# CARMA Tutorial

Zikun Yang

8/23/2022

## Introduction

This document describes a complete walk through the usage of the package 'CARMA' with an application to computing the posterior inclusion probability (PIP) of variants at loci of interest. In this document, we will illustrate typical fine-mapping studies with two types of datasets:

- Summary statistics based on individual level phenotype and genotype data, and in-sample linkage disequilibrium (LD) matrix.
- Summary statistics generated by meta-analysis, and LD matrix extracted from reference panels.

First, download the example datasets from GitHub repository ZikunY/CARMA. The directory path of the demo data should be under the repo folder of the git clone unless the user setwd to the git clone directory.

```
git clone https://github.com/ZikunY/CARMA.git
```

## Individual level data

### Simulating data

We simulate individual level data for the purpose of this demonstration. We use the R package 'sim1000G' (Dimitromanolakis et al. 2019) to simulate genotypes based on the 1000 Genomes Project data (phase 3, European population). The phenotype is simulated through a Gaussian regression model with the simulated genotypes $\boldsymbol{X}$:

$$\boldsymbol{y} = \boldsymbol{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

where $\boldsymbol{\beta}$ is a sparse coefficient vector such as $\beta_i \neq 0$ if the $i$th SNP is a causal SNP, and $\boldsymbol{\epsilon}$ is the standard Gaussian error. The probability of a variant being causal is computed based on the linear predictor $\boldsymbol{w}_i'\boldsymbol{\theta}$, where $\boldsymbol{w}_i$ is the vector of annotations associated with the $i$th SNP and $\boldsymbol{\theta}$ is the coefficients vector of the annotations.

### Example of locus chr1: 200,937,832-201,937,832

In this section, we use the simulated data based on the locus chr1:200937832-201937832. We computed the summary statistics (Z-scores) and the LD matrix. The pre-determined causal SNPs are the **287**th, **1275**th, and **2572**th SNPs at the locus (the left, middle, right red points respectively).

As shown in the figure below, one of the causal SNP has few highly correlated SNPs and larger Z-scores, whereas the other two SNPs are highly correlated to the surrounding SNPs with similar values of Z-scores.

**Running CARMA without annotations**  We run CARMA without annotations first. The input format of CARMA is the list class. We use the "CARMA'' function in the package, which is designed to run in-sample data. As recommended in the paper, we choose the dimensional hyperparameter $\eta$ as $1/\sqrt{p}$, where $p$ is the total number of SNPs at the locus. Notice that without annotations, all SNPs have identical prior probabilities of being causal generated by the Poisson prior distribution, which assigns prior probability on the model size and provides false discovery control.
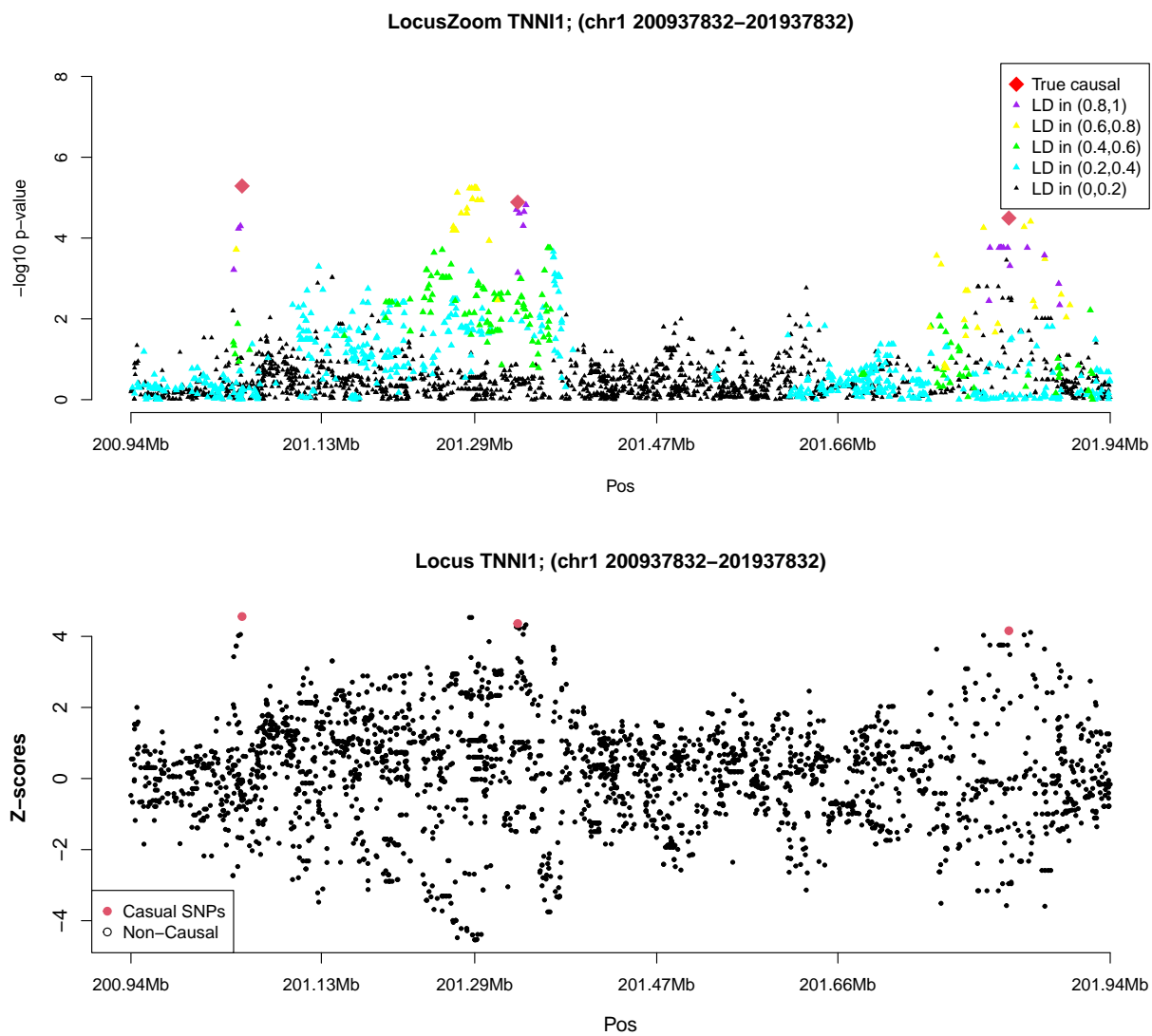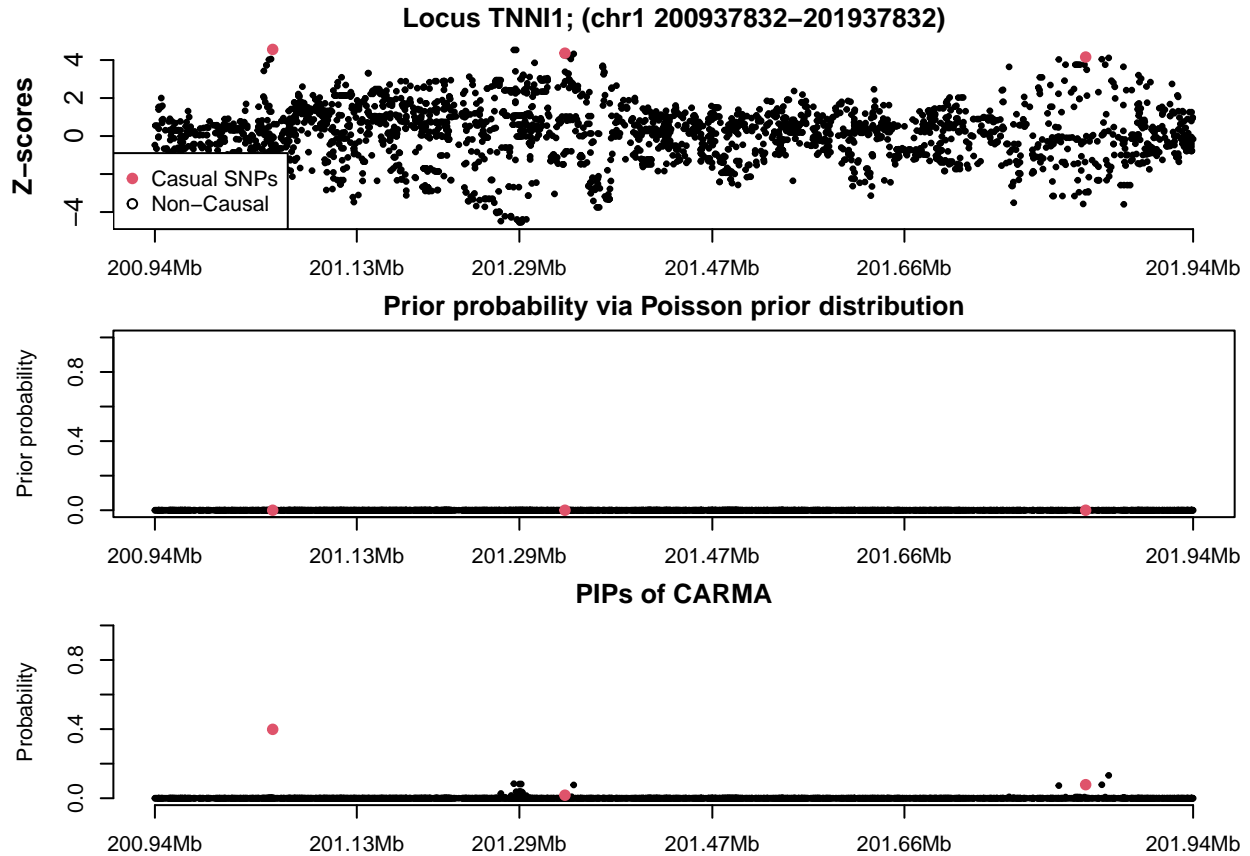
Figure 1: LocusZoom and Z-score plots for locus chr1:200,937,832-201,937,832

```
setwd('CARMA') ### setting up the working directory or the wd where the data are stored
data<-readRDS('Sample_data/in-sample_data.RData')
z.list<-list()
ld.list<-list()
lambda.list<-list()
z.list[[1]]<-data$Z
ld.list[[1]]<-data$LD
lambda.list[[1]]<-1/sqrt(nrow(ld.list[[1]]))
CARMA.results<-CARMA(z.list,ld.list,lambda.list=lambda.list)
```

We can check the results.



```
##       SNPs.index Causal.status Z.scores        PIPs
## 287          287  True causal  4.558480  0.39867786
## 2603        2603    Non-causal  4.113919  0.13225303
## 1095        1095    Non-causal  4.531850  0.08410905
## 1120        1120    Non-causal -4.536846  0.08223600
## 1131        1131    Non-causal -4.536846  0.08223600
## 2572        2572  True causal  4.158553  0.07841921
## 2591        2591    Non-causal  4.042399  0.07810110
## 1304        1304    Non-causal  4.325996  0.07686544
## 2530        2530    Non-causal  4.030837  0.07314266
## 1121        1121    Non-causal -4.536846  0.03893160
## 1125        1125    Non-causal -4.536846  0.03893160
## 1105        1105    Non-causal  4.531850  0.03777435
## 1134        1134    Non-causal -4.525274  0.03660785
## 1110        1110    Non-causal -4.400549  0.03462106
```

```
## 1141         1141    Non-causal -4.386612 0.03256794
## 1057         1057    Non-causal -4.478988 0.02743057
## 1268         1268    Non-causal  4.265880 0.01807633
## 1272         1272    Non-causal  4.360784 0.01750991
## 1273         1273    Non-causal  4.360784 0.01750991
## 1275         1275   True causal  4.360784 0.01750991
```

We can observe that the 287th SNP (the causal SNP at left), which is a true causal SNP with a larger Z-score comparing to other highly correlated SNPs, received a medium PIP value. On the other hand, the other two causal SNPs, which are highly correlated to other surrounding SNPs with similar Z-scores, shared the PIPs with the highly correlated SNPs. We can also check the credible sets and credible models.

```
CARMA.results[[1]]$`Credible set`[[2]]
```

```
## list()
```

```
CARMA.results[[1]]$`Credible model`[[3]]
```
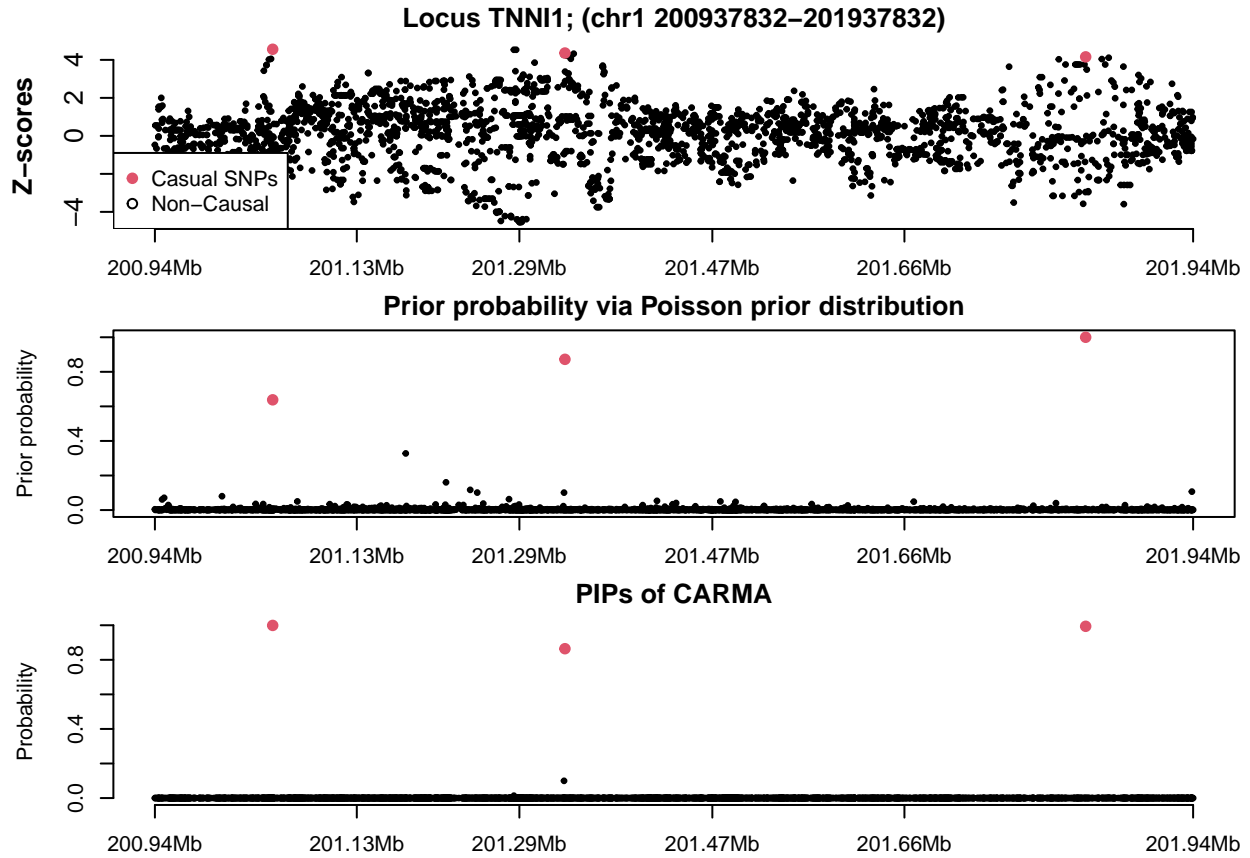
```
##  [1]  287 2603 2591 2530 1120 1121 1125 1131 1095 1105 1134 2572 1057 1304 1109
## [16] 1110 1141 1158 1272 1273 1275 1086
```

Due to relatively weak signal strength, none of the signals have enough PIPs to formulate credible sets. On the other hand, credible model still identified 22 candidate SNPs, which include all three true causal SNPs.

**Running CARMA with annotations**   We can include functional annotations into CARMA:

```
data<-readRDS('Sample_data/in-sample_data.RData')
z.list<-list()
ld.list<-list()
lambda.list<-list()
annot.list<-list()
z.list[[1]]<-data$Z
ld.list[[1]]<-data$LD
lambda.list[[1]]<-1/sqrt(nrow(ld.list[[1]]))
annot.list[[1]]<-data$Annotations
CARMA.results<-CARMA(z.list,ld.list,lambda.list=lambda.list,w.list=annot.list)
```

We can first check the resulting PIPs. This time, we include the prior probability of a variant being causal estimated by CARMA.

**Locus TNNI1; (chr1 200937832–201937832)**

**Prior probability via Poisson prior distribution**

**PIPs of CARMA**

```
##      SNPs.index Causal.status  Z.scores         PIPs
## 287         287   True causal  4.558480 0.9990413503
## 2572       2572   True causal  4.158553 0.9939951828
## 1275       1275   True causal  4.360784 0.8643504788
## 1273       1273    Non-causal  4.360784 0.0996357305
## 1095       1095    Non-causal  4.531850 0.0137849561
## 2603       2603    Non-causal  4.113919 0.0030364038
## 1272       1272    Non-causal  4.360784 0.0027148441
## 1131       1131    Non-causal -4.536846 0.0025171780
## 1121       1121    Non-causal -4.536846 0.0024274259
## 1304       1304    Non-causal  4.325996 0.0023687987
## 1109       1109    Non-causal -4.400549 0.0021916386
## 1827       1827    Non-causal -2.575304 0.0014537779
## 1294       1294    Non-causal  4.056951 0.0009512645
## 1110       1110    Non-causal -4.400549 0.0009001549
## 2521       2521    Non-causal -3.157733 0.0008416967
## 1105       1105    Non-causal  4.531850 0.0008121812
## 2591       2591    Non-causal  4.042399 0.0007778996
## 1125       1125    Non-causal -4.536846 0.0006670210
## 1120       1120    Non-causal -4.536846 0.0006559526
## 2571       2571    Non-causal -2.959334 0.0006052759
```

As observed from the figure above, the inclusion of annotations helps CARMA distinguish the true causal variants from the highly correlated SNPs, such as the 1275th and 2572th SNP, which in the absence of functional annotations cannot be distinguished from other highly correlated SNPs. Also, the 287th SNP receives a larger PIP this time. We can also examine the credible sets and credible models of CARMA.

```r
CARMA.results[[1]]$`Credible set`[[2]]
```

```
## [[1]]
## [1] 287
##
## [[2]]
## [1] 2572
##
## [[3]]
## [1] 1275 1273 1095 1272 1131 1121 1304 1109 1294
```

```r
CARMA.results[[1]]$`Credible model`[[3]]
```

```
## [1]  287 1275 2572 1273
```

The numbers of included SNPs in credible models have been reduced significantly. Also, the credible sets strengthened by the annotations identified the three true causals.

## Summary statistics and LD matrix extracted from reference panels

Usually, individual level data are not available in large GWAS studies. Instead, summary statistics are made available and an external LD matrix is used. These complex meta-analysis settings create inconsistencies between summary statistics and LD values which can lead to biased PIP values.

We use summary statistics from a meta-analysis for Alzheimer's disease (AD) (Jansen et al. 2019). The meta-analysis of AD is based on clinically diagnosed AD and AD-by-proxy with 71,880 cases and 383,378 controls of European ancestry. The clinically diagnosed AD case-control data are from 3 consortia (PGC-ALZ, IGAP, and ADSP), and the AD-by-proxy data are based on 376,113 individuals of European ancestry from UK BioBank (UKBB). We use the LD matrix extracted from the UKBB. For the CARMA model, we include 187 annotations provided by PolyFun plus PolyFun prior causal probability (Weissbrod et al. 2020).

### Demonstration with the loci ADAMTS4 and CR1

We illustrate CARMA on two loci, ADAMTS4 and CR1 on chromosome 1. We extract data at locus ADAMTS4/CR1, and extract the corresponding LD matrices from the UKBB (provided by PolyFun).

### Sample of data at the locus ADAMTS4

```
##        uniqID.a1a2 CHR        BP A1 A2        SNP          Z         P  Nsum
## 1 1:160656603_A_T   1 160656603  A  T   rs6702441  0.14921734 0.8813821 429975
## 2 1:160657127_T_C   1 160657127  T  C rs143426473  0.81166395 0.4169845 435185
## 3 1:160657137_G_A   1 160657137  G  A  rs11589131  0.04962457 0.9604216 429757
## 4 1:160657197_C_G   1 160657197  C  G  rs10908797  1.75450624 0.0793438 377075
## 5 1:160657356_C_T   1 160657356  C  T rs145169682 -0.28401355 0.7764000  17477
## 6 1:160658364_G_A   1 160658364  G  A   rs7539434  0.23514690 0.8140947 380902
##      Neff  dir       EAF          BETA          SE
## 1 423496.9 ?+-+ 0.3695570  0.0003359043 0.002251107
## 2 428659.9 ?+++ 0.0337234  0.0048561053 0.005982901
## 3 423280.9 ?+-+ 0.3693410  0.0001117523 0.002251954
## 4 375814.5 ??++ 0.0967510  0.0068457320 0.003901800
## 5  17477.0 ???- 0.0117955 -0.0140704510 0.049541479
## 6 379634.8 ??++ 0.0280567  0.0016341896 0.006949654
```

### Sample of data at the locus CR1

```
##        uniqID.a1a2 CHR        BP A1 A2        SNP          Z         P  Nsum
## 1 1:207287187_T_C   1 207287187  T  C   rs2808470  0.75499216 0.4502537 433909
## 2 1:207288258_C_T   1 207288258  C  T rs147553990 -0.62678527 0.5308000 364527
```

```
## 3 1:207288297_T_C   1 207288297   T   C   rs17020983   1.13570436 0.2560803 434723
## 4 1:207288309_T_G   1 207288309   T   G   rs79498904   0.87863355 0.3796000 364051
## 5 1:207288392_G_A   1 207288392   G   A   rs17020993   1.10788037 0.2679135 436498
## 6 1:207288897_T_C   1 207288897   T   C   rs12031629   0.09394084 0.9251562  71639
##      Neff  dir      EAF           BETA          SE
## 1 427395.4 ?-++ 0.19396400  0.0020652587 0.002735470
## 2 364527.0 ??-? 0.00332615 -0.0127494442 0.020341008
## 3 428202.0 ?-++ 0.09601170  0.0041656405 0.003667892
## 4 364051.0 ??+? 0.01279640  0.0091614442 0.010426923
## 5 429961.1 ?+++ 0.15706800  0.0032833906 0.002963669
## 6  71639.0 ?-?+ 0.42993500  0.0005013039 0.005336379
```

From the AD data we use Z-scores. Notice that the sample size values in the column "Nsum" can vary from 9,703 to 444,006 depending on which datasets are included in the meta-analyses of the AD study.
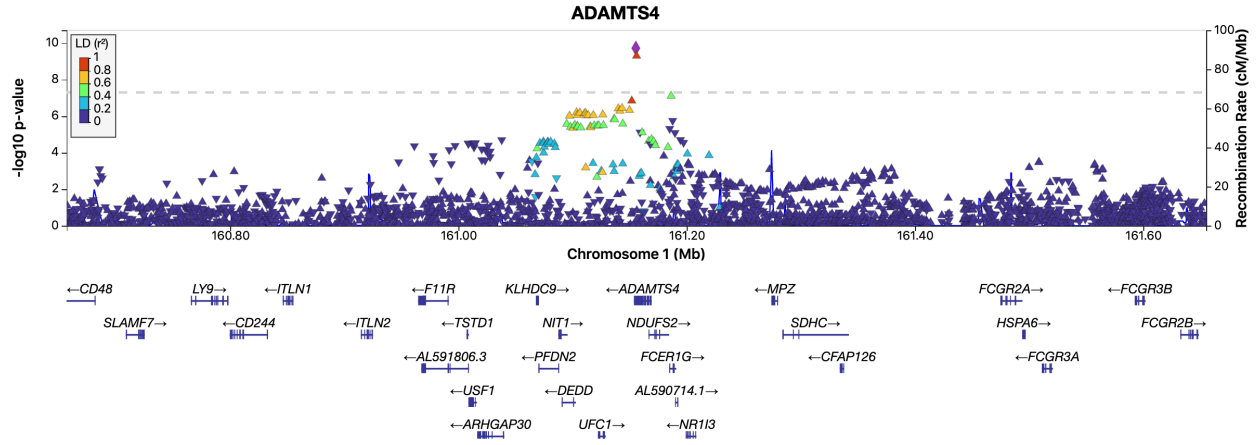


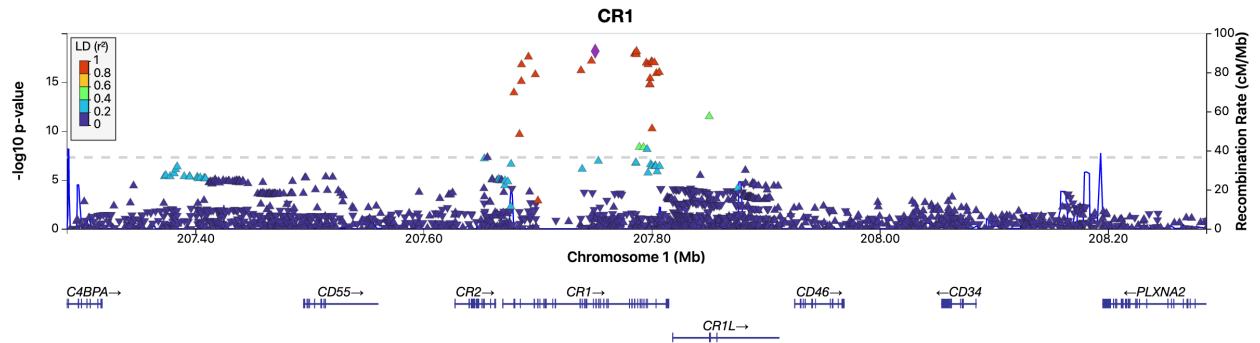Figure 2: LocusZoom plot for ADAMTS4.



Figure 3: LocusZoom plot for CR1.

Next we run CARMA with two settings: 1. without annotations, and 2. with annotations as described above. We use the function "CARMA_fixed_sigma" to run the meta-analysis with the external LD. The hyperparameter $\eta$ is set at 1 due to the existence of discrepancies between Z/LD.
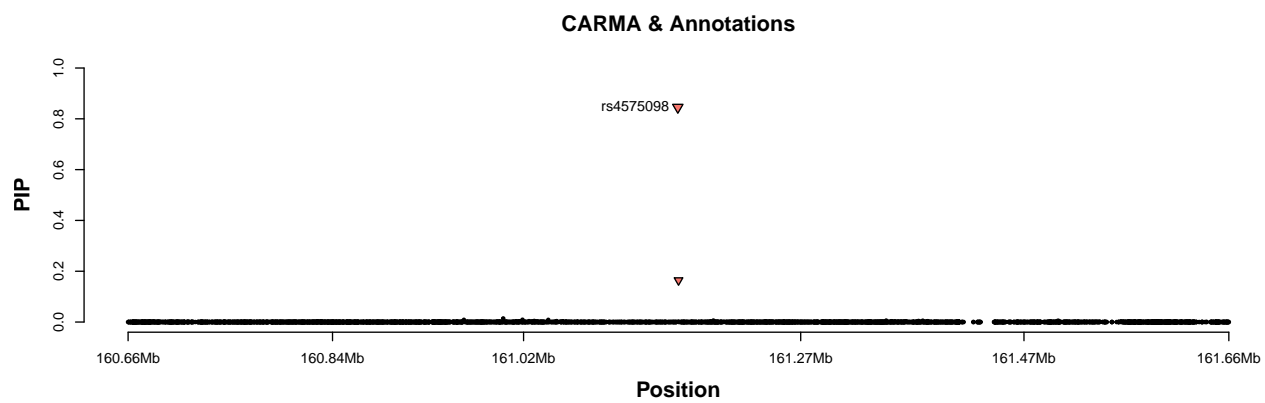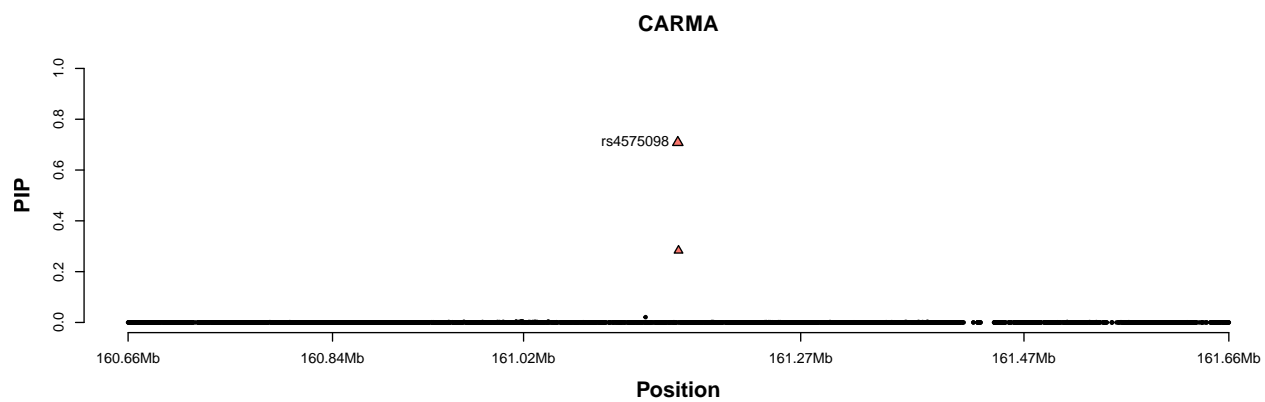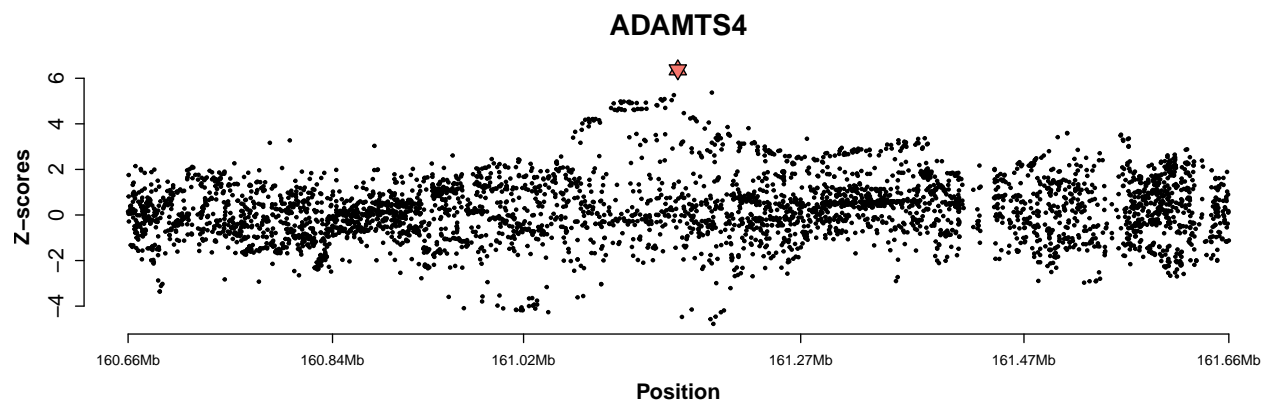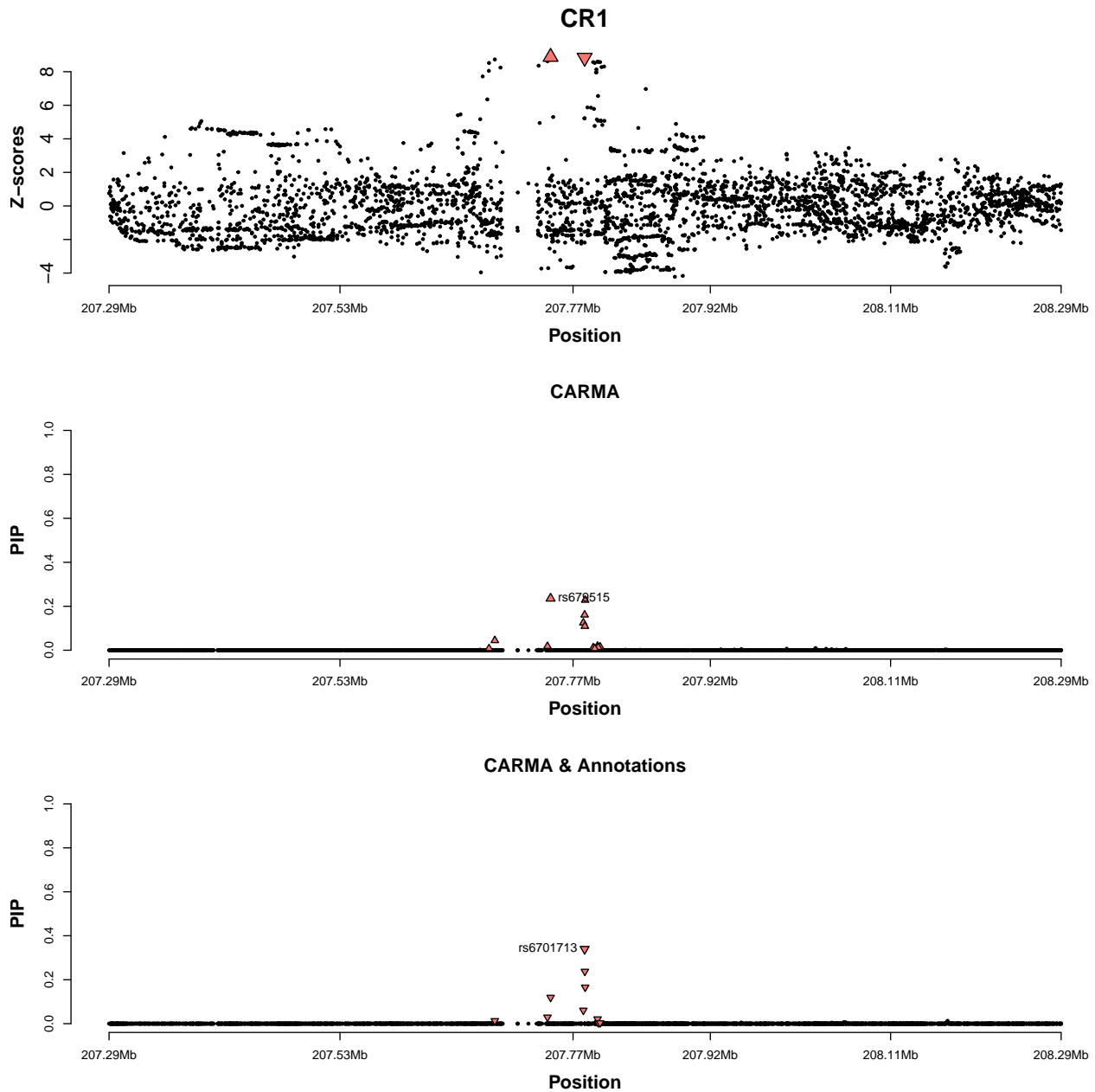
7

```r
Data_ADAMTS4<-readRDS('Sample_data/ADAMTS4.RData')
Data_CR1<-readRDS('Sample_data/CR1.RData')
z.list<-list()
ld.list<-list()
z.list[[1]]<-Data_ADAMTS4$`Meta-data`$Z
z.list[[2]]<-Data_CR1$`Meta-data`$Z
ld.list[[1]]<-Data_ADAMTS4$LD
ld.list[[2]]<-Data_CR1$LD
CARMA.results_no_annot<-CARMA_fixed_sigma(z.list,ld.list)
######With annotations
######The first 6 column of annotation file include location information etc.
annot.list<-list()
annot.list[[1]]<-as.matrix(cbind(1,Data_ADAMTS4$Annotations[,-(1:6)]))
annot.list[[2]]<-as.matrix(cbind(1,Data_CR1$Annotations[,-(1:6)]))
CARMA.results_annot<-CARMA_fixed_sigma(z.list,ld.list,w.list=annot.list)
```

First, we examine the PIPs estimated by CARMA.

**ADAMTS4**

**CARMA**

rs4575098

**CARMA & Annotations**

rs4575098

In the figure above, the credible sets are highlighted by colored shapes. Next, we can examine the SNPs included in the credible sets. For simplicity we only show the credible sets when including functional annotations.

```
## [1] "This is the first credible set of the locus ADAMTS4"
##      CHR       BP A1 A2       SNP        Z       PIPs
## 2184   1 161155392  A  G  rs4575098 6.369505 0.8453322
## 2186   1 161156033  A  C rs11585858 6.217633 0.1642777

## [1] "This is the first credible set of the locus CR1"
##      CHR       BP A1 A2       SNP        Z       PIPs
## 1563   1 207786289  A  G  rs6701713 8.837089 0.339477118
## 1567   1 207786542  A  G  rs2093761 8.791985 0.237695266
## 1570   1 207786828  A  G  rs2093760 8.877538 0.165632427
## 1426   1 207750568  T  C   rs679515 8.877562 0.119234143
```

```
## 1559    1 207784968   A   G   rs3818361 8.807716 0.060534523
## 1409    1 207747296   A   G   rs1752684 8.619752 0.029224223
## 1621    1 207799874   A   C rs10863420 8.601581 0.021014536
## 1343    1 207692049   A   G   rs6656401 8.728364 0.013134075
## 1626    1 207800555   T   C   rs1408078 8.582403 0.003682269
## 1629    1 207802552   A   C   rs4844610 8.572777 0.003369866
```

We can also examine the credible models.

```
## [1] "This is the credible model of the locus ADAMTS4"
```

```
##       CHR        BP A1 A2        SNP        Z       PIPs
## 2184    1 161155392  A  G  rs4575098 6.369505 0.8453322
## 2186    1 161156033  A  C rs11585858 6.217633 0.1642777
```

```
## [1] "This is the credible model of the locus CR1"
```

```
##       CHR        BP A1 A2       SNP        Z        PIPs
## 1563    1 207786289  A  G rs6701713 8.837089 0.33947712
## 1567    1 207786542  A  G rs2093761 8.791985 0.23769527
## 1570    1 207786828  A  G rs2093760 8.877538 0.16563243
## 1426    1 207750568  T  C  rs679515 8.877562 0.11923414
## 1559    1 207784968  A  G rs3818361 8.807716 0.06053452
```

## References

Dimitromanolakis, Apostolos, Jingxiong Xu, Agnieszka Krol, and Laurent Briollais. 2019. "sim1000G: A User-Friendly Genetic Variant Simulator in r for Unrelated Individuals and Family-Based Designs." *BMC Bioinformatics* 20 (1): 26.

Jansen, Iris E, Jeanne E Savage, Kyoko Watanabe, Julien Bryois, Dylan M Williams, Stacy Steinberg, Julia Sealock, et al. 2019. "Genome-Wide Meta-Analysis Identifies New Loci and Functional Pathways Influencing Alzheimer's Disease Risk." *Nature Genetics* 51 (3): 404–13.

Weissbrod, Omer, Farhad Hormozdiari, Christian Benner, Ran Cui, Jacob Ulirsch, Steven Gazal, Armin P Schoech, et al. 2020. "Functionally Informed Fine-Mapping and Polygenic Localization of Complex Trait Heritability." *Nature Genetics*, 1–9.