

# Large-scale genomic investigation of the gut microbiome on Parkinson's disease etiology - Code

Rodrigo Sardon

11/27/2021

## Preprocessing GWAS Summary Datasets

The first step of this pipeline was to identify datasets we'd be performing our association studies on. Go to <https://mibiogen.gcc.rug.nl/> to find the repository of datasets provided by Kurilshikov et. al, 2021 (<https://pubmed.ncbi.nlm.nih.gov/33462485/>). On this website, you'll find 6 links:

1. MBG.allHits.p1e4.txt -> top hit variants ( $p < 1e-4$ ) for any level of the bacterial hierarchy
2. MiBioGen\_QmbQTL\_summary\_phylum.zip (2.4 GB) -> summary statistics of bacterial phyla
3. MiBioGen\_QmbQTL\_summary\_class.zip (4.4 GB) -> summary statistics of bacterial classes
4. MiBioGen\_QmbQTL\_summary\_order.zip (5.4 GB) -> summary statistics of bacterial orders
5. MiBioGen\_QmbQTL\_summary\_family.zip (9.5 GB) -> summary statistics of bacterial families
6. MiBioGen\_QmbQTL\_summary\_genus.zip (35.0 GB) -> summary statistics of bacterial genera

In our case, we wanted explore the bacterial genera taxonomy. In latter analyses, we hope to run family, order, class, and phylum taxonomies of bacteria. When you unzip the MiBioGen\_QmbQTL\_summary\_genus.zip, it will look something like this:

- MiBioGen\_QmbQTL\_summary\_genus:
  - genus.Clostridiuminnocuumgroup.id.14397.summary.txt.gz
  - genus.Eubacteriumbrachygroup.id.11296.summary.txt.gz

Unzip the subfolders and you get a table like this:

bac	chr	bp	rsID	ref.allele	eff.allele	beta	SE	Z weighted	P weighted	N	Cohorts
name	5	71186626	rs6890185	C	T	0.113	0.023	4.868	1.122e-06	4166	20

Now perform these data pre-processing steps (on sample data):

### Import “reticulate” to incorporate python functionalities.

```
knitr::opts_chunk$set(tidy.opts=list(width.cutoff=80), tidy=TRUE)

library(reticulate)
use_virtualenv("base")
use_python("/Volumes/T7Touch/Applications/anaconda3/bin/python")
```

Each .gz file of the sum stats will give a txt file, convert these into csv files.

```
knitr::opts_chunk$set(tidy.opts=list(width.cutoff=80), tidy=TRUE)

sumstats_root <- "/Users/rodrigosandon/Documents/GitHub/ParkisonsMicrobiomeStudy/MB-PD_Association_Study"

source_python("Utilities.py")
out_paths <- txt_to_csv_files_in_root(sumstats_root)
# If this doesn't work, run the py file itself with the process execution
```

Add a “p” column to each sum stat csv file.

```
source_python("Utilities.py")

knitr::opts_chunk$set(tidy.opts=list(width.cutoff=80), tidy=TRUE)

add_p_col_to_df <- function(bac_sumstat_path, new_col_name, out_path) {

  df <- read.csv(bac_sumstat_path, header = TRUE, sep = ",")
  df[new_col_name] <- 2*pnorm(-abs(df$beta/df$SE))
  eles_of_bac_path <- strsplit(bac_sumstat_path, "/")
  name_of_bac_sumstat <- eles_of_bac_path[[1]][length(eles_of_bac_path[[1]])]
  # ~finding the name original name of the bac sumstat file

  new_name <- paste("addedP_", name_of_bac_sumstat, sep = "")
  new_path <- paste(out_path, new_name, sep = "")
  #print(new_path)
  # Export
  write.csv(df, new_path, row.names = FALSE, quote = FALSE)

}

create_lst_of_file_paths <- function(root_path, files_endswith_str) {
  files <- list.files(path = root_path, pattern = files_endswith_str, full.names = TRUE,
                     recursive = FALSE)
  return (files)
}

add_p_col_to_csvs_in_root <- function(root_path, out_path, files_endswith_str,
                                     new_col_name, omit_csvs_that_contain) {
  dir.create(out_path, showWarnings = FALSE) #uncomment if dir hasn't been created yet, else, keep comm
  files <- create_lst_of_file_paths(root_path, files_endswith_str)
  # Reminder: current files in the list should be in csv format

  for (bac_sumstat_path in files) {
    print(paste("Working on ...", bac_sumstat_path))
    start.time <- Sys.time()
    if (grepl(omit_csvs_that_contain, bac_sumstat_path, fixed = TRUE) == FALSE) {
      add_p_col_to_df(bac_sumstat_path, new_col_name, out_path)
    }
    end.time <- Sys.time()
    print(paste("Time to process: ", end.time - start.time))
  }
}
```

```

ROOT_PATH <- "/Users/rodrigosandon/Documents/GitHub/ParkisonsMicrobiomeStudy/MB-PD_Association_Study_Pipe
OUT_PATH <- "/Users/rodrigosandon/Documents/GitHub/ParkisonsMicrobiomeStudy/MB-PD_Association_Study_Pipe

#perform this function only for the bacterial genera in which the names are known
add_p_col_to_csvs_in_root(ROOT_PATH, OUT_PATH, files_endswith_str = "*.csv",
                          new_col_name = "pDerived", omit_csvs_that_contain = "unknown")

```

## Polygenic Risk Score Analysis

Make a txt file of all of the genera you will be performing PRS on

```

#Identify the path where you've modified the sumstats
sumstats_path <- "/Users/rodrigosandon/Documents/GitHub/ParkisonsMicrobiomeStudy/MB-PD_Association_Study_Pipe

#Define the group of files to identify within ^ this root path
look_for <- "addedP_"

#Identify path where you'll locate the txt file
txt_out <- "/Users/rodrigosandon/Documents/GitHub/ParkisonsMicrobiomeStudy/MB-PD_Association_Study_Pipe

source_python("Utilities.py")

files_lst <- find_paths_startswith(sumstats_path, look_for)
listdir_to_txt_file(files_lst, txt_out)

```

Using PRSice.R, perform PRS and acquire the output file

```

### Polygenic risk score analyses of 119 bacteria genres versus PD risk

## Make a list of summary stats file
ls *csv > /Users/rodrigosandon/Documents/GitHub/ParkisonsMicrobiomeStudy/MB-PD_Association_Study_Pipe

## Format these files
cat genres.txt | while read LINE
do
  echo $LINE
  sed 's/\\/ /g' $LINE | sed 's/,/ /g' > temp.txt
  awk '{print $0"\t"$2:"$3}' temp.txt | sed 's/chr\:bp/ID/' > $LINE.temp_formatted.txt
  rm temp.txt
done

## Identify independent risk SNPs using our in-house LD reference data for European populations (/data/

Rscript /data/LNG/pdMeta5v2/leaveOneOutPrsice/PRSice_linux/PRSice.R --cov-file /data/LNG/saraB/WGS/noag
Rscript /data/LNG/pdMeta5v2/leaveOneOutPrsice/PRSice_linux/PRSice.R --cov-file /data/LNG/saraB/WGS/noag

## Remove NeuroX individuals & extract nominated variants
cat genres_formatted_list.txt | while read LINE
do

```

```

plink --bfile /data/LNG/saraB/concat_HARDCALLS_PD_september_2018_no_cousins --remove-fam NeuroX.txt --e
done

## Make score files
cat genuses_formatted_list.txt | while read LINE
do
awk '{print $14, $6, $7}' addedPgenus.$LINE.summary.txt.csv.temp_formatted.txt | sed '1d' > $LINE.toscore
done

## Make sure score files have 3 expected fields rather than 2
cat genuses_formatted_list.txt | while read LINE
do
grep ":" $LINE.toscore.txt > true_$LINE.toscore.txt
done

## Calculate scores
cat genuses_formatted_list.txt | while read LINE
do
plink --bfile pruned_$LINE --score $LINE.toscore.txt --make-bed --out pruned_$LINE
done

```

## Run PRS (logistic regression) in R

```

#install.packages("data.table")
library("data.table")
listOfProfiles <- read.table("genera_formatted_list.txt", header = T)
names(listOfProfiles) <- c("id")
covs1 <- fread("/data/LNG/saraB/WGS/noage_toPRsice_phenosAndCovs_renamed.tab", header = T)
covs2 <- fread("/data/LNG/saraB/concat_HARDCALLS_PD_september_2018_no_cousins.fam", header = F)
colnames(covs2) <- c("FID", "IID", "MAT", "PAT", "SEX", "PHENO")
covsfinal <- merge(covs1, covs2, by = "FID")
covsfinal$CASE <- covsfinal$PHENO.x - 1
outPut <- matrix(ncol = 4, nrow = length(listOfProfiles$id), NA)
colnames(outPut) <- c("genus", "b", "se", "p")
for(i in 1:length(listOfProfiles$id))
{
  profileName <- as.character(listOfProfiles$id[i])
  profile <- fread(file = paste(profileName, ".profile", sep = ""), header = T)
  profile$index <- paste(profile$FID, profile$IID, sep = "")
  data <- merge(covsfinal, profile, by = "index")
  meanControls <- mean(data$SCORE[data$CASE == 0])
  sdControls <- sd(data$SCORE[data$CASE == 0])
  data$zSCORE <- (data$SCORE - meanControls)/sdControls
  grsTest <- glm(CASE ~ zSCORE + SEX + PC1 + PC2 + PC3 + PC4 + PC5 + PC6 + PC7 + PC8 + PC9 + PC10 + D
  beta <- summary(grsTest)$coefficients["zSCORE", "Estimate"]
  se <- summary(grsTest)$coefficients["zSCORE", "Std. Error"]
  p <- summary(grsTest)$coefficients["zSCORE", "Pr(>|z|)"]
  outPut[i,1] <- profileName
  outPut[i,2] <- beta
  outPut[i,3] <- se
  outPut[i,4] <- p
}
write.table(outPut, "Genus_PRS.tab", quote = F, sep = "\t", row.names = F)

```

```
## SAMPLE RESULT:

# Call:
# glm(formula = CASE ~ zSCORE + SEX + PC1 + PC2 + PC3 + PC4 + PC5 +
#      PC6 + PC7 + PC8 + PC9 + PC10, family = "binomial", data = data)
#
# Deviance Residuals:
#      Min       1Q   Median       3Q      Max
# -1.919  -1.003  -0.810   1.278   1.796
#
# Coefficients:
#              Estimate Std. Error z value Pr(>|z|)
# (Intercept)   0.43040    0.03974  10.831 < 2e-16 ***
# zSCORE        0.03021    0.01424   2.121  0.0339 *
# SEX          -0.58575    0.02603 -22.505 < 2e-16 ***
# PC1          -35.38600    2.73911 -12.919 < 2e-16 ***
# PC2           50.17593    2.86877  17.490 < 2e-16 ***
# PC3           10.63757    2.68575   3.961 7.47e-05 ***
# PC4           0.63048    2.65991   0.237  0.8126
# PC5           16.19265    2.70187   5.993 2.06e-09 ***
# PC6          -24.20283    2.74525  -8.816 < 2e-16 ***
# PC7           1.61607    2.62627   0.615  0.5383
# PC8           12.66905    2.70987   4.675 2.94e-06 ***
# PC9           -5.96044    2.64009  -2.258  0.0240 *
# PC10          -16.18338    2.65955  -6.085 1.16e-09 ***
# ---
# Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
#
# (Dispersion parameter for binomial family taken to be 1)
#
#      Null deviance: 35275  on 26385  degrees of freedom
# Residual deviance: 34124  on 26373  degrees of freedom
# AIC: 34150
#
# Number of Fisher Scoring iterations: 4
```

## Linkage Disequilibrium Score Regression (LDSC) Analysis

### Data Preprocessing

```
pre_ldsc_chrPosRs <- source_python("/Users/rodrigossandon/Documents/GitHub/ParkisonsMicrobiomeStudy/MB-P/
# chrPosRs.tab available upon request
# reformatted2_META5_all_with_rsid.txt available upon request
#
```

### LDSC

```
"""
Note: Make sure you are running in a Python 2 environment.
"""
import pandas as pd
```

```

import csv
import os
import time
import os.path
from os import path

# the PD file for LD SR run is : /Volumes/T7Touch/NIHSummer2021/Data/LDSR_analysis/nalls_onlyRsIDs.sums

# Ex: /Volumes/Passport/119Bacs_addedP/addedPgenus..Clostridiuminnocuumgroup.id.14397.summary.txt.csv

def formatSumStatsForBac(csvPath):
    newFileName = "formatted." + csvPath.split("/") [4].replace(
        ".txt", "").replace("..", ".").replace("addedPgenus", "addedP.genus")
    newFilePath = "/Volumes/Passport/formatted119Bacs_addedP/%s" % (
        newFileName)

    if path.exists(newFilePath) == False:
        df = pd.read_csv(csvPath)
        newDf = df[['rsID', 'eff.allele', 'ref.allele',
                    'beta', 'SE', 'N', 'pDerived']].copy()
        newDf.insert(3, 'Zscore', newDf['beta'] / newDf['SE'], True)
        newDf = newDf.drop(['beta', 'SE'], axis=1)
        newDf = newDf.rename(columns={
            'rsID': 'snpid', 'pDerived': 'P-value', 'eff.allele': 'A1', 'ref.allele':

            # Now formatted.addedP.genus.Clostridiuminnocuumgroup.id.14397.summary.csv

            newDf.to_csv(newFilePath, index=False)
    else:
        print("File %s already exists" % (newFilePath))
    return newFilePath

def CSVtoTXT(csv_file, txtFile):
    with open(txtFile, "w") as my_output_file:
        with open(csv_file, "r") as my_input_file:
            [my_output_file.write(" ".join(row)+'\n')
             for row in csv.reader(my_input_file)]
        my_output_file.close()

# ex csv file now: /Volumes/Passport/formatted119Bacs_addedP/formatted.addedP.genus.Clostridiuminnocuum

def mungeDataCall(csv_file):

    txtFilePath = csv_file.replace(".csv", ".txt")
    # /Volumes/Passport/formatted119Bacs_addedP/formatted.addedP.genus.Clostridiuminnocuumgroup.id.1439

    newFilePath = txtFilePath.replace(
        "formatted119Bacs_addedP", "munge119Bacs_output").replace(".txt", "")

    if path.exists(txtFilePath) == False:

```

```

        CSVtoTXT(csv_file, txtFilePath)
    else:
        print("File %s already munged" % (newFilePath))
        # os.chdir("/Users/rodrigossandon/ldsc")

    cmd = "./munge_sumstats.py \
        --sumstats %s \
        --out %s \
        --merge-alleles /Volumes/T7Touch/NIHSummer2021/Data/LDSR_analysis/ldsc/w_hm3.snplist" % (txtFil

    if path.exists(newFilePath + ".sumstats.gz") == False: # only if .gz file don't exist
        os.system(cmd)

    return newFilePath

# ex gz bac name munged: /Volumes/Passport/munge119Bacs_output/formatted.addedP.genus.Clostridiuminnocu

def LDscore_regression(munged_bac_output):
    print("munged_bac_output: ", munged_bac_output)
    LDSR_outName = "/Volumes/T7Touch/NIHSummer2021/Code/LDSC_analysis2/results/%s_ldscResults" % (
        munged_bac_output.split("/") [4].split(".")[3]) # <--only bac name

    # os.chdir("/Users/rodrigossandon/ldsc")

    cmd = "./ldsc.py \
        --rg %s,/Volumes/T7Touch/NIHSummer2021/Code/LDSC_analysis2/munged_META5_all_with_rsid.sumstats.
        --ref-ld-chr /Volumes/T7Touch/NIHSummer2021/Data/LDSR_analysis/ldsc/eur_w_ld_chr/ \
        --w-ld-chr /Volumes/T7Touch/NIHSummer2021/Data/LDSR_analysis/ldsc/eur_w_ld_chr/ \
        --out %s " % (munged_bac_output + ".sumstats.gz", LDSR_outName)

    os.system(cmd)

    return LDSR_outName

###MAIN###

masterDir = "/Volumes/Passport/119Bacs_addedP/"

for root, dirs, files in os.walk(masterDir):
    for name in files:
        start = time.time()
        bacPathToProcess = os.path.join(root, name)
        print("Processing", bacPathToProcess)

        newFilePath1 = formatSumStatsForBac(bacPathToProcess)

        newFilePath2 = mungeDataCall(newFilePath1)

        LDSR_outName = LDscore_regression(newFilePath2)

```

```

print("Results can be found in", LDSR_outName)
end = time.time()
print("Time to perform LDSR:", (end-start)/60, "mins")

# BEFORE RUNNING FILE, MUNGE NEW PD GWAS w/ w_hm3 snplist
cmd = "./munge_sumstats.py \
      --sumstats /Volumes/T7Touch/NIHSummer2021/Code/LDSC_analysis2/reformatted2_META5_all_with_rsid.
      --out /Volumes/T7Touch/NIHSummer2021/Code/LDSC_analysis2/munged_META5_all_with_rsid \
      --merge-alleles /Volumes/T7Touch/NIHSummer2021/Data/LDSR_analysis/ldsc/w_hm3.snplist"

```

## Bayesian Colocalization

```

# addedPgenus.CandidatusSoleaferrea.id.11350.summary.txt.csv on request
# META5_all_with_rsid.txt upon request
# nallsEtAl2019_excluding23andMe_allVariants.tab from Nalls et. al 2019.

library("data.table")
install.packages("coloc")
source("https://bioconductor.org/biocLite.R")
biocLite("snpStats")
install.packages("robustbase")
library("robustbase")
library("coloc")
library("tidyverse")

#The file that contains SNPs with significance of 5E-8 for both disease and phenotype
masterDir <- "/Volumes/T7Touch/NIHSummer2021/Code/Colocalization/p_smr_multi_mrQTL_belowFDR"
files <- list.files(path=masterDir, pattern="*.tab", full.names = TRUE, recursive = FALSE)

#Run this same pipeline for other bacterial taxa
Candida_instruments <- fread(file="/Volumes/Passport/119Bacs_addedP/addedPgenus.CandidatusSoleaferrea.i

for (k in files) {
  TOCOLOC <- read.table(k, header = T, sep="\t")
  nameOfTest <- str_replace(str_split(k, "/", n = Inf, simplify = TRUE)[[8]], ".tab", "")

  #Just need a subset of file that contains SNPs with significance of 5E-8 for both disease and phenotype
  SNPlist <- subset(TOCOLOC, select=c(topSNP, topSNP_bp, topSNP_chr))

  #For each SNP, search through phenotype summary stats to get 1MB of variants b4 and after SNP, this i
  for(i in 1:length(SNPlist$topSNP)) {
    thisSNP <- SNPlist$topSNP[i]
    thisChr <- SNPlist$topSNP_chr[i]
    thisBP <- SNPlist$topSNP_bp[i]
    thisBpLow <- thisBP - 1000000
    thisBpHigh <- thisBP + 1000000
    getRegions <- subset(Candida_instruments, chr == thisChr & bp > thisBpLow & bp < thisBpHigh)
    print(paste("Count is:", i))
    print(paste("Size of region for ",thisSNP," in ",thisChr,"_",thisBP, " : ", nrow(getRegions), sep="
    dir.create(paste("/Volumes/T7Touch/NIHSummer2021/Code/Colocalization/coloc_mrQTL_regions/",nameOfTest
    write.table(getRegions, paste("/Volumes/T7Touch/NIHSummer2021/Code/Colocalization/coloc_mrQTL_region
  }

```



```

}

#(PD vs Candida)
##Bayesian colocalisation analysis

#####Excluding some columns and changing name of columns of region files#####
library("dplyr")
masterDir <- "/Volumes/T7Touch/NIHSummer2021/Code/Colocolization/regionsV2"
files <- list.files(path=masterDir, pattern="*.tab", full.names = TRUE, recursive = FALSE)

for (i in files) {
  #/Volumes/T7Touch/NIHSummer2021/Code/Colocolization/regions/rs489567_region.tab
  data <- fread(file=i,header=T,sep="\t")
  df <- select(data, `chr:pos`, rsID, eff.allele, ref.allele, `beta.x`, SE, N, Ncohorts, pDerived)
  names(df) <- c("SNP", "rsID", "A1", "A2", "beta", "SE", "N", "Ncohorts", "P-value")
  write.table(df, i, quote = F, sep = "\t", row.names = F)
}

#####reading one type of Dx sum stats and renaming col so that merging is possible#####

dataPD <- fread(file="/Volumes/T7Touch/NIHSummer2021/Code/Colocolization/META5_all_with_rsid.txt",header=T,sep="\t")
names(dataPD)[17] <- "rsID"
write.table(dataPD, "/Volumes/T7Touch/NIHSummer2021/Code/Colocolization/META5_all_with_rsid.tab", quote=F, sep="\t", row.names=F)

#####reading one type of Dx sum stats but no change of col bc rsID column name already exists for nalls_PD
nalls_PD <- fread("/Volumes/T7Touch/NIHSummer2021/Code/Colocolization/nallsEtA12019_excluding23andMe_all_regions.tab", header=T, sep="\t")
#####
nalls_PD$SNP <- gsub("chr","",nalls_PD$SNP) #substituting a pattern of characters in a specific column in a data frame

#####Merging all info on candida data with even more info from the Dx sum stats by rsID#####

masterDir <- "/Volumes/T7Touch/NIHSummer2021/Code/Colocolization/regionsFullStats"
files <- list.files(path=masterDir, pattern="*.tab", full.names = TRUE, recursive = FALSE)

for (i in files) {
  print(paste("Working on...",i))
  data <- fread(file=i,header=T,sep="\t")
  print("here")
  df <- merge(data, dataPD, by="rsID")
  print("here")
  write.table(df, i, quote = F, sep = "\t", row.names = F)
}

#####Running the colocalization#####

```

```

#We subset data frames and rename them and let the function coloc.abf do all the work
SNPlist <- fread(file = paste("/Volumes/T7Touch/NIHSummer2021/Code/Colocolization/hitsBelow5E-5ForPD/hi

regions <- c(SNPlist$ID)

for (i in regions) {
  #/Volumes/T7Touch/NIHSummer2021/Code/Colocolization/regions
  thisRegion <- i
  print(i)
  data <- fread(file = paste("/Volumes/T7Touch/NIHSummer2021/Code/Colocolization/regionsFullStats/",this
  #print(names(data))
  df2_Candida_coloc <- data[,c("SNP","beta","P-value.x")]
  df1_PD_coloc <- data[,c("SNP","b","p","Freq1")]
  df2_Candida_coloc$var <- var(df2_Candida_coloc$beta)
  df1_PD_coloc$var <- var(df1_PD_coloc$b)
  names(df1_PD_coloc) <- c("snp","beta","pvalues","MAF","varbeta")
  names(df2_Candida_coloc) <- c("snp","beta","pvalues","varbeta")

  #Remove duplicates in dataset 1
  df1_PD_coloc <- df1_PD_coloc[!duplicated(df1_PD_coloc$snp), ]
  df2_Candida_coloc <- df2_Candida_coloc[!duplicated(df2_Candida_coloc$snp), ]

  df1_PD_coloc$type <- "cc"
  df2_Candida_coloc$type <- "quant"
  df1_PD_coloc$s <- 0.0590704464
  df2_Candida_coloc$sdY <- 1
  results <- coloc.abf(df1_PD_coloc, df2_Candida_coloc, MAF = NULL, p1 = 1e-02, p2 = 1e-02, p12 = 1e-02)
  #results <- coloc.abf(df1_PD_coloc, df2_Candida_coloc, MAF = NULL, p1 = 1e-04, p2 = 1e-04, p12 = 1e-04)
  sink(file = paste("/Volumes/T7Touch/NIHSummer2021/Code/Colocolization/PD_coloc/",thisRegion,"_results.txt"),
  print(results$summary)
  sink()
  write.table(results$results, file = paste("/Volumes/T7Touch/NIHSummer2021/Code/Colocolization/PD_coloc/",thisRegion,"_results.txt"),
}

#####
#####POST COLOCALIZATION: FINDING COLOCALIZED REGIONS THAT EXPLAIN PHENOTYPE AND DX ETIOLOGIES VIA SIMILARITY#####
#####

masterDir <- "/Volumes/T7Touch/NIHSummer2021/Code/Colocolization/PD_coloc"
files <- list.files(path=masterDir, pattern="*Summary.txt", full.names = TRUE, recursive = FALSE)

for (i in files) {
  #print(paste("Working on...",i))
  data <- fread(file=i,header=T,sep=" ")
  #print(paste("The PP.H4.abf for this region is:",data$`PP.H4.abf`))
  print(as.double(data$PP.H4.abf))
  if (as.double(data$PP.H0.abf) > 0.95) {
    print("PP.H0.abf HIT")
  }
  if (as.double(data$PP.H1.abf) > 0.95) {
    print("PP.H1.abf HIT")
  }
  if (as.double(data$PP.H2.abf) > 0.95) {
    print("PP.H2.abf HIT")
  }
}

```

```

}
if (as.double(data$PP.H3.abf) > 0.95) {
  print("PP.H3.abf HIT")
}
if (as.double(data$PP.H4.abf) > 0.95) {
  print("PP.H4.abf HIT")
}
}

```

## Generalised Summary-data-based Mendelian Randomisation

### Analysis for each bacterial taxa hit

```

library("gsmr")
library("dplyr")

allo_sumstats <- fread("/Volumes/T7Touch/NIHSummer2021/Code/GSMR/ScoreFilesButMergedWSumstats/Alloprevotella.id.")
cand_sumstats <- fread("/Volumes/T7Touch/NIHSummer2021/Code/GSMR/ScoreFilesButMergedWSumstats/CandidatusSoleaferus.id.")

meta5_pd <- fread("/Volumes/T7Touch/NIHSummer2021/Code/Colocolization/META5_all_with_rsid.tab", header=T)

merge1 <- merge(allo_sumstats, meta5_pd, by="rsID")
hits_allo <- merge1[merge1$`P-value` < 0.0005, ]
hits_allo2 <- hits_allo[hits_allo$`p_bac` < 0.0005, ]
write.table(hits_allo2, file = "/Volumes/T7Touch/NIHSummer2021/Code/GSMR/ScoreFilesButMergedWSumstats/Alloprevotella_hits.txt", as.is=T)

merge2 <- merge(cand_sumstats, meta5_pd, by="rsID")
hits_cand <- merge2[merge2$`P-value` < 0.0005, ]
hits_cand2 <- hits_cand[hits_cand$`p_bac` < 0.0005, ]
write.table(hits_cand2, file = "/Volumes/T7Touch/NIHSummer2021/Code/GSMR/ScoreFilesButMergedWSumstats/CandidatusSoleaferus_hits.txt", as.is=T)

allo <- fread("/Volumes/T7Touch/NIHSummer2021/Code/GSMR/ScoreFilesButMergedWSumstats/Alloprevotella.id.")
cand <- fread("/Volumes/T7Touch/NIHSummer2021/Code/GSMR/ScoreFilesButMergedWSumstats/CandidatusSoleaferus.id.")

allo_sub <- allo[,c("rsID", "eff.allele", "ref.allele", "beta.x", "SE", "p_bac", "N", "Effect", "StdErr")]
cand_sub <- cand[,c("rsID", "eff.allele", "ref.allele", "beta.x", "SE", "p_bac", "N", "Effect", "StdErr")]

allo_sub <- setNames(allo_sub, c("SNP", "a1", "a2", "bzx", "bzx_se", "bzx_pval", "bzx_n", "bzy", "bzy_se", "bzy_n"))
cand_sub <- setNames(cand_sub, c("SNP", "a1", "a2", "bzx", "bzx_se", "bzx_pval", "bzx_n", "bzy", "bzy_se", "bzy_n"))

write.table(allo_sub, file = "/Volumes/T7Touch/NIHSummer2021/Code/GSMR/ScoreFilesButMergedWSumstats/Alloprevotella_sub.txt", as.is=T)
write.table(cand_sub, file = "/Volumes/T7Touch/NIHSummer2021/Code/GSMR/ScoreFilesButMergedWSumstats/CandidatusSoleaferus_sub.txt", as.is=T)

# Save the genetic variants and effect alleles in a text file using R
write.table(allo_sub[,c(1,2)], "/Volumes/T7Touch/NIHSummer2021/Code/GSMR/allo_sub.allele", col.names=F, as.is=T)
write.table(cand_sub[,c(1,2)], "/Volumes/T7Touch/NIHSummer2021/Code/GSMR/cand_sub.allele", col.names=F, as.is=T)

# Extract the genotype data from a GWAS dataset using GCTA
#gcta64 --bfile gsmr_example --extract gsmr_example_snps.allele --update-ref-allele gsmr_example_snps.allele --r2 0.2 --output gsmr_example_snps.gcta
#ALLOPREVOTELLA

```

```

bzx = allo_sub$std_bzx      # SNP effects on the risk factor
bzx_se = allo_sub$std_bzx_se # standard errors of bzx
bzx_pval = allo_sub$bzx_pval # p-values for bzx
bzy = allo_sub$bzy          # SNP effects on the disease
bzy_se = allo_sub$bzy_se    # standard errors of bzy
bzy_pval = allo_sub$bzy_pval # p-values for bzy
n_ref = 7703                # Sample size of the reference sample
gwas_thresh = 5e-8          # GWAS threshold to select SNPs as the instruments for the GSMR analysis
single_snp_heidi_thresh = 0.01 # p-value threshold for single-SNP-based HEIDI-outlier analysis
multi_snp_heidi_thresh = 0.01 # p-value threshold for multi-SNP-based HEIDI-outlier analysis
nsnps_thresh = 10           # the minimum number of instruments required for the GSMR analysis
heidi_outlier_flag = T       # flag for HEIDI-outlier analysis
ld_r2_thresh = 0.05          # LD r2 threshold to remove SNPs in high LD
ld_fdr_thresh = 0.05         # FDR threshold to remove the chance correlations between the SNP instruments
gsmr2_beta = 0               # 0 - the original HEIDI-outlier method; 1 - the new HEIDI-outlier method that is cu
gsmr_results = gsmr(bzx, bzx_se, bzx_pval, bzy, bzy_se, bzy_pval, ldrho, snp_coeff_id, n_ref, heidi_out
filtered_index=gsmr_results$used_index
cat("The estimated effect of the exposure on outcome: ",gsmr_results$bxy)

##Getting formatted summ stats

bac1 <- fread("/Volumes/T7Touch/NIHSummer2021/Code/GSMR/bacSumStats_GCTA-COJO_format/Alloprevotella.id.
bac2 <- fread("/Volumes/T7Touch/NIHSummer2021/Code/GSMR/bacSumStats_GCTA-COJO_format/CandidatusSoleafer

bac1_sub <- bac1[,c("rsID","eff.allele","ref.allele","beta.x","SE","p_bac","N")]
bac2_sub <- bac2[,c("rsID","eff.allele","ref.allele","beta.x","SE","p_bac","N")]

bac1_sub <- setNames(bac1_sub, c("SNP", "A1", "A2", "b", "se", "p", "N"))
bac2_sub <- setNames(bac2_sub, c("SNP", "A1", "A2", "b", "se", "p", "N"))

write.table(bac1_sub, "/Volumes/T7Touch/NIHSummer2021/Code/GSMR/bacSumStats_GCTA-COJO_format/Alloprevot
write.table(bac2_sub, "/Volumes/T7Touch/NIHSummer2021/Code/GSMR/bacSumStats_GCTA-COJO_format/Candidatus

```

## Merging SNPs for all bacterial taxa and disease

```

library("dplyr")
library("data.table")
library("stringr")
#install.packages("jbb")
library("jbb")
#####

bac_mrqtlfolder_paths <- c()
bac_mrqtlfolder_paths <-c(bac_mrqtlfolder_paths,"/Volumes/T7Touch/NIHSummer2021/Data/MR-QTL_bacteria/SM
bac_mrqtlfolder_paths <-c(bac_mrqtlfolder_paths,"/Volumes/T7Touch/NIHSummer2021/Data/MR-QTL_bacteria/SM

pd_mrqtlfolder_path_root_path <- "/Volumes/T7Touch/NIHSummer2021/Data/MR-QTL_bacteria/MR-QTL/"

results_rootpath <- "/Volumes/T7Touch/NIHSummer2021/Data/MR-QTL_bacteria/results3/"

dir.create(results_rootpath)

```

```

bacsMRQTLS_pdMRQTLS_mergeFDR <- function(bac_mrqtlfolder_paths,pd_mrqt_root_path,results_rootpath) {
  for (i in bac_mrqtlfolder_paths){
    bac_mrqtls <- list.files(path=bac_mrqtlfolder_paths, pattern="*.msmr", full.names = TRUE, recursive=TRUE)
    nameBac <- str_split(str_split(i, "/", n = Inf, simplify = TRUE)[[7]], "_", n = Inf, simplify = TRUE)[[1]]
    #print(paste("Working on",nameBac))

    dir.create(paste(results_rootpath,nameBac,"_PDmerged",sep=""))
    for (k in bac_mrqtls) { #going through each .msmr file
      #example name: /Volumes/T7Touch/NIHSummer2021/Data/MR-QTL_bacteria/SMR_results_CS/CS_chromatin_blood_Bryo
      subMRQTL_name <- str_replace(str_split(k, "/", n = Inf, simplify = TRUE)[[8]], paste(nameBac,"_",sep=""),
      paste(nameBac,"_PDmerged",sep=""))

      pdMRQTL_filepath <- paste(pd_mrqt_root_path,"PD_",subMRQTL_name,sep="") #so i dont need to
      #print(paste("MRQTL focused on:",pdMRQTL_filepath))

      #possibleError <- tryCatch(fread(pdMRQTL_filepath, header=T, sep="\t"), error=function(e) {next})
      if (file.exists(pdMRQTL_filepath) == T) {
        #REAL WORK
        bac_mrqt <- fread(k, header=T, sep="\t")
        pd_mrqt <- fread(pdMRQTL_filepath, header=T, sep="\t") #now have both csvs opened

        #/Volumes/T7Touch/NIHSummer2021/Data/MR-QTL_bacteria/MR-QTL/PD_chromatin_blood_Bryo_SMR_allCh
        nameofMRQTL <- str_replace(paste(nameBac,"_",subMRQTL_name,sep=""),".msmr","")

        newdf <- merge(bac_mrqt,pd_mrqt,by="topSNP") #merging w/ mrqt
        #new dir for just this mrqt cuz will have p filtered and p unfiltered results

        dir.create(paste(results_rootpath,nameBac,"_PDmerged","/",nameofMRQTL,sep=""))

        write.table(newdf,paste(results_rootpath,nameBac,"_PDmerged/",nameofMRQTL,"/", "FDR_",nameofMRQTL),as.is=T,
        sep=";",col.names=c("topSNP","p_SMR_multi","p_SMR_multi.y","p_SMR_multi.x","p_SMR_multi.y.y","p_SMR_multi.y.x","p_SMR_multi.x.x"),
        row.names=FALSE)

        #do some p_SMR filtering
        #newdf$p_fdr_bac <- p.adjust(newdf$p_SMR_multi.x, method = "fdr")
        #newdf$p_fdr_pd <- p.adjust(newdf$p_SMR_multi.y, method = "fdr")
        #hits1 <- newdf[newdf$p_fdr_bac < 0.05, ]
        #hits2 <- hits1[hits1$p_fdr_pd < 0.05, ]
        #alpha <- 0.05
        #hits1 <- newdf[newdf$p_SMR_multi.x < alpha, ]
        #hits2 <- hits1[hits1$p_SMR_multi.y < alpha, ]

        print(paste(nameofMRQTL,"N Rows for MRQTL after p filtering:",nrow(hits2)))
        if (nrow(hits2) > 0) {
          print("#####HIT#####")
        }

        write.table(hits2,paste(results_rootpath,nameBac,"_PDmerged/",nameofMRQTL,"/",toString(alpha),".csv"),as.is=T,
        sep=";",col.names=c("topSNP","p_SMR_multi","p_SMR_multi.y","p_SMR_multi.x","p_SMR_multi.y.y","p_SMR_multi.y.x","p_SMR_multi.x.x"),
        row.names=FALSE)
      }
    }
  }
}
bacsMRQTLS_pdMRQTLS_mergeFDR(bac_mrqtlfolder_paths,pd_mrqt_root_path,results_rootpath)

```

## Including Plots

You can also embed plots, for example:

Note that the `echo = FALSE` parameter was added to the code chunk to prevent printing of the R code that generated the plot.