

A real example for performing GWAS summary-level data based MR analysis with MRAPSS package

HU Xianghong

5/8/2020

Introduction

MR-APSS is a unified approach for Mendelian randomization accounting for pleiotropy, sample overlap and selection bias using genome-wide summary statistics. The MRAPSS package implement the MR-APSS approach to test for the causal effects between an exposure and a outcome disease.

We illustrate how to analyze GWAS summary level data using the MRAPSS software by an real example, i.e. BMI (exposure) and T2D (outcome). The MRAPSS analysis comprises five steps:

- Step 1: Download GWAS summary-level data from public resources
- Step 2: Format data
- Step 3: Harmonise datasets and estimate nuisance parameters
- Step 4: IVs selection and LD clumping
- Step 5: Fit MRAPSS

Step 3 requires the LD score files at [link](#). You can also use the LD scores calculated by yourself.

Step 0: Installtion and load packages

```
#install.packages("devtools")
devtools::install_github("YangLabHKUST/MRAPSS")
```

```
library(MRAPSS)
library(readr)
```

Step 1: Download GWAS summary-level data from public resources

To begin, set your working directory use `setwd()`.

Download the GWAS summary statistics at links for [BMI](#)(Watanabe et.al. (2019) [PMID: 25673413]) and [T2D](#)(Xue et al. (2015) [PMID: 31427789]).

You can also download the datasets [here](#)

Uncompress and rename the files as “BMI_ukb.txt” and “T2D.txt”, then read the datasets into R:

```
BMI_raw <- readr::read_delim("BMI_ukb.txt", "\t", escape_double = FALSE,
                             trim_ws = TRUE, progress = F)
#> Parsed with column specification:
#> cols(
#>   .default = col_double(),
#>   SNP = col_character(),
#>   A1 = col_character(),
```

```

#> TEST = col_character(),
#> A2 = col_character(),
#> SNPID_UKB = col_character(),
#> A1_UKB = col_character(),
#> A2_UKB = col_character()
#> )
#> See spec(...) for full column specifications.

T2D_raw <- readr::read_delim("T2D.txt", " ",
                           escape_double = FALSE,
                           trim_ws = TRUE, progress = F)
#> Parsed with column specification:
#> cols(
#>   CHR = col_double(),
#>   BP = col_double(),
#>   SNP = col_character(),
#>   A1 = col_character(),
#>   A2 = col_character(),
#>   frq_A1 = col_double(),
#>   b = col_double(),
#>   se = col_double(),
#>   P = col_double(),
#>   N = col_double()
#> )

```

Step 2: Format summary statistics

Format the summary-level data to have the following columns by `format_data()`:

- SNP: rs number
- A1: effect allele
- A2: the other allele
- Z: Z score
- chi2: χ^2 statistics
- P: pvalue
- N: sample size

```

BMI = format_data(BMI_raw,
                  snp_col = "SNPID_UKB",
                  b_col = "BETA",
                  se_col = "SE",
                  freq_col = "MAF_UKB",
                  A1_col = "A1",
                  A2_col = "A2",
                  p_col = "P",
                  n_col = "NMISS",
                  info_col = "INFO_UKB")
#> Begin formatting ....
#> The raw dataset has 10599054 dat lines.
#> Remove SNPs in MHC region ...
#> Merge SNPs with the hapmap3 snplist ...

```

```

#> Remove SNPs with imputation info less than 0.9 ...
#> Remove ambiguous SNPs ...
#> Remove SNPs with alleles not matched with the hapmap3 snplist
#> Remove SNPs with p value < 0 or p value > 1
#> Infer Z score from b/se ...
#> Remove SNPs with chi2 > chi2_max ...
#> The formatted data has 1036975 dat lines.

T2D = format_data(T2D_raw,
                  snp_col = "SNP",
                  b_col = "b",
                  se_col = "se",
                  freq_col = "frq_A1",
                  A1_col = "A1",
                  A2_col = "A2",
                  p_col = "P",
                  n_col = "N")

#> Begin formatting ....
#> The raw dataset has 5053015 dat lines.
#> Remove SNPs in MHC region ...
#> Merge SNPs with the hapmap3 snplist ...
#> Remove SNPs with p value < 0 or p value > 1
#> Infer Z score from b/se ...
#> Remove SNPs with chi2 > chi2_max ...
#> The formatted data has 936381 dat lines.

```

Have a look at the formatted datasets:

```

head(BMI)
#>      SNP A1 A2      Z      N      chi2      P
#> 1 rs1000000 A  G -0.2879593 385270 0.08292058 0.77340
#> 2 rs10000010 C  T -0.7299169 376552 0.53277868 0.46540
#> 3 rs1000002 T  C -1.9962652 385336 3.98507464 0.04592
#> 4 rs10000023 G  T  2.5662432 371893 6.58560413 0.01028
#> 5 rs1000003 G  A -1.5807309 383709 2.49871017 0.11380
#> 6 rs10000033 C  T  1.4280409 382629 2.03930068 0.15310

head(T2D)
#>      SNP A1 A2      Z      N      chi2      P
#> 1 rs1000000 A  G -0.2857143 573633 0.08163265 0.77320
#> 2 rs10000010 C  T  0.4931507 587226 0.24319760 0.61580
#> 3 rs1000002 T  C -2.5512821 579347 6.50904011 0.01035
#> 4 rs10000023 G  T  0.3000000 557936 0.09000000 0.76390
#> 5 rs1000003 G  A -1.6388889 577450 2.68595679 0.10140
#> 6 rs10000033 C  T  1.9230769 576206 3.69822485 0.05519

```

Note: `format_data()` will try to interpretate the raw datasets with user specified column names, for example specify “SNP” in “BMI.txt” as “SNP_col”. At the same time, it will also conduct the following quality control procedures.

- extract SNPs in the set of HapMap 3 list with minor allele frequency > 0.05.
- remove SNPs with alleles not in (G,C,T,A).
- remove SNPs with ambiguous alleles (G/C or A/T) or other false alleles (A/A T/T,G/G or C/C).
- exclud SNPs in the complex Major Histocompatibility Region (Chromosome 6, 26Mb-34Mb).

- remove SNPs with $\chi^2 > \text{chi2}_{max}$. The default value for chi2_{max} is 80.

Step 3: Harmonise the formatted datasets and estimate nuisance parameters

This step is designed for merging and harmonizing the formatted datasets from step 2 to make sure effect sizes for the same SNP correspond to the same allele for the exposure and outcome.

Then estimate the variance-covariance matrix Ω in background model and the residual correlation parameter ρ due to sample overlap through implementation of LD score regression.

The analysis for this step can be accomplished by function `est_paras()`:

```
paras = est_paras(dat1 = BMI,
                  dat2 = T2D,
                  trait1.name = "BMI",
                  trait2.name = "T2D",
                  ldscore.dir = "./eur_w_ld_chr")

#> Merge data1 and dat2 by SNP ...
#> Harmonise the direction of SNP effects of trait 1 and trait 2
#> Read in LD scores ...
#> Add LD scores to the harmonised data sets...
#> The Harmonised dataset will also be used for MR analysis
#> Begin estimation of Sigma and Omega using LDSC ...
#> Estimate heritability for trait 1 ...
#> Using two-step estimator with cutoff at 30.
#> Mean Chi2:2.9163.
#> Intercept: 1.1528(0.0251).
#> Total Observed scale h2:0.2093(0.0064).
#> Estimate heritability for trait 2 ...
#> Using two-step estimator with cutoff at 30.
#> Mean Chi2:1.654.
#> Intercept: 1.0911 (0.0168).
#> Total Observed scale h2:0.0445 (0.0024).
#> Estimate genetic covariance ...
#> Using two-step estimator with cutoff at 30.
#> Intercept: 0.1623 (0.0146).
#> Total Observed scale gencov: 0.0528 (0.0028).
#> Time elapsed: 127.348
```

Here the argument “ldscore.dir” specifies the path to LD score files.

Try the following commands to see the estimates:

```
paras$Omega
#>           [,1]           [,2]
#> [1,] 1.783664e-07 4.502337e-08
#> [2,] 4.502337e-08 3.790286e-08

paras$Sigma_err
#>           [,1]           [,2]
#> [1,] 1.1528175 0.1623214
#> [2,] 0.1623214 1.0911408
```

Note : the non-diagonal elements of `Sigma_err` is the estimate of ρ by the intercept of cross-trait LD score regression. The diagonal elements of `Sigma_err` is the intercept of single-trait LD score regressions. The

intercept can be fixed at 1 by specify “h2.intercept = T” in function “est_params()”.

The harmonized dataset will be used for LD clumping.

```
head(paras$dat)
#>      SNP A1 A2      b.exp      b.out      se.exp      se.out
#> 1 rs1000000 A G -0.0004639259 -0.0003772375 0.001611081 0.001320331
#> 2 rs10000010 C T -0.0011894904 0.0006435422 0.001629624 0.001304960
#> 3 rs1000002 T C -0.0032158697 -0.0033518849 0.001610943 0.001313804
#> 4 rs10000023 G T 0.0042081267 0.0004016327 0.001639800 0.001338776
#> 5 rs1000003 G A -0.0025518606 -0.0021567128 0.001614355 0.001315960
#> 6 rs10000033 C T 0.0023086159 0.0025334233 0.001616632 0.001317380
#>      pval.exp pval.out      L2
#> 1 0.77340 0.77320 20.90684
#> 2 0.46540 0.61580 23.52035
#> 3 0.04592 0.01035 46.00095
#> 4 0.01028 0.76390 25.62940
#> 5 0.11380 0.10140 37.94950
#> 6 0.15310 0.05519 10.62294
```

Step 4: IVs selection and LD clumping

Specify the IV selection threshold to obtain a selected dataset satisfying “pval.exp < Threshold”. Then apply the “clump()” function to the selected SNPs, which uses the PLINK clump method to extract a data frame for a set of nerarly independent SNPs.

```
Threshold = 5e-05      # IV selection threshold
MRdat = clump(subset(paras$dat,
                    pval.exp < Threshold),
              SNP_col = "SNP",
              pval_col = "pval.exp",
              clump_kb = 1000,
              clump_r2 = 0.001)
#> API: public: http://gwas-api.mrcieu.ac.uk/
#> Clumping ZKNb17, 25694 variants
#> Removing 24466 of 25694 variants due to LD with other variants or absence from LD reference panel
```

Note: by default, clump() pefroms LD clumping through API, which means you don’t need to install PLINK tools in your machine (see dependencies in <https://github.com/MRCIEU/ieugwasr>). For sure, you can do LD clumping locally with PLINK, (see ?clump()).

Step 5: Fit MRAPSS

Fit MRAPSS when parameters estimates (Sigma_err and Omega) and summary statistics of clumped SNPs (MRdat) are ready.

```
MRres = MRAPSS(MRdat,
               exposure = "BMI",
               outcome = "T2D",
               Sigma_err = paras$Sigma_err,
               Omega = paras$Omega ,
               Threshold = Threshold)
#> *****
#> MR test results of BMI on T2D :
#> MR-APSS: beta = 0.2427 beta.se = 0.0241 pval = 8.3998e-24 #SNPs= 1228
```

```
#> Correlation parameter (rho) due to sample overlap : 0.1623214
#> Proportion of effective IVs with foreground signals: 0.2819996
#> Variance component (Omega) for background model =
#>      [,1]      [,2]
#> [1,] 1.783664e-07 4.502337e-08
#> [2,] 4.502337e-08 3.790286e-08
#> Variance component (Lambda) for foreground model =
#>      [,1]      [,2]
#> [1,] 1.25785e-06 0.000000e+00
#> [2,] 0.00000e+00 2.549167e-07
#> *****
```

Visualize MRAPSS analysis results by “MRplot()”:

```
MRplot(MRres, exposure = "BMI", outcome = "T2D")
```

