

piMASS User Manual v0.9

Yongtao Guan
Baylor College of Medicine

January 20, 2011

Contents

1	What piMASS Can Do?	2
2	Input File Formats	2
2.1	Mean Genotype File Format	2
2.2	Phenotype File Format	2
2.2.1	Normal Quantile Transformation	3
2.3	SNP Location File Format	3
2.4	Use of Multiple Genotype And Phenotype Files	4
2.4.1	The Strand Issue	4
3	Running piMASS	4
4	Output Files	6
4.1	Log file: <code>prefix.log</code>	6
4.2	SNP information file: <code>prefix.snp.txt</code>	6
4.3	Sampling path file: <code>prefix.path.txt</code>	6
4.4	Sampled model file: <code>prefix.gamma.txt</code>	6
4.5	Sampled model file: <code>prefix.mcmc.txt</code>	6

1 What piMASS Can Do?

The software is developed and maintained by Yongtao Guan, based on the work of Guan and Stephens (to appear, *Annals of Applied Statistics*, 2011), please refer to the paper for the details of the method. piMASS is designed to do the following:

- Perform multi-SNP association analysis either genome-wide, or on a small region.
- Estimate proportion of phenotypic variance explained by the genotypes.
- Estimate the posterior of the number of SNPs that have effects on phenotype.
- Report the marginal posterior inclusion probabilities of each SNP, a measure of the strength of the marginal association.

2 Input File Formats

The users should prepare genotype files and phenotypes, described below. Both types of files allow to have missing values. However, individuals who have either missing genotype or missing phenotype will be excluded from the study. Therefore, we recommend users to perform imputation on the original genotype to fill in the missing values. An alternative, but less ideal, solution is to fill in a missing genotype with the mean genotype of that SNP. For files that have multiple columns, different columns can be separated by either a coma, a blank space, or a semi-colon.

2.1 Mean Genotype File Format

The genotype file is identical in format to the mean genotype file used in BIMBAM , not surprisingly. The first column of the mean genotype files is the SNP ID, the second and third columns are allele types with minor allele first. The rest columns are the mean genotypes of different individuals – numbers between 0 and 2 that represents the (posterior) mean genotype, or dosage of the minor allele. An example of mean genotypes file of two SNPs and three individuals follows.

```
rs1, A, T, 0.02, 0.80, 1.50  
rs2, G, C, 0.98, 0.04, 1.00
```

2.2 Phenotype File Format

In the phenotype input file, each line is a number indicating the phenotype value for each individual in turn, in the same order as in the genotype file. Missing phenotypes should be denoted as NA. The number of lines should be equal to the number of individuals in genotype file (N), otherwise the program will either throw away the values after N or append “NA” at the end to observe N values. In either case, a warning will be printed.

Example phenotype file with 3 individuals:

1.2
NA
-0.3

If the phenotypes are binary (e.g. in a case-control study) then the format is the same, but each entry should be 0, 1 or NA. It does not matter which group is denoted 0 and which denoted 1.

2.2.1 Normal Quantile Transformation

For quantitative traits an important assumption underlying the methods implemented in `piMASS` is that the phenotype has a normal distribution within each genotype class. We suggest using a normal quantile transformation to transform the phenotype to be normal before running `piMASS`. For example, in R, this can be accomplished using `xtransformed = qqnorm(x,plot.it=F)$x`.

This quantile transformation does not fully solve the problem (it ensures that the phenotype is normal overall, but not necessarily normal within each genotype class). However, with the small effect sizes typical in genetic association studies it appears to be a simple sensible way to guard against strong departures from modeling assumptions. If you have other covariates that may be important predictors of phenotype (e.g. Age, Sex) we suggest first regressing the phenotype on these covariates using standard multiple linear regression software, and then running `piMASS` on the residuals from this regression (after applying a normal quantile transformation to these residuals).

2.3 SNP Location File Format

The file contains three columns: the first column is the SNP ID, the second column is its physical location, and the third column contains its chromosome number. Note, it is OK if the rows are not ordered according to position, but the file must contain all the SNPs in the genotype files. If the genotype files contain SNPs across different chromosome, `piMASS` will sort SNPs based on its chromosome and position.

Example file:

```
rs1, 1200, 1  
rs2, 4000, 1  
rs3, 3320, 1
```

Note: This file is strictly needed only if the order of the SNPs in the genotype file is not the same as the order of their physical locations along the chromosome, or if multiple genotype and phenotype files are used (see below). To align SNPs in the correct order is important for the sampling procedure because we propose exchange nearby SNPs to facilitate the mixing among correlated SNPs.

2.4 Use of Multiple Genotype And Phenotype Files

In some cases it may be convenient to provide genotypes (and corresponding phenotypes) in multiple files. For example, in a genome-wide study, it may be helpful to have one genotype file containing the case data, and a second genotype file containing the control data. Or one genotype files containing individuals in the first stage of the study and the another contains the second.

When using multiple genotype files `piMASS` does not require that the same SNPs be present in both files (although if the same SNP is present in both files then the SNP identifier should be the same in both files, to convey this information). However, SNPs missing in one of the files will cause the SNP to be excluded from the study because of the missingness.

When using multiple genotype files, the user must also provide multiple phenotype files, with each phenotype file corresponding to the individuals in a genotype file. The exception to this is that, for a case/control study, the case phenotypes can be specified by `-p 1` and control phenotypes can be specified by `-p z`.

2.4.1 The Strand Issue

When merging genotypes from different studies, there arises the issue of whether or not the genotypes for a SNP were obtained on the same strand. In some cases this can be checked easily: for example, if a SNP in one study is A/G, and in the other is T/C, we infer that the two studies used different strands, and we can flip one of the SNPs to correct this. `piMASS` performs these kinds of flip automatically. However, if a SNP is A/T, or C/G, one cannot tell whether the strandedness is the same or different across studies without external information. In this situation, `piMASS` assumes that genotypes for a single SNP in multiple input files refer to the same strand.

Note: if genotypes at a SNP are not compatible with the SNP being bi-allelic, even after strand flips (as might happen when multiple genotypes are used, see below), then the SNP is considered to be “bad” and `piMASS` will exclude the SNP from the study.

3 Running piMASS

First some general comments:

- `piMASS` is a command line based program. The command should be typed in a terminal window, in the directory in which `piMASS` executable exists.
- The command line should be all on one line: the line-break in the example is only because the line is too large to fit on one page.
- Unless otherwise stated, the “options” (`-g -p -pos -o`, etc.) are all case-sensitive.

Now we illustrate how to use `piMASS` through examples.

1. A single genotype file and a single phenotype file

```
./pimass -g cohort.txt -p pheno.txt -w 10000 -s 100000 -o pref -num 10
```

The command line will run MCMC with burn-in steps 1000 and sampling steps 100000, every 10 steps a sample will be recorded. The output file names will begin with `pref`.

2. Multiple genotype files and multiple phenotype files

```
./pimass -g cohort1.txt -p pheno1.txt -g cohort2.txt  
-p pheno2.txt -w 100000 -s 1000000 -o pref -pos pos.txt -num 10
```

This command line takes two genotype files and two phenotype files, merge them based on the SNP ID and sort them according to their locations in position files. In this example `piMASS` will run 100k warm-up steps and follow by 1M sampling steps, every 100 steps record states sampled.

3. Binary phenotypes

```
./pimass -g case_mgt.txt -p 1 -g ctrl_mgt.txt -p z -pos pos.txt  
-o pref -w 10000 -s 1000000 -num 100 -cc
```

This command line asks `piMASS` to take two genotype files. The `-p 1` assign all individuals in the matching genotype (`case_mgt.txt` in the example) as 1, and `'-p z'` assign all individuals in the matching genotype (`ctrl_mgt.txt` in the example) as 0. `piMASS` will run 10k warm-up steps and follow by 1M sampling steps, every 100 steps record states sampled. The `-cc` option tells `piMASS` that this data has binary phenotypes.

4. Setup other parameters

```
./pimass -g cohort.txt -p pheno.txt -w 10000 -s 100000 -o pref -num 10  
-hmin 0.01 -hmax 0.5 -pmin 1 -pmax 1000 -smin 1 -smax 200
```

This command line is identical to the first command line in the list except that it specifies ranges for the hyper-parameters h and p and the restrictions on the number of SNPs in the model. Specifically, it specifies that minimum and maximum of the h is 0.01 and 0.5 respectively, the minimum and maximum of π is 1 1000 out of total number of SNPs, respectively. In addition, it restricts the minimum and maximum number of SNPs in a model to be 1 and 200.

4 Output Files

piMASS will create output files in a directory named `output/`. (If this directory does not exist then it will be created.) Output files will be produced, each with a name beginning with “prefix” that was specified by the `-o` option. We now describe the contents of these output files.

4.1 Log file: `prefix.log`

A log file includes details of the run parameters used and any warnings generated. When sending in a bug report, it is important to include the log file as an attachment.

4.2 SNP information file: `prefix.snp.txt`

This file contains 10 columns: the SNP rsID, minor allele, major allele, minor allele frequency, variance estimates of the SNP genotypes, chromosome, position, numerical ID, $\log_{10}(\text{single SNP BF})$ based on default prior, and the effect size estimates. Among them, the numerical ID is used in the `prefix.gamma.txt` to record which SNPs are in the model.

4.3 Sampling path file: `prefix.path.txt`

This file contains 9 columns, each row corresponding a sampled states (excluding gamma). The first column is the h , the second column is the re-estimated h based on sampling posterior effect sizes, the rest columns are $\log_{10}(\pi)$, number of SNPs, acceptance ratio for the local proposals, acceptance ratio for the long-range proposals, the derived prior on effect sizes, $\log(\text{BF})$ of the model, and the log-prior of the model computed from π and number of SNPs in the model.

4.4 Sampled model file: `prefix.gamma.txt`

The first column of the file is the number of SNPs in the model, the rest columns are SNP IDs, when the number of SNPs is smaller than the maximum allowed number of SNPs, NAs are appended. The mapping between a numerical ID and its rsID can be found in `prefix.snp.txt` file.

4.5 Sampled model file: `prefix.mcmc.txt`

For quantitative phenotypes, it contains SNP ID, chromosome, position, estimates of the posterior inclusion probabilities based on simple counting, estimates of the posterior inclusion probabilities based on Rao-Blackwellisation, naive estimates of the posterior effect size, and Rao-Blackwellised estimates of the posterior effect size.

For binary phenotype (when `-cc` is used), the output file is different in that the two columns of the Rao-Blackwellised estimates are no longer there.

Appendix A: piMASS Options

Unless otherwise stated, *arg* stands for a string, *num* stands for a number.

FILE I/O RELATED OPTIONS:

- -g *arg* can use multiple times, must pair with -p.
- -p *arg* can use multiple times, must pair with -g. *arg* can be a file name; z or 1, which indicates the pairing genotype individuals have phenotype 0 or 1.
- -pos *arg* can use multiple times. *arg* is a file name.
- -o *arg* *arg* will be the prefix of all output files, the random seed will be used by default.
- -cc calc bf of logit regression on binary phenotype.

PARAMETER SETUP:

- -w(warm) *num* specify number of burn-in steps for MCMC.
- -s(step) *num* specify number of sampling steps for MCMC.
- -num *num* specify thinning, record one states in every *num* steps.
- -r *num* specify random seed, system time by default.
- -hmin *num* specify minimum value for h.
- -hmax *num* specify maximum value for h.
- -pmin *num* specify minimum value for p.
- -pmax *num* specify maximum value for p.
- -smin *num* specify minimum value for number of SNPs in the model.
- -smax *num* specify maximum value for number of SNPs in the model.
- -nstart *num* specify number of SNPs (with top marginal signal) in the model to begin with.

OTHER OPTIONS:

- -v(ver) print version and citation
- -h(help) print this help
- -exclude-maf *num* exclude SNPs whose maf < *num* , default 0.01.
- -exclude-nopos *num* exclude SNPs that has no position information, 1 = yes (default), 0 = no
- -silence no terminal output.