, colour = hwe, alpha = alpha)) +
  geom_point(size = 0.5) +
  scale_colour_manual(name = 'Hardy Weinberg \n Equilibrium',
                     values = c(Pass = "#377EB8", Fail = "#E41A1C")) +
  scale_alpha_continuous(guide = "none", range = c(0.8, 0.2)) +
  theme_bw() + theme(legend.position = 'bottom')

hweplot

detach("package:ggtern", unload=TRUE)

ggsave('results/plots/hweplot.png', plot = hweplot, height = 4, width = 6, units = 'in')

```

Sample QC

Call Rate

A low genotyping call rate in a sample can be indicative of poor DNA sample quality, so samples with a call rate < xxx% are excluded from further analysis.

Filtering samples on a call rate of 95% can be done in PLINK by typing the following at the shell prompt:

```
plink \
  --bfile work/habshd_hwe \
  --keep-allele-order \
  --mind 0.05 \
  --make-bed --out work/habshd_sampleQC
```

Sex Discordance

Samples with discordance between self-reported and genetically predicted sex likely have errors in sample handling, such as sample swaps. Predicted sex can be determined by calculating X chromosome heterozygosity using an F test, because biological men have one X chromosome and women have two. An F value of ~ 0.99 indicates males, and an F value of ~ 0.03 indicates females. Furthermore, checking X chromosome heterozygosity may reveal sex chromosome anomalies (~ 0.28 in reported females; ~ 0.35 in males).

Since sex discordance may be due to sample swaps or to incorrect phenotyping, sex discordant samples should generally be removed unless a swap can be reliably resolved.

Identification of individuals with discordant sex can be done in PLINK 1.9 by typing the following at the shell prompt, which will produce a list of individuals with discordant sex data.

```
plink \
  --bfile resources/HABSHD/genotypes/HABLE_GSA_20230418a_FINAL \
  --check-sex --out work/HABLE_GSA_20230418a
```

```
plink \
  --bfile resources/HABSHD/genotypes/HABLE_GSA_20220602_FINAL \
  --check-sex --out work/HABLE_GSA_20220602
```

```
awk 'FNR==1 && NR==1 {print; next} FNR>1 {print}' work/HABLE_GSA_20220602.sexcheck work/HABLE_GSA_20220602.sexcheck
```

```
## ---- Read in Data ----##
sexcheck.raw <- read_table2('work/habshd_sexcheck.txt')

## recode sex variables
sexcheck <- sexcheck.raw %>%
  mutate(PEDSEX = recode(PEDSEX, '2' = 'Female', '1' = 'Male'))

## Exclude samples with no sex inconsistencies
sex_exclude.samples <- sexcheck %>%
```

```
filter(STATUS == 'PROBLEM') %>%
mutate(PEDSEX = recode(PEDSEX, '2' = 'Female', '1' = 'Male'))
```

The following plot (Fig. @ref(fig:sexplot)) displays the X Chromosome heterozygosity for self reported sex, with samples with problems highlighted in red. Table @ref(tab:sextab) displays the individule records that should be excluded from further downstream analysis.

```
sexcheck.p <- ggplot(data = sexcheck, aes(x = as.factor(PEDSEX), y = F, colour = STATUS, s
geom_jitter() +
scale_color_manual(values = c( "#377EB8", "#E41A1C")) +
theme_bw() + labs(x = 'Self reported sex', y = 'X CHR Heterozygosity (F)') + theme(legende

ggsave('results/plots/sexcheck.png', plot = sexcheck.p, height = 4, width = 6, units = 'in
```

Pruning

Pruning is typically done to remove linkage disequilibrium (LD) between SNPs, which is often a necessary step in various genetic analyses to ensure the independence of markers and is necessary for estimating heterozygosity, relatedness, and population stratification.

```
plink \
--bfile work/habshd_sampleQC \
--indep-pairwise 50 5 0.2 \
--out work/indepSNP
```

Heterozygosity check

Insufficient heterozygosity can indicate inbreeding or other family substructures, while excessive heterozygosity may indicate poor sample quality.

Individuals with outlying heterozygosity rates can be identified in PLINK 1.9 by typing the following command at the shell prompt:

```
plink \
--bfile work/habshd_sampleQC \
--extract work/indepSNP.prune.in \
--het --out work/habshd
```

This produces a file containing Method-of-moments F coefficient estimates, which can be used to calculate the observed heterozygosity rate in each individual. Analysis is performed using an LD pruned snplist.

We calculate a heterozygosity similarly using observed and expected counts from the PLINK output $[(\text{Observed} - \text{Expected})/N]$ and exclude samples that are ± 3 sd from the cohort mean.

```
## ---- Read in Data ----##
het.raw <- read_table('work/habshd.het')

## calculate heterozygosity
het <- het.raw %>%
  rename(O = `O(HOM)`, E = `E(HOM)`, N = `N(NM)`) %>%
  mutate(Het = (N - O) / N)

## Calculate exclusion thresholds
upper.het <- mean(het$Het) + sd(het$Het)*3
lower.het <- mean(het$Het) - sd(het$Het)*3

## Exclusion of samples
het <- het %>%
  mutate(exclude = ifelse(Het >= upper.het | Het <= lower.het, TRUE, FALSE))

het_exclude.samples <- het %>% filter(exclude == TRUE)
```

Figure @ref(fig:plothet) displays the distribution of heterozygosity in xxx. Samples with excessive ($\text{Het} > \text{xxx}$) or deficient ($\text{Het} < \text{xxx}$) heterozygosity are colored red. Table @ref(tab:het) displays samples that are to be excluded.

```
heterozygosity.p <- ggplot(het, aes(x = Het, fill = exclude)) + geom_histogram(binwidth =
  geom_vline(xintercept = upper.het, colour = 'red', linetype = 2) +
  geom_vline(xintercept = lower.het, colour = 'red', linetype = 2) +
  theme_bw() + scale_fill_manual(values = c("#377EB8", "#E41A1C")) +
  theme(legend.position = 'bottom') +
  labs(x = 'Heterozygosity')

ggsave('results/plots/heterozygosity.png', plot = heterozygosity.p, height = 4, width = 6,
```

Cryptic Relatedness

Population based cohorts are often limited to unrelated individuals as associations statistics often assume independence across individuals. Closely related samples will share more of their genome and are likely to be more phenotypically similar than two individuals chosen randomly from the population. A common measure of relatedness is identity by descent (IBD),

where a kinship correlation coefficient (π -hat) greater 0.1 suggests that samples maybe related or duplicates samples.

```
# IBD relationship table
# https://github.com/WheelerLab/GWAS_QC/blob/master/example_pipelines/QC%20Analysis%20-%20

rel_tab <- tibble(relationship = c("unrelated", "identical-twins",
                                "parent-child", "full-siblings",
                                "half-siblings", "grandparent-grandchild",
                                "avuncular", "half-avuncular",
                                "first-cousin", "half-first-cousin",
                                "half-sibling-first-cousin"),
  pi_hat = c(0, 1, 0.5, 0.5, 0.25, 0.25, 0.25, 0.125, 0.125, 0.0625, 0.375),
  z0 = c(1, 0, 0, 0.25, 0.5, 0.5, 0.5, 0.75, 0.75, 0.875, 0.375),
  z1 = c(0, 0, 1, 0.5, 0.5, 0.5, 0.5, 0.25, 0.25, 0.125, 0.5),
  z2 = c(0, 1, 0, 0.25, 0, 0, 0, 0, 0, 0, 0.125)
)

dup_relationships <- c("grandparent-grandchild", "avuncular", "half-avuncular")
rel_tab_filt <- rel_tab %>%
  filter(relationship %nin% dup_relationships) %>%
  mutate(relationship = ifelse(relationship == "half-siblings", "2nd degree",
                              ifelse(relationship == "first-cousin",
                                      "3rd degree", relationship)))
```

Identifying duplicated or related samples can be done in PLINK 1.9 by typing the following command at the shell prompt.

```
plink \
  --bfile work/habshd_sampleQC \
  --extract work/indepSNP.prune.in \
  --genome --min 0.05 --out work/habshd.ibd

# select samples with kinship coefficients > 0.1875
# https://link.springer.com/protocol/10.1007/978-1-60327-367-1_19
pi_hat_thres = 0.1875

# Find closest match
closest <- function(vals, ref) {
  fc <- Vectorize(function(x) {
    ref[which.min(abs(ref - x))]
  })
}
```

```

    }) #finds closest
    fc(vals)
  }

# Iteratively Remove related samples
remove_samples <- function(ibdcoeff, fam, msg = "closely related to") {
  fam_fi <- fam %>%
    mutate(FI = paste0(FID, "--tempsep--", IID)) %>%
    mutate(status = ifelse(status > 2, 0.5, status))

  ibdcoeff %<>%
    mutate(FI1 = paste0(FID1, "--tempsep--", IID1),
           FI2 = paste0(FID2, "--tempsep--", IID2))
  related_samples <- NULL
  excluded <- c()
  fam_table <- tibble(FID = c("deleteme"),
                     IID = c("deleteme"),
                     Related = c("deleteme"))
  while (nrow(ibdcoeff) > 0) {
    test_tab <- plyr::count(c(ibdcoeff$FI1, ibdcoeff$FI2))
    if (!("x" %in% names(test_tab))) {
      print(ibdcoeff)
    }
    sample.counts <- plyr::count(c(ibdcoeff$FI1, ibdcoeff$FI2)) %>%
      as_tibble %>%
      rename(FI = x) %>%
      mutate(FI = as.character(FI)) %>%
      inner_join(fam_fi, by = "FI") %>%
      arrange(desc(qc_failed), status, desc(freq))
    rm.sample <- sample.counts[[1, "FI"]]
    id_ <- str_split(rm.sample, "--tempsep--")[[1]]
    fid <- id_[1]
    iid <- id_[2]
    remtxt <- sprintf("%s %i other samples.",
                     msg,
                     sample.counts[[1, "freq"]])
    message(paste("Removing sample", iid, remtxt))
    ft <- tibble(FID = fid, IID = iid, Related = remtxt)
    fam_table <- fam_table %>%
      bind_rows(ft)
    ibdcoeff <- ibdcoeff[ibdcoeff$FI1 != rm.sample &

```

```

        ibdcoeff$FI2 != rm.sample, ]
    related_samples <- c(as.character(rm.sample), related_samples)
  }
  return(
    list(related_samples = related_samples,
         fam_table = filter(fam_table, Related != "deleteme"),
         exclude_samples = tibble(FI = as.character(related_samples)) %>%
           separate(FI, c("FID", "IID"), sep = "_-_-tempsep_-_-"))
  )
}

# Import data
fam <- "work/habshd_sampleQC.fam" %>%
  read_table(col_types = "cc---i", col_names = c("FID", "IID", "status")) %>%
  mutate(qc_failed = FALSE)

relatedness.raw = read_table("work/habshd_ibd.genome")

ibdcoeff <- relatedness.raw %>%
  filter(PI_HAT > pi_hat_thres) %>%
  mutate(
    pi_hat = closest(PI_HAT, rel_tab_filt$pi_hat),
    z0 = closest(Z0, rel_tab_filt$z0),
    z1 = closest(Z1, rel_tab_filt$z1),
    z2 = closest(Z2, rel_tab_filt$z2),
  ) %>%
  left_join(rel_tab_filt)

ibdcoeff_unrelated <- remove_samples(ibdcoeff, fam)

```

The following histogram (Fig. @ref(fig:kinplot)) shows the distribution of proportion of IBD sharing (pi-hat in PLINK; PropIBD in KING) between all pairs.

```

ggplot(relatedness, aes(x = PI_HAT)) +
  geom_histogram(binwidth = 0.01, fill = "#377EB8") +
  scale_y_continuous(trans = 'log10', breaks = scales::trans_breaks("log10", function(x) r
  coord_cartesian(xlim = c(min(relatedness$PI_HAT) - 0.05, 1)) +
  annotation_logticks() +
  theme_bw() +
  labs(x = "IBD Sharing (pi-hat in PLINK)") +
  geom_vline(xintercept = pi_hat_thres,

```



```
colour = "red", linetype = 2)
```

The following plot (Fig. @ref(fig:relplot)) shows the xxx by the proportion of loci where individuals share zero alleles (Z0), where the proportion of IBD sharing is greater than 0.05. In family based studies, pairs are colored by IBD relationship. Table @ref(tab:ibdfail) displays the individuals where the kinship coefficient was greater than xxx in population based studies OR how were duplicates in family based studies.

```
ggplot(ibdcoeff, aes(x = Z0, y = Z1, color = relationship)) +
  geom_point() +
  labs(x = 'P(IBD=0)', y = "P(IBD=0)") +
  theme_bw()
```

Population Substructure

After excluding population outliers from the dataset, population substructure will remain due to the presence of genetic diversity within apparently homogenous populations. Within a single ethnic population, even subtle degrees of population stratification can bias results due to differences in allele frequencies between subpopulations. Principal components based on the observed genotypes in the dataset of interest can be used to capture information on substructure and be included as covariates in downstream analysis.

To obtain the principal components for the sample dataset after population outliers have been removed, type the following PLINK 1.9 commands at the shell prompt to generate the principal component eigenvalue and eigenvector files.

```
plink \
  --bfile work/habshd_sampleQC \
  --extract work/indepSNP.prune.in \
  --pca 10 \
  --out work/habshd

# PCA file from plink

zscore = function(x){(x - mean(x, na.rm = TRUE)) / sd(x, na.rm = TRUE)}

# Read in eigenvectors and z-score transform
pca <- read_delim('work/habshd.eigenvec',
  delim = " ", col_names = c("FID", "IID", paste0("PC", 1:10)),
  col_types = cols(.default = "d", FID = "c", IID = "d")) %>%
```

```

      mutate_at(paste0("PC", 1:10), zscore) %>%
left_join(habshd %>% select(med_id, race), by = c('IID' = "med_id"))

# read in eigenvalues
eigenval.raw <- parse_number(read_lines('work/habshd.eigenval'))

eigenval <- tibble(eigenval = eigenval.raw,
                  PC = 1:length(eigenval.raw)) %>%
  mutate(PVE = round(eigenval / sum(eigenval), 3)) %>%
  select(PC, eigenval, PVE)

```

Scree Plot

The below scree plot (Fig. @ref(fig:ScreePlotStrat)) shows the amount of variation retained by each principal component (Left) and the cumulative proportion of variance explained by each principal component (Right).

```

#Include the number of PC for where the cumulative PVE is 95%
PC.inc <- findInterval(0.95, cumsum(eigenval$PVE)) + 1

## ---- Plot scree plot of proportion of variance explained by Principal components ---- #
p1 <- ggplot(data = eigenval, aes(x = PC, y = PVE, group = factor(1))) +
  geom_point(colour = '#377EB8') + geom_path(colour = '#377EB8') +
  labs(x = 'Principal Components') +
  theme_bw() + coord_cartesian(ylim = c(0,1), default = T)

p2 <- ggplot(data = eigenval, aes(x=PC, y=cumsum(PVE), group = factor(1))) +
  geom_point(colour = '#377EB8') + geom_path(colour = '#377EB8') +
  labs(x = 'Principal Components', y = 'cumulative PVE') +
  theme_bw() + coord_cartesian(ylim = c(0,1), default = T) +
  geom_hline(yintercept = 0.95, colour = '#E41A1C', linetype = 2)

gridExtra::grid.arrange(p1, p2, ncol = 2)

```

Population substructure

The following plots show the population structure of xxx based on the first two (Fig. @ref(fig:2PCstrat)) and three (Fig. @ref(fig:3PCstrat))) principal components compared with the reference populations from 1000 Genomes.

```
## Plot Superpopulations, PC1 + PC2
ggplot(pca, aes(x = PC2, y = PC1, color = race)) +
  geom_point() +
  theme_bw() + theme(legend.position = 'right')
```

Exclude Samples

```
bind_rows(
  sex_exclude.samples,
  het_exclude.samples,
  ibdcoeff_unrelated$exclude_samples %>% mutate_at(c('FID', 'IID'), as.numeric)
) %>%
  select(FID, IID) %>%
  distinct(FID, IID) %>%
  write_tsv('work/habshd.ExcludeSamples.tsv', col_names = F)
```

```
plink \
  --bfile work/habshd_sampleQC \
  --keep-allele-order \
  --remove work/habshd.ExcludeSamples.tsv \
  --make-bed --out work/habshd_gwas
```

GWAS

TBD

Import

```
fam <- read_table('work/habshd_gwas.fam', col_names = c('FID', 'IID'))
habshd <- read_csv("work/habshd_pheno.csv") %>% distinct(med_id, .keep_all = T)
pc <- pca <- read_delim('work/habshd.eigenvec',
                        delim = " ", col_names = c("FID", "IID", paste0("PC", 1:10)),
                        col_types = cols(.default = "d", FID = "c", IID = "d"))
```

Phenotype File

```
pheno <- fam %>%
  left_join(select(habshd, med_id, cdr_sum), by = c('IID' = 'med_id'))

pheno %>%
  write_tsv('work/habshd_gwas.pheno', col_names = F)
```

Covariate File

```
covar <- fam %>%
  left_join(select(habshd, med_id, age, id_gender), by = c('IID' = 'med_id')) %>%
  left_join(select(pc, IID, PC1, PC2, PC3, PC4), by = 'IID')

covar %>%
  write_tsv('work/habshd_gwas.covar', col_names = F)
```

GWAS

```
plink \  
  --bfile work/habshd_gwas \  
  --pheno work/habshd_gwas.pheno \  
  --covar work/habshd_gwas.covar \  
  --linear hide-covar \  
  --out results/habshd_cdr_gwas
```

Manhattan Plot

```
gwas.raw <- read_table('results/habshd_cdr_gwas.assoc.linear') %>%  
  select(-X10)  
  
ggman::ggman(gwas.raw, snp = "SNP", bp = "BP", chrom = "CHR", pvalue = "P", relative.posit  
  theme_bw()
```

GWAS-SS

Polygenic risk scores, Two-sample Mendelian Randomization, and Genetic Correlation methods require the use of summary statistics from genome-wide association studies, including single nucleotide polymorphisms (SNPs), beta coefficients, standard errors, p-values, and allele frequencies. However, the historical lack of standards for data content and file formats in GWAS summary statistics has resulted in heterogeneous data sets. To address this issue, standardizing and harmonizing the GWAS summary statistics is crucial before conducting MR analyses. The [GWAS Catalog](#) and [OpenGWAS](#) platforms have developed formats such as GWAS-SSF (Hayhurst et al. 2022) and GWAS-VCF (Lyon et al. 2021) to facilitate sharing of GWAS SumStats. Tools like [MungeSumstats](#) (Murphy et al 2021) and GWAS2VCF (Lyon et al. 2021) are available that provide rapid standardization and quality control of GWAS SumStats.

AD GWAS

We download the International Genomics of Alzheimer's Project (IGAP) Alzheimer's disease GWAS of Kunkle et al. Nat Genet, 2019. from the GWAS catalogue. These summary statistics correspond to the meta-analysis results obtained in stage 1 including genotyped and imputed data (11,480,632 variants, phase 1 integrated release 3, March 2012) of 21,982 Alzheimer's disease cases and 41,944 cognitively normal controls.

The Summary statistics consists of the following information for each SNP and its association to Alzheimer's disease based on meta-analysis in the publication mentioned below.

- Chromosome: Chromosome of the SNP (Build 37, Assembly Hg19)
- Position: Position of the SNP (Build 37, Assembly Hg19)
- MarkerName: SNP rsID or chromosome:position:I/D if rsID not available. I/D indicates indel or deletion respectively.
- Effect_allele: Reference allele (coded allele)
- Non_Effect_allele: Non reference allele (non coded allele)
- Beta: Overall estimated effect size for the effect allele
- SE: Overall standard error for effect size estimate
- Pvalue: Meta-analysis Pvalue using regression coefficients (beta and standard error)

```
curl https://ftp.ebi.ac.uk/pub/databases/gwas/summary_statistics/GCST007001-GCST008000/GCS
```

MungeSumstats

The MungeSumstats package is designed to facilitate the standardisation of GWAS summary statistics.

```
library(tidyverse)

-- Attaching core tidyverse packages ----- tidyverse 2.0.0 --
v dplyr      1.1.3      v readr      2.1.4
v forcats    1.0.0      v stringr    1.5.0
v ggplot2    3.4.4      v tibble     3.2.1
v lubridate  1.9.3      v tidyr      1.3.0
v purrr      1.0.2
-- Conflicts ----- tidyverse_conflicts() --
x dplyr::filter() masks stats::filter()
x dplyr::lag()     masks stats::lag()
i Use the conflicted package (<http://conflicted.r-lib.org/>) to force all conflicts to become
available

load.raw <- read_table('resources/Kunkle_etal_Stage1_results.txt')

load <- load.raw %>%
  filter(nchar(Effect_allele) == 1 & nchar(Non_Effect_allele) == 1) %>%
  mutate(
    Ncaas = 21982,
    Nctrl = 41944,
    N = 63926
  )

reformatted <-
  MungeSumstats::format_sumstats(path=load,
                                  ref_genome="GRCh37",
                                  dbSNP = 144,
                                  return_data = TRUE
                                ) %>%
  as_tibble()

write_tsv(reformatted, 'work/Kunkle2019load.tsv.gz')
```

Manhattan Plot

Genetic Ancestry

Prinicpal Component Analysis

```
#| eval: false
plink \
  --keep-allele-order \
  --bfile resources/genetic_epi/imputed_1kG_merged \
  --fam resources/genetic_epi/imputed_1kG_merged_fixed.fam \
  --pca 10 \
  --within resources/genetic_epi/1kG_pops.txt \
  --pca-clusters resources/genetic_epi/1kG_pops_unique.txt \
  --out work/imputed_1kG_merged
```

ADMIXTURE

```
#| eval: false
admixture -P -s 42 habshd_merged_gnomad-hgdp-1kg.hg38.bed 12 -j1
```


Part III

Polygenic Risk Scores

3 Polygenic Risk Scores

This is a book created from markdown and executable code.

See¹³ for additional discussion of literate programming.

Part IV

Mendelian Randomization

4 Mendelian Randomization

This is a book created from markdown and executable code.

See¹³ for additional discussion of literate programming.

References

1. Purcell S, Neale B, Todd-Brown K, et al. [PLINK: A Tool Set for Whole-Genome Association and Population-Based Linkage Analyses](#). *The American Journal of Human Genetics*. 2007;81(3):559-575.
2. Danecek P, Bonfield JK, Liddle J, et al. [Twelve years of SAMtools and BCFtools](#). *GigaScience*. 2021;10(2):giab008.
3. Murphy AE, Schilder BM, Skene NG. [MungeSumstats: A Bioconductor package for the standardisation and quality control of many GWAS summary statistics](#). *Bioinformatics*. 2021;37(23):btab665-.
4. Alexander DH, Novembre J, Lange K. [Fast model-based estimation of ancestry in unrelated individuals](#). *Genome Research*. 2009;19(9):1655-1664.
5. Maples BK, Gravel S, Kenny EE, Bustamante CD. [RFMix: A Discriminative Modeling Approach for Rapid and Robust Local-Ancestry Inference](#). *The American Journal of Human Genetics*. 2013;93(2):278-288.
6. Bulik-Sullivan B, Finucane HK, Anttila V, et al. [An atlas of genetic correlations across human diseases and traits](#). *Nature Genetics*. 2015;47(11):1236-1241.
7. Ning Z, Pawitan Y, Shen X. [High-definition likelihood inference of genetic correlations across human complex traits](#). *Nature Genetics*. 2020;52(8):859-864.
8. Grotzinger AD, Rhemtulla M, Vlaming R de, et al. [Genomic structural equation modelling provides insights into the multivariate genetic architecture of complex traits](#). *Nature human behaviour*. 2019;3(5):513-525.
9. Choi SW, O'Reilly PF. [PRSice-2: Polygenic Risk Score software for biobank-scale data](#). *GigaScience*. 2019;8(7).
10. Choi SW, García-González J, Ruan Y, et al. [PRSet: Pathway-based polygenic risk score analyses and software](#). *PLOS Genetics*. 2023;19(2):e1010624.
11. Ruan Y, Lin YF, Feng YCA, et al. [Improving polygenic prediction in ancestrally diverse populations](#). *Nature Genetics*. 2022;54(5):573-580.
12. Hemani G, Zheng J, Elsworth B, et al. [The MR-Base platform supports systematic causal inference across the human phenome](#). *eLife*. 2018;7:e34408.
13. Knuth DE. [Literate programming](#). *Comput J*. 1984;27(2):97-111.