# User Manual for mrMLM

## 1. Introduction

mrMLM, a R package, aims to provide a user-friendly interface to conduct genome-wide association study (GWAS) via a multi-locus random-SNP-effect mixed linear model (mrMLM) methodology and to visualize its results. It works on the platforms of Windows, Linux and MacOS. The GUI is based on available add-on package RGtk2, via the aid of another package gWidgetsRGtk2. The visualization of results is based on package qqman, such as Manhattan and QQ plots.

## 2. Installation

### 2.1 Install GTK+

You may need to install GTK+ before installing RGtk2, because RGtk2 depends on GTK+.

For **Windows** user, you do as below:
Download GTK+ here
(http://sourceforge.net/projects/gladewin32/files/gtk%2B-win32-runtime/2.10.11/gtk-2.10.11-win32-1.exe).
Run the resulting file (gtk-2.10.11-win32-1.exe), which is an automated installer that will help you complete the installation of Gtk2 libraries.

For **Mac OS** users, you do as below:
Download GTK+ here (http://sourceforge.net/projects/gtk-osx/files/latest/download).
Extract and run the resulting file (gtk-osx-docbook-1.2.tar.gz).

For **Linux** users, you do as below:
You may or may not upgrade the GTK libraries depending on your distribution.
There are more details on RGtk2 at RGtk2's home page (http://www.ggobi.org/rgtk2/).

### 2.2 Install R

Download R from CRAN (https://cran.r-project.org/) and install it by running the downloaded file.

### 2.3 Install the R packages

The following R packages are needed: RGtk2, cairoDevice, gWidgets, gWidgetsRGtk2, RGtk2Extras and qqman, which can be downloaded from CRAN (https://cran.r-project.org/). Install them in order, as some depend on others. Within R environment, these packages can be installed directly using the below command:
install.packages(pkgs=c("RGtk2","cairoDevice","gWidgets","gWidgetsRGtk2","RGtk2Extras","qqman"))

## 2.4 Install mrMLM

The mrMLM package is freely available at the CRAN (https://cran.r-project.org/web/packages/mrMLM/index.html or soyzhang@mail.hzau.edu.cn or soyzhang@hotmail.com), you can download or request this R software. Within R environment, the mrMLM software can be installed directly using the below command:

install.packages(pkgs="mrMLM")

# 3. Running

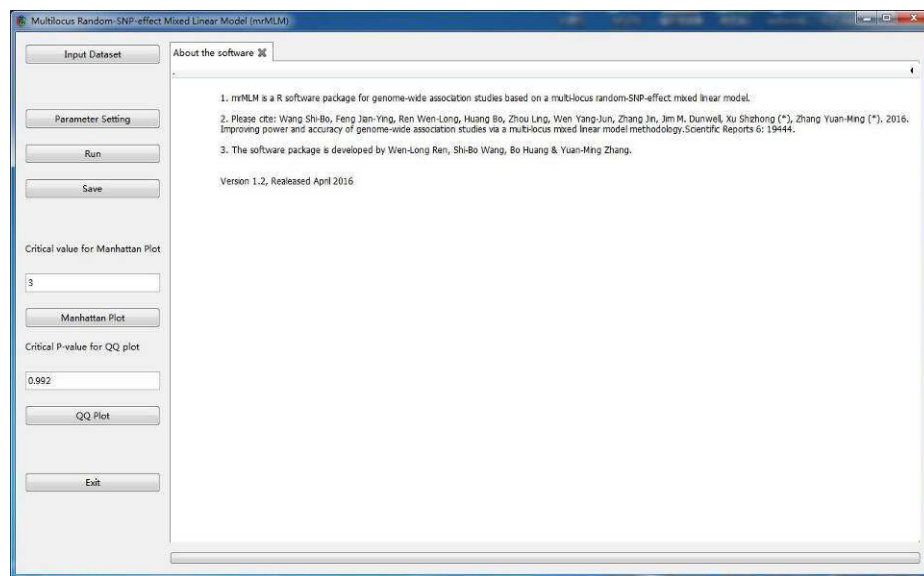The **RUN** steps are described as below. Within R environment, launch the mrMLM by command: library(mrMLM), then the following dialog will appear.



**Figure 1**. The GUI for the mrMLM

To restart the GUI, the command mrMLM() can be issued.

## 3.1 Input Dataset

Use the **Input Dataset** button to input dataset files, and then a dialog box will be appeared. In the dialog box, there are four steps. First, users select the dataset formats, which include **mrMLM numeric** format, **mrMLM character** format and **hapmap** format used in the TASSEL software. Then, use the **Genotype** and **Phenotype** buttons to input the genotypic and phenotypic datasets, respectively. Once one file is successfully uploaded, one tabbed page is added to the software notebook. Third, two things will be implemented in this step. One is to sort the individuals between the genotypic and phenotypic files and all the common individuals between the two files are selected to be analyzed in the further analyses. Another is to transfer the character genotypes into the numeric genotypes if the genotypes are character. Once users press the **DO** button, the two things will be conducted. Once the two files will be successfully uploaded, two tabbed pages (Genotype and Phenotype) will be added to the software notebook. Finally, use the **Kinship** and

**Population Structure** buttons to input the kinship and population structure matrices, respectively. If one file is successfully uploaded, the corresponding data page will be added to the notebook. Note that the **Kinship** and **Population Structure** buttons have two options. For the **kinship** button, one is to directly upload the kinship matrix and another is to calculate the kinship matrix in this software. For the **Population Structure** button, the population structure matrix may be not included in the mixed linear model of the GWAS if it has no effect on GWAS. If not, it should be included in the mixed model.

**Note:** About the **input file formats** in details, please see **Direction 1** in the end of the manual.
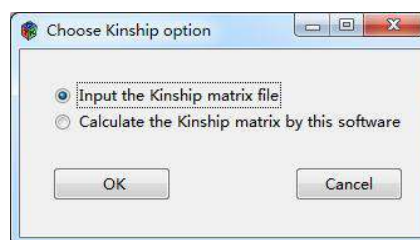


**Figure 2.** The **Input Dataset** dialog



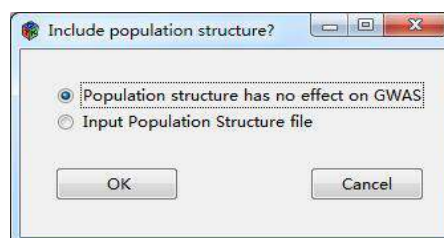**Figure 3.** The **Kinship** dialog



**Figure 4.** The **population structure** dialog

### 3.2 Run Program

Use the **Parameter Setting** button to set parameters before run the program. In the rMLM,

"Critical P-value in Manhattan plot" is set at $0.05/m_e$, where $m_e$ is the effective number of markers (please see Wang et al. Scientific Reports 2016, 6: 19444). All the SNPs significantly associated with the trait are marked by the value of $0.05/m_e$ in the last column of the **Result1** file. If users want to change the **Critical P-value** in Manhattan plot, a new $-\log_{10}$(P-value) may be input in the box above the **Manhattan Plot** button.

"Search radius of candidate gene (kb)" means to keep the one marker with the least P-value, and to delete all the other markers within the radius of the associated marker with the least P-value. Use the **Run** button to execute the software. If the program runs, a progress bar with the "**Please be patient...**" words will appear in the bottom of the interface. If the program finished, a bar with the "**All done.**" will appear.
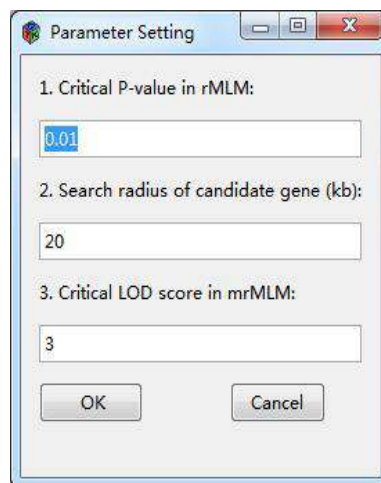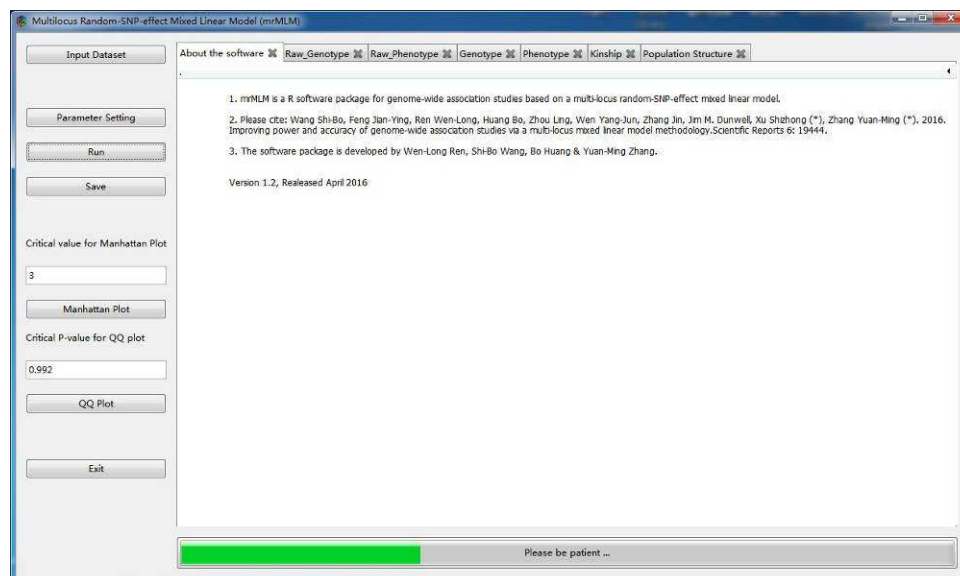


**Figure 5.** The **Parameter Setting** dialog



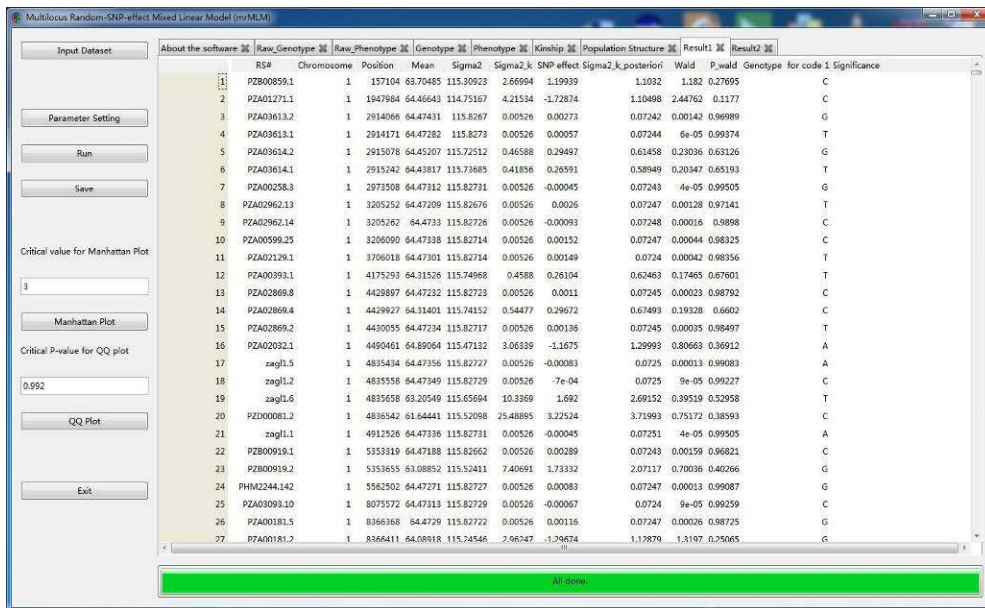**Figure 6.** A **running** program interface

**Figure 7.** A finished program interface (the **rMLM Results:** Result1)
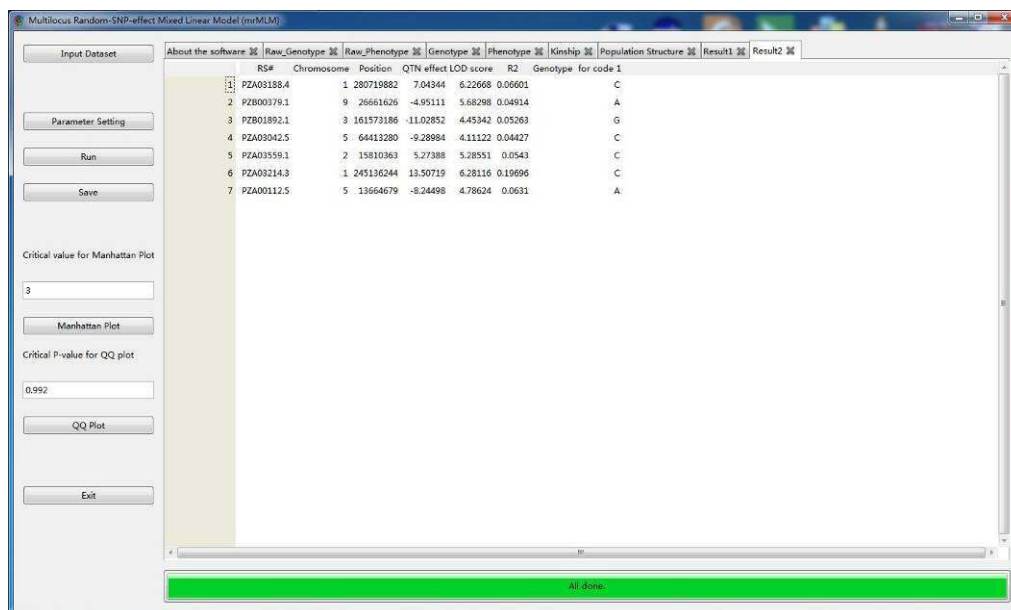


**Figure 8.** A finished program interface (the **mrMLM Results:** Result2)

### 3.3 Save results

Use **Save** button to save the results as **\*.csv** files. The **Results in the rMLM** are saved as Result1.csv and the **Results in the mrMLM** are saved as Result2.csv. If click **OK** button (after the **Save** button), a dialog is used to choose the pathway and the file name for the saving files.

**Note:** About the **explanation of Result1 and Result2** in details, please see **Direction 2** in the end of the manual.
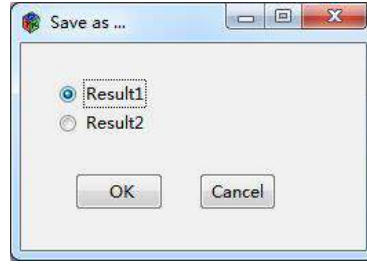
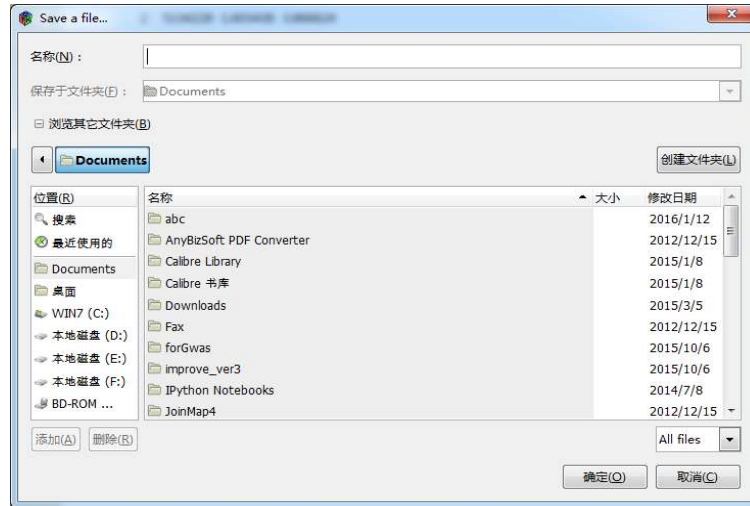**Figure 9.** The result **Save** dialog for the rMLM and mrMLM methods



**Figure 10.** The **Save** dialog

**Warning:** It is better not to include other languages except English in the pathway and file name. Otherwise, there may be something wrong.

# 4. Visualization of Results

If the program have finished, you may have the result visualization. Before use the **Manhattan Plot** button, please set the critical value for $-\log_{10}(P)$, which is defaulted the value of $-\log(0.05/m_e)$, where $m_e$ is the effective number of markers. Of course, users may change this value. Before press the **QQ Plot** button, please set the critical P-value for QQ plot, which is defaulted the value of 0.992. This is because the P-values are a mixture of a $\chi^2$ distribution with one degree of freedom and a point mass at one. Note that users may also change this 0.992 based on yourself results.
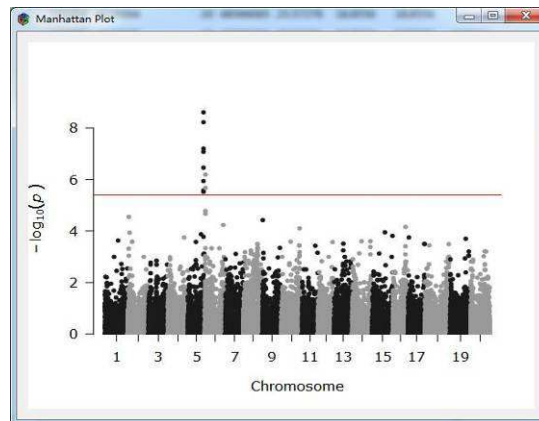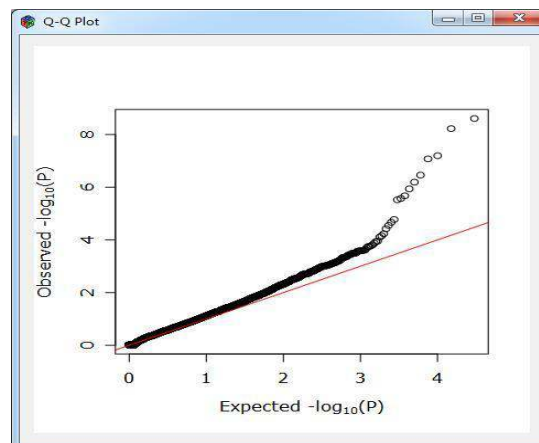
**Figure 11.** The **Manhattan Plot**



**Figure 12.** The **QQ Plot**

# Directions

**Direction 1: Explanation of input files in details**

**D1.1 The Genotypic file**

The **Genotypic** file should be a **\*.csv** format file.

**D1.1.1 mrMLM numeric format**

The first column, named "**rs#**" in the first row, stands for marker ID. The second column, named "**chrom**" in the first row, stands for chromosome. The third column, named "**pos**" in the first row, stands for the position (bp) of SNP in the chromosome. The fourth column, named "**genotype for code 1**" in the first row, stands for reference bases. And each of the remaining columns stands for one individual. In their first rows, the individual names are appeared. For each marker, homozygous genotypes are expressed by 1 and -1, respectively, and the heterozygous and missing genotypes are indicated by zero. Note that the genotypes with code **1** will be listed in the **Result** files.

**Figure D1.1.1** The **genotypic** file with **mrMLM numeric** format

## D1.1.2 mrMLM character format

The first three columns are same as those in the "**D1.1.1 mrMLM numeric format**". The marker values are character, such as **A, T, C, G** and **N**, and the other notations are heterozygous genotypes. The "**N**" indicates missing. The first rows from the fourth to last columns are individual code.



**Figure D1.1.2** The **genotypic** file with **mrMLM character** format

## D1.1.3 Hapmap (TASSEL) format

Please see the TASSEL software in details. Here we introduce simply. The first eleven columns describe the specific information of markers and individuals, and their column names must be **"rs#"**, **"alleles"**, **"chrom"**, **"pos"**, **"strand"**, **"assembly#"**, **"center"**, **"protLSID"**, **"assayLSID"**, **"panel"** and **"QCcode"**. In the **"rs#"** (1)**, "chrom"** (3) **and "pos"** (4) columns, the information for each marker (row) must be listed, and their meanings are same as those in the above. The values for marker genotypes should be character, such as **AA, TT, CC, GG, NN, AC** and **AG**, where the "**NN**" indicates missing or unknown genotypes. In the 2 and 5 to 11 columns, the **no available** information must be marked by **"NA"**. All the individuals will be showed from the 12 to last columns, the first element in each column is individual ID (name) and the others are the genotypes (character).

**Figure D1.1.3** The **genotypic** file with **Hapmap (TASSEL)** format

## D1.2 The Phenotypic file

The **Phenotypic** file should be a **\*.csv** format file. The first column stands for individual ID, such as 33-16, Nov-38 and 4226. The second column is phenotypic values for the trait. Note that the phenotypic file includes only one trait. The first element in the first column must be **"<Phenotype>" or "<Trait>"**.



**Figure D1.2.** The **Phenotypic** file

## D1.3 The Kinship file

The **Kinship** file should be a **\*.csv** file. The number of rows (or columns) equals to the number of the common individuals between the phenotypic and genotypic datasets. If the Kinship matrix is calculated by this software, we calculate only the Kinship matrix between the common individuals. If the Kinship matrix has been obtained and uploaded from a known file, it is possible that the number and order of individuals in the known file are not consistent with those of the common (valid) individuals in the further analysis. At this situation, the software will change the known K matrix in order that the number and order of new K matrix match the number and order of common (valid) individuals in the Phenotypic and Genotypic files. In the known K matrix, the **first element** in the first column must be **the number of valid individuals**, such as 263; and the other elements in the first column are individual ID (names).

**Figure D1.3.** The **Kinship** file

**D1.4 The Population Structure file**

The **Population Structure** file is a **\*.csv** file. Using the **Structure** software, the population structure matrix may be calculated. The first column stands for the valid individual ID (names). The first and second elements in the first column must be **"<Covariate>"** and **"<Trait>"** respectively. If population structure matrix has **k** columns, please input all the **k** columns. In the second row, it must be **"Q1"**, **"Q2"**, ... , **"Qk"** following the **"<Trait>"**.

If the **Population Structure** file has been obtained and uploaded from a known file, it is possible that the number and order of individuals in the known file are not consistent with those of the common (valid) individuals in the further analysis. At this situation, the software will change the known matrix in order that the number and order of new **Population Structure** matrix match the number and order of common (valid) individuals in the Phenotypic and Genotypic files.



**Figure D1.4.** The **Population Structure** file

**Direction 2: Explanation of Result1 and Result2 in details**

**D2.1 Explanation of Result1 file**

The **Result1** table with twelve columns shows the results from the rMLM (random-SNP-effect mixed linear model) method. The corresponding column names are as follows: reference sequence number (rs#, marker name), chromosome, marker's position (bp) in the chromosome, population mean value (Mean), residual variance ($\sigma^2$, Sigma2), priori variance of the kth SNP effect ($\phi_k^2$, Sigma2_k), SNP effect ($\gamma_k$, Effect), posteriori variance of SNP effect ($\mathrm{var}(\gamma_k)$, Sigma2_k_posteriori), Wald test statistic value, the P-value of Wald test, significance and genotype for code 1, respectively. In the significance column, only significant markers under the critical value $0.05/m_e$ are marked.

| | A | B | C | D | E | F | G | H | I | J | K | L | M |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | RS# | Chromosome | Position | Mean | Sigma2 | Sigma2_k | SNP effect | Sigma2_k_posteriori | Wald | P_wald | Genotype for code 1 | Significance | |
| 2 | PZB00859.1 | 1 | 157104 | 63.70485 | 115.3092 | 2.66994 | 1.19939 | | 1.1032 | 1.182 | 0.27695 C | | |
| 3 | PZA01271.1 | 1 | 1947984 | 64.46643 | 114.7517 | 4.21534 | -1.72874 | | 1.10498 | 2.44762 | 0.1177 C | | |
| 4 | PZA03613.2 | 1 | 2914066 | 64.47431 | 115.8267 | 0.00526 | 0.00273 | | 0.07242 | 0.00142 | 0.96989 G | | |
| 5 | PZA03613.1 | 1 | 2914171 | 64.47282 | 115.8273 | 0.00526 | 0.00057 | | 0.07244 | 6.00E-05 | 0.99374 T | | |
| 6 | PZA03614.2 | 1 | 2915078 | 64.45207 | 115.7251 | 0.46588 | 0.29497 | | 0.61458 | 0.23036 | 0.63126 G | | |
| 7 | PZA03614.1 | 1 | 2915242 | 64.43817 | 115.7369 | 0.41856 | 0.26591 | | 0.58949 | 0.20347 | 0.65193 T | | |
| 8 | PZA00258.3 | 1 | 2973508 | 64.47312 | 115.8273 | 0.00526 | -0.00045 | | 0.07243 | 4.00E-05 | 0.99505 G | | |
| 9 | PZA02962.13 | 1 | 3205252 | 64.47209 | 115.8268 | 0.00526 | 0.0026 | | 0.07247 | 0.00128 | 0.97141 T | | |
| 10 | PZA02962.14 | 1 | 3205262 | 64.4733 | 115.8273 | 0.00526 | -0.00093 | | 0.07248 | 0.00016 | 0.9898 C | | |
| 11 | PZA00599.25 | 1 | 3206090 | 64.47338 | 115.8271 | 0.00526 | 0.00152 | | 0.07247 | 0.00044 | 0.98325 C | | |
| 12 | PZA02129.1 | 1 | 3706018 | 64.47301 | 115.8271 | 0.00526 | 0.00149 | | 0.0724 | 0.00042 | 0.98356 T | | |
| 13 | PZA00393.1 | 1 | 4175293 | 64.31526 | 115.7497 | 0.4588 | 0.26104 | | 0.62463 | 0.17465 | 0.67601 T | | |
| 14 | PZA02869.8 | 1 | 4429897 | 64.47232 | 115.8272 | 0.00526 | 0.0011 | | 0.07245 | 0.00023 | 0.98792 C | | |
| 15 | PZA02869.4 | 1 | 4429927 | 64.31401 | 115.7415 | 0.54477 | 0.29672 | | 0.67493 | 0.19328 | 0.8602 C | | |

**Figure D2.1.** Results in the rMLM (**Result1**)

**D2.2 Explanation of Result2 file**

The Result2 table with seven columns shows the final results of the mrMLM (multi-locus random-SNP-effect mixed linear model) method. The corresponding column names are as follows: reference sequence number (rs#, marker names), chromosome, marker's position (bp) in the chromosome, QTN effect, LOD score, the proportion of phenotypic variance explained by the putative QTN, and genotype for code 1, respectively.

| | A | B | C | D | E | F | G | H |
|---|---|---|---|---|---|---|---|---|
| 1 | RS# | Chromosom | Position | QTN effect | LOD score | R2 | Genotype | for code 1 |
| 2 | PZA03188.4 | 1 | 280719882 | 7.04344 | 6.22668 | 0.06601 | C | |
| 3 | PZB00379.1 | 9 | 26661626 | -4.95111 | 5.68298 | 0.04914 | A | |
| 4 | PZB01892.1 | 3 | 161573186 | -11.02852 | 4.45342 | 0.05263 | G | |
| 5 | PZA03042.5 | 5 | 64413280 | -9.28984 | 4.11122 | 0.04427 | C | |
| 6 | PZA03559.1 | 2 | 15810363 | 5.27388 | 5.28551 | 0.0543 | C | |
| 7 | PZA03214.3 | 1 | 245136244 | 13.50719 | 6.28116 | 0.19696 | C | |
| 8 | PZA00112.5 | 5 | 13664679 | -8.24498 | 4.78624 | 0.0631 | A | |

**Figure D2.2.** Results in the mrMLM (**Result2**)

**Reference**

Wang Shi-Bo, Feng Jian-Ying, Ren Wen-Long, Huang Bo, Zhou Ling, Wen Yang-Jun, Zhang Jin, Jim M. Dunwell, Xu Shizhong (*), Zhang Yuan-Ming (*). 2016. Improving power and accuracy of genome-wide association studies via a multi-locus mixed linear model methodology. Scientific Reports 6: 19444.