**Laboratory Note**

# Genetic Epistasis
# VII - Assessing Algorithm TEAM
**LN-7-2014**

**Ricardo Pinho and Rui Camacho**
**FEUP**
Rua Dr Roberto Frias, s/n,
4200-465 PORTO
Portugal
Fax: (+351) 22 508 1440
e-mail: ei09045@fe.up.pt
www : http://www.fe.up.pt/∼ei09045
rcamacho@fe.up.pt
www : http://www.fe.up.pt/∼rcamacho

May 2014

**Abstract**

In this lab note, the algorithm TEAM is presented. TEAM is an exhaustive algorithm that works by updating contingency tables and the minimum spanning tree made from SNPs. The results obtained show an increase in Power by population size, allele frequency, and odds ratio. There is also an increase in Type 1 Error Rate with population size, but not a clear indicator for allele frequency. The scalability of the algorithm is questionable, considering that there is a big increase in the running time required by data sets with different population sizes, which is not relevant for these experiments but may be problematic for larger data sets.

# 1 Introduction

Tree-based Epistasis Association Mapping) TEAM [ZHZW10] is an exhaustive algorithm that computes all two-locus pairs to obtain a permutation test, which is applicable to all statistical relevancy tests, due to the contingency table generated. TEAM also uses Family-wise error rate (FWER) and false discovery rate (FDR) to control error rate using the permutation test, which is better than Bonferroni correction but also more computationally expensive. The algorithm builds a minimum spanning tree containing SNPs as nodes and the edges represent the genotype difference between two SNPs. This is used to update the contingency tables, allowing for a pruning of many individuals.

The algorithm receives the SNPs genotypes and phenotypes of each individual, creating a specified number of permutations. The contingency tables for each single-locus are generated. The minimum spanning tree is built, using the different genotypes associated to each edge. The tree is then updated for each leaf node with the information related to the contingency table for genotype relation between SNPs. The test values are then calculated, using the contingency tables.

|       | $S_1$ | $S_2$ | $S_3$ | $S_4$ | $S_5$ | $S_6$ | $S_7$ | $S_8$ | $S_9$ | $S_{10}$ |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|----------|
| $X_1$ | 0     | 0     | 0     | 1     | 2     | 0     | 2     | 0     | 2     | 0        |
| $X_2$ | 2     | 0     | 0     | 2     | 0     | 2     | 0     | 1     | 2     | 1        |
| $X_3$ | 2     | 2     | 0     | 1     | 2     | 2     | 0     | 1     | 1     | 0        |
| $X_4$ | 0     | 2     | 2     | 0     | 0     | 0     | 0     | 1     | 0     | 1        |
| $X_5$ | 2     | 1     | 2     | 0     | 1     | 2     | 0     | 1     | 0     | 2        |
| $Y_1$ | 1     | 1     | 1     | 0     | 1     | 0     | 1     | 1     | 1     | 0        |
| $Y_2$ | 0     | 0     | 0     | 1     | 1     | 0     | 1     | 0     | 1     | 0        |
| $Y_3$ | 1     | 0     | 1     | 1     | 1     | 0     | 1     | 0     | 1     | 0        |

Table 1: An example of the input data, consisting of 5 SNPs $X_1...X_6$, the original Phenotype $Y_1$, and two permutations $Y_2,Y_3$ for 10 individuals $S_1,...,S_{10}$.

## 1.1 Input files

The input consists of 2 files: a file containing the genotype information and another containing the phenotype information for a number of individuals.

| | $X_{i=0}$ | | | $X_{i=1}$ | | | $X_{i=2}$ | | | Total |
|---|---|---|---|---|---|---|---|---|---|---|
| | $X_{i=0}$ | $X_{i=1}$ | $X_{i=2}$ | $X_{j=0}$ | $X_{j=1}$ | $X_{j=2}$ | $X_{j=0}$ | $X_{j=1}$ | $X_{j=2}$ | |
| $Y_{k=0}$ | Event $a_1$ | Event $a_2$ | Event $a_3$ | Event $b_1$ | Event $b_2$ | Event $b_3$ | Event $e_1$ | Event $e_2$ | Event $e_3$ | |
| $Y_{k=1}$ | Event $c_1$ | Event $c_2$ | Event $c_3$ | Event $d_1$ | Event $d_2$ | Event $d_3$ | Event $f_1$ | Event $f_2$ | Event $f_3$ | |
| Total | | | | | | | | | | $M$ |

Table 2: The contingency table between two SNPs $X_i$ and $X_j$ for a given phenotype $Y_k$. M refers to the total amount of individuals.

(a) Genotype

| 0011001121 |
|---|
| 1212111121 |
| 1001000102 |
| 2202121111 |

(b) Phenotype

| 0000000010 |
|---|

Table 3: An example of the input file containing genotype and phenotype information with 4 SNPs and 10 individuals. Genotype 0,1,2 corresponds to homozygous dominant, heterozygous, and homozygous recessive. The phenotype 0,1 corresponds to control and case respectively.

## 1.2 Output files

The output consists of a list of every SNP pair and the relevant test score. The score can be calculated for any statistics defined in the contingency table. In this experiment, the test score corresponds to chi-square statistic.

## 1.3 Parameters

The customizable parameters are as follows:

- *individual* - The number of individuals in the data. In this case it is dependent on the data set parameters.

- *SNPs* - Number of SNPs in the data. In this case it is dependent on the data sets ( which is fixed at 300).

- *permutation* - Number of permutations used in the significant test.

- $fdr_threshold$ - The FDR threshold for significance.

# 2 Experimental Settings

The datasets used in the experiments are characterized in Lab Note 1. The computer used for this experiments used the 64-bit Ubuntu 13.10 operating system, with an Intel(R) Core(TM)2 Quad CPU Q6600 2.40GHz processor and 8,00 GB of RAM memory.

Team contains a C++ program that takes as parameters the genotype file, the phenotype file, the number of individuals, number of SNPs, number of permutations for the significance test and FDR threshold. The number of permutations is set to 100 and the FDR threshold is set to 1.

# 3 Results

The algorithm only outputs pairwise relations between SNPs. Due to this, only epistasis detection will be evaluated.

The Power observed in Figure 1 show a maximum value of 8% of Power for data sets with 500 individuals, 65% for 1000 individuals and 95% for 2000 individuals. There is a big correlation between the Power and the size of the data sets. However, for frequencies smaller than 0.1 there is a near 0% Power for most configurations, with the exception of data sets with 2000 individuals and 0.05 minor allele frequency. These values also increase with allele frequency, with the exception of 0.5 allele frequencies.
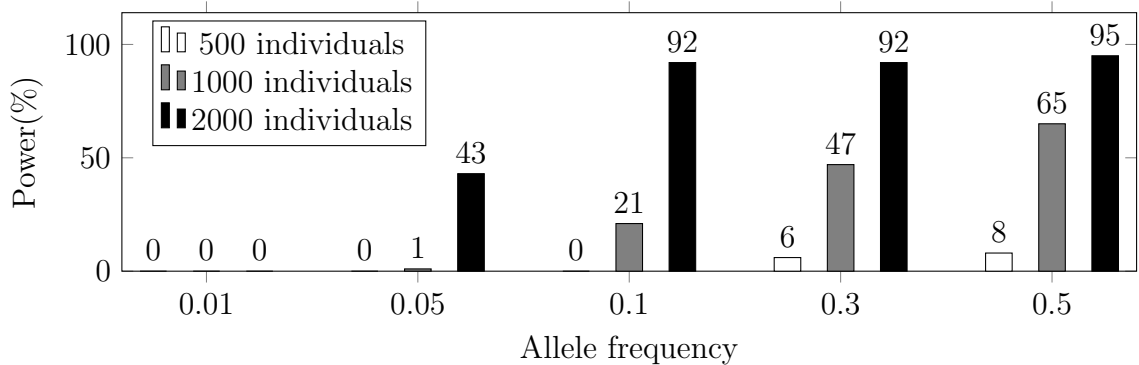


Figure 1: Power by allele frequency. For each frequency, three sizes of data sets were used to measure the Power, with an odds ratio of 2.0 and prevalence of 0.02. The Power is measured by the amount of data sets where the ground truth was amongst the most relevant results, out of all 100 data sets.

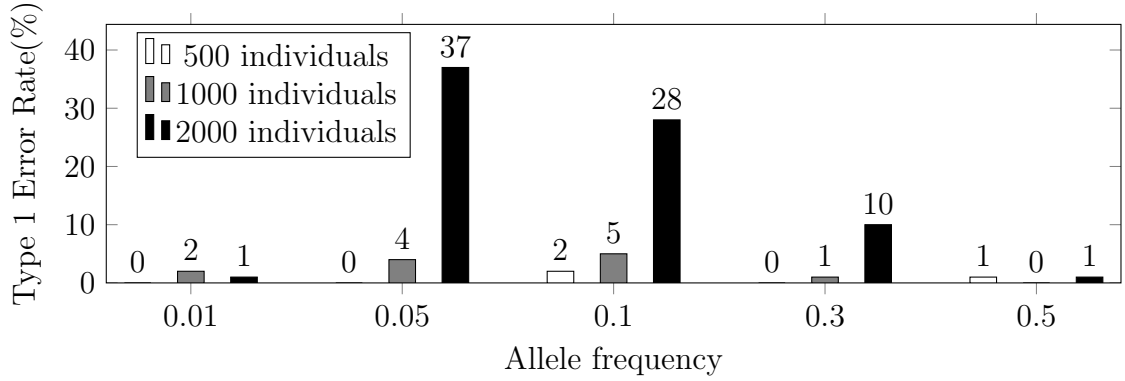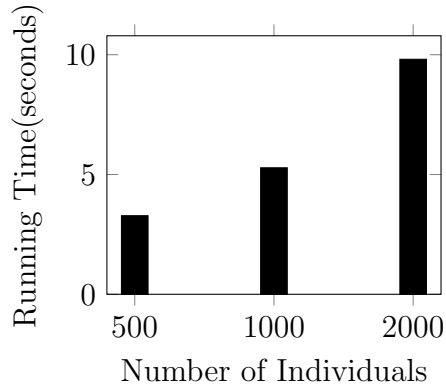The Type 1 Error Rate in Figure 2 has an interesting pattern, clearly

Figure 2: Type 1 Error Rate by allele frequency and population size. The Type 1 Error Rate is measured by the amount of data sets where the false positives were amongst the most relevant results, out of all 100 data sets.
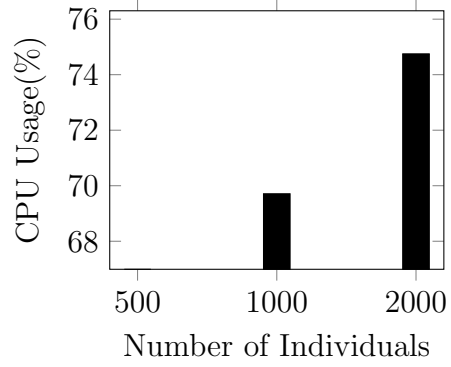
showing a growth in error rate with the population size. However, the error does not necessarily increase with allele frequency, revealing a maximum of 37% in data sets with 0.05 allele frequency and 2000 individuals. There is also a decrease in data sets with higher allele frequencies in data sets with 2000 individuals. Therefore the relation between error rate and allele frequency is undetermined.

There is a 10% difference in CPU usage (b) and 7 seconds in running time (a), with a maximum of 74% and 10 seconds, respectively. The memory usage increases from 162 MB to 228 MB, which is a 40% increase. The most relevant increase is the running time because the running time for 2000 individuals is the triple of the running time for 500 individuals, which is a problem for big data sets.
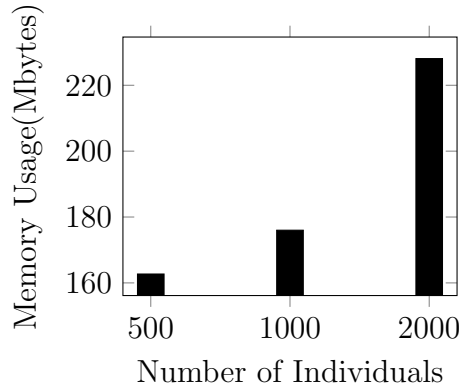
There is a clear increase in Power for odds ratio increase in Figure 5, especially from a 1.1 to a 1.5 odds ratio and population increase in Figure 4, with emphasis on the difference between 1000 and 2000 individuals The prevalence test shows very little difference between disease prevalence in Figure 6 and the allele frequency shows a growth with the increase in minor allele frequency.

(a) Average running time.

(b) Average CPU usage.



(c) Average memory usage.

Figure 3: Comparison of scalability measures between different sized data sets. The data sets have a minor allele frequency is 0.5, 2.0 odds ratio, 0.02 prevalence.

# 4  Summary

TEAM is an exhaustive algorithm that uses permutation tests to generate contingency tables, which then can be applied to any relevancy test. The results show an increase in Power related to the increase in population size and allele frequency. The scalability test shows that the running time of data sets with the highest population size is the triple of the running time for data sets with the lowest population size. The Type 1 Error Rate increases with the population size, but the relation between error rate and allele frequency is undetermined. The results of data set configurations by population and allele frequency confirm the previously discussed results. The odds ratio increase yields a clear increase in Power, but the prevalence increase shows nearly the same Power.

# References

[ZHZW10]  Xiang Zhang, Shunping Huang, Fei Zou, and Wei Wang. TEAM: efficient two-locus epistasis tests in human genome-wide association study. *Bioinformatics (Oxford, England)*, 26:i217–i227, 2010.

# A  Bar Graphs



Figure 4: Distribution of the Power by population. The allele frequency is 0.1, the odds ratio is 2.0, and the prevalence is 0.02.

Figure 5: Distribution of the Power by odds ratios. The allele frequency is 0.1, the number of individuals is 2000, and the prevalence is 0.02.
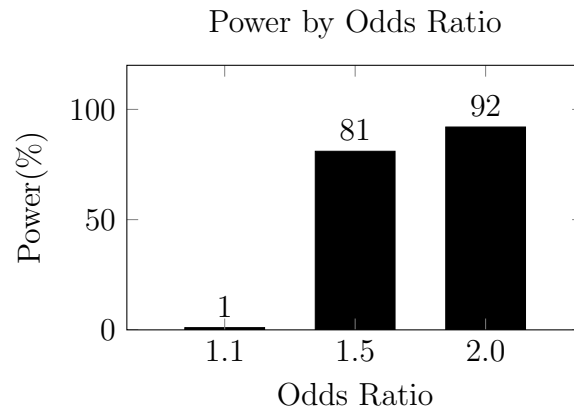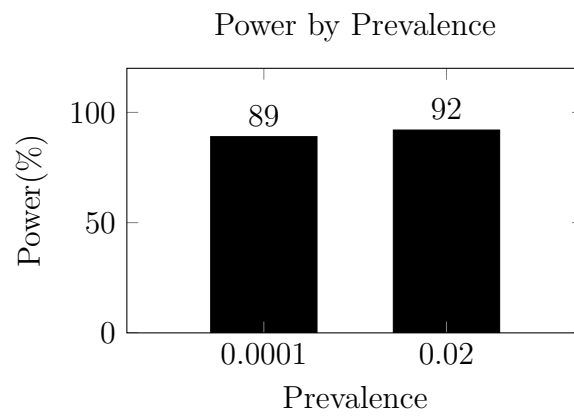


Figure 6: Distribution of the Power by prevalence. The allele frequency is 0.1, the number of individuals is 2000, and the odds ratio is 2.0.
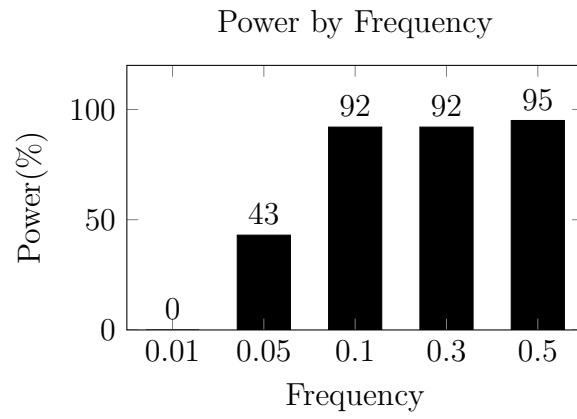
Figure 7: Distribution of the averaged Power by allele frequency. The number of individuals is 2000, the odds ratio is 2.0, and the prevalence is 0.02.

# B  Table of Results

Table 4: A table containing the percentage of true positives and false positives in each configuration. The first column contains the description of the configuration. The second and third columns contain the number of datasets with true positives and false positives respectively, out of all 100 data sets per configuration.

| Configuration* | TP (%) | FP (%) |
|---|---|---|
| 0.5,500,I,2.0,0.02 | 8 | 1 |
| 0.5,500,I,2.0,0.0001 | 9 | 0 |
| 0.5,500,I,1.5,0.02 | 3 | 0 |
| 0.5,500,I,1.5,0.0001 | 0 | 1 |
| 0.5,500,I,1.1,0.02 | 0 | 0 |
| 0.5,500,I,1.1,0.0001 | 0 | 3 |
| 0.5,2000,I,2.0,0.02 | 95 | 1 |
| 0.5,2000,I,2.0,0.0001 | 100 | 22 |
| 0.5,2000,I,1.5,0.02 | 93 | 7 |
| 0.5,2000,I,1.5,0.0001 | 47 | 1 |
| 0.5,2000,I,1.1,0.02 | 10 | 3 |
| 0.5,2000,I,1.1,0.0001 | 3 | 2 |
| 0.5,1000,I,2.0,0.02 | 65 | 0 |
| 0.5,1000,I,2.0,0.0001 | 79 | 5 |
| 0.5,1000,I,1.5,0.02 | 53 | 3 |
| 0.5,1000,I,1.5,0.0001 | 4 | 0 |
| 0.5,1000,I,1.1,0.02 | 0 | 1 |
| 0.5,1000,I,1.1,0.0001 | 0 | 3 |
| 0.3,500,I,2.0,0.02 | 6 | 0 |
| 0.3,500,I,2.0,0.0001 | 26 | 3 |
| 0.3,500,I,1.5,0.02 | 2 | 0 |
| 0.3,500,I,1.5,0.0001 | 5 | 1 |
| 0.3,500,I,1.1,0.02 | 0 | 0 |
| 0.3,500,I,1.1,0.0001 | 0 | 4 |
| 0.3,2000,I,2.0,0.02 | 92 | 10 |
| 0.3,2000,I,2.0,0.0001 | 100 | 56 |
| 0.3,2000,I,1.5,0.02 | 95 | 15 |
| 0.3,2000,I,1.5,0.0001 | 98 | 10 |
| 0.3,2000,I,1.1,0.02 | 2 | 1 |
| 0.3,2000,I,1.1,0.0001 | 2 | 4 |
| 0.3,1000,I,2.0,0.02 | 47 | 1 |

| | | |
|---|---|---|
| 0.3,1000,I,2.0,0.0001 | 100 | 12 |
| 0.3,1000,I,1.5,0.02 | 40 | 3 |
| 0.3,1000,I,1.5,0.0001 | 49 | 5 |
| 0.3,1000,I,1.1,0.02 | 0 | 0 |
| 0.3,1000,I,1.1,0.0001 | 0 | 1 |
| 0.1,500,I,2.0,0.02 | 0 | 2 |
| 0.1,500,I,2.0,0.0001 | 1 | 1 |
| 0.1,500,I,1.5,0.02 | 1 | 0 |
| 0.1,500,I,1.5,0.0001 | 0 | 1 |
| 0.1,500,I,1.1,0.02 | 0 | 1 |
| 0.1,500,I,1.1,0.0001 | 0 | 1 |
| 0.1,2000,I,2.0,0.02 | 92 | 28 |
| 0.1,2000,I,2.0,0.0001 | 89 | 20 |
| 0.1,2000,I,1.5,0.02 | 81 | 5 |
| 0.1,2000,I,1.5,0.0001 | 42 | 3 |
| 0.1,2000,I,1.1,0.02 | 1 | 3 |
| 0.1,2000,I,1.1,0.0001 | 5 | 2 |
| 0.1,1000,I,2.0,0.02 | 21 | 5 |
| 0.1,1000,I,2.0,0.0001 | 12 | 6 |
| 0.1,1000,I,1.5,0.02 | 5 | 0 |
| 0.1,1000,I,1.5,0.0001 | 1 | 1 |
| 0.1,1000,I,1.1,0.02 | 0 | 2 |
| 0.1,1000,I,1.1,0.0001 | 0 | 4 |
| 0.05,500,I,2.0,0.02 | 0 | 0 |
| 0.05,500,I,2.0,0.0001 | 0 | 1 |
| 0.05,500,I,1.5,0.02 | 0 | 0 |
| 0.05,500,I,1.5,0.0001 | 0 | 3 |
| 0.05,500,I,1.1,0.02 | 0 | 1 |
| 0.05,500,I,1.1,0.0001 | 0 | 1 |
| 0.05,2000,I,2.0,0.02 | 43 | 37 |
| 0.05,2000,I,2.0,0.0001 | 57 | 21 |
| 0.05,2000,I,1.5,0.02 | 40 | 24 |
| 0.05,2000,I,1.5,0.0001 | 43 | 19 |
| 0.05,2000,I,1.1,0.02 | 0 | 3 |
| 0.05,2000,I,1.1,0.0001 | 0 | 3 |
| 0.05,1000,I,2.0,0.02 | 1 | 4 |
| 0.05,1000,I,2.0,0.0001 | 3 | 2 |
| 0.05,1000,I,1.5,0.02 | 0 | 1 |

| | | |
|---|---|---|
| 0.05,1000,I,1.5,0.0001 | 1 | 3 |
| 0.05,1000,I,1.1,0.02 | 0 | 1 |
| 0.05,1000,I,1.1,0.0001 | 0 | 6 |
| 0.01,500,I,2.0,0.02 | 0 | 0 |
| 0.01,500,I,2.0,0.0001 | 0 | 2 |
| 0.01,500,I,1.5,0.02 | 0 | 0 |
| 0.01,500,I,1.5,0.0001 | 0 | 4 |
| 0.01,500,I,1.1,0.02 | 0 | 1 |
| 0.01,500,I,1.1,0.0001 | 0 | 0 |
| 0.01,2000,I,2.0,0.02 | 0 | 1 |
| 0.01,2000,I,2.0,0.0001 | 0 | 5 |
| 0.01,2000,I,1.5,0.02 | 0 | 3 |
| 0.01,2000,I,1.5,0.0001 | 0 | 4 |
| 0.01,2000,I,1.1,0.02 | 0 | 1 |
| 0.01,2000,I,1.1,0.0001 | 0 | 2 |
| 0.01,1000,I,2.0,0.02 | 0 | 2 |
| 0.01,1000,I,2.0,0.0001 | 0 | 1 |
| 0.01,1000,I,1.5,0.02 | 0 | 0 |
| 0.01,1000,I,1.5,0.0001 | 0 | 2 |
| 0.01,1000,I,1.1,0.02 | 0 | 2 |
| 0.01,1000,I,1.1,0.0001 | 0 | 3 |

*MAF,POP,MOD,OR,PREV where MAF represents the minor allele frequency, POP is the number of individuals, MOD is the used model (with or without main effect and with or without epistasis effect), OR is the odds ratio and PREV is the prevalence of the disease.

Table 5: A table containing the running time, cpu usage and memory usage in each configuration.

| Configuration* | Running Time (s) | CPU Usage (%) | Memory Usage (KB) |
|---|---|---|---|
| 0.5,500,I,2.0,0.02 | 3.28 | 66.99 | 166590.64 |
| 0.5,500,I,2.0,0.0001 | 3.81 | 54.75 | 166590.28 |
| 0.5,500,I,1.5,0.02 | 3.07 | 74.40 | 166590.44 |
| 0.5,500,I,1.5,0.0001 | 3.76 | 57.98 | 166590.60 |
| 0.5,500,I,1.1,0.02 | 3.08 | 68.52 | 161592.60 |
| 0.5,500,I,1.1,0.0001 | 3.91 | 55.00 | 166590.72 |
| 0.5,2000,I,2.0,0.02 | 9.81 | 74.75 | 233543.92 |
| 0.5,2000,I,2.0,0.0001 | 11.00 | 72.09 | 233802.28 |

| | | | |
|---|---|---|---|
| 0.5,2000,I,1.5,0.02 | 9.83 | 72.85 | 233535.72 |
| 0.5,2000,I,1.5,0.0001 | 10.98 | 66.89 | 233821.76 |
| 0.5,2000,I,1.1,0.02 | 9.82 | 73.74 | 233562.12 |
| 0.5,2000,I,1.1,0.0001 | 10.99 | 69.46 | 233832.84 |
| 0.5,1000,I,2.0,0.02 | 5.28 | 69.71 | 181210.92 |
| 0.5,1000,I,2.0,0.0001 | 6.08 | 68.02 | 181210.72 |
| 0.5,1000,I,1.5,0.02 | 5.53 | 66.72 | 181210.60 |
| 0.5,1000,I,1.5,0.0001 | 6.10 | 66.35 | 181210.64 |
| 0.5,1000,I,1.1,0.02 | 5.40 | 68.68 | 181210.64 |
| 0.5,1000,I,1.1,0.0001 | 6.09 | 65.91 | 181210.84 |
| 0.3,500,I,2.0,0.02 | 3.12 | 71.53 | 166590.44 |
| 0.3,500,I,2.0,0.0001 | 3.79 | 56.19 | 166590.60 |
| 0.3,500,I,1.5,0.02 | 3.13 | 70.93 | 166590.72 |
| 0.3,500,I,1.5,0.0001 | 3.77 | 56.60 | 166590.72 |
| 0.3,500,I,1.1,0.02 | 3.08 | 72.65 | 166590.68 |
| 0.3,500,I,1.1,0.0001 | 3.78 | 56.01 | 166590.40 |
| 0.3,2000,I,2.0,0.02 | 9.84 | 72.54 | 233557.00 |
| 0.3,2000,I,2.0,0.0001 | 10.94 | 73.45 | 233778.48 |
| 0.3,2000,I,1.5,0.02 | 9.92 | 72.03 | 233546.36 |
| 0.3,2000,I,1.5,0.0001 | 10.95 | 73.75 | 233801.92 |
| 0.3,2000,I,1.1,0.02 | 9.95 | 72.35 | 233546.48 |
| 0.3,2000,I,1.1,0.0001 | 11.00 | 70.49 | 233828.96 |
| 0.3,1000,I,2.0,0.02 | 5.34 | 67.05 | 181210.88 |
| 0.3,1000,I,2.0,0.0001 | 6.09 | 63.97 | 181210.80 |
| 0.3,1000,I,1.5,0.02 | 5.35 | 69.00 | 181210.56 |
| 0.3,1000,I,1.5,0.0001 | 6.12 | 63.37 | 181210.80 |
| 0.3,1000,I,1.1,0.02 | 5.44 | 67.27 | 181210.44 |
| 0.3,1000,I,1.1,0.0001 | 6.11 | 65.06 | 181210.68 |
| 0.1,500,I,2.0,0.02 | 3.28 | 65.33 | 166590.76 |
| 0.1,500,I,2.0,0.0001 | 3.78 | 56.18 | 166590.60 |
| 0.1,500,I,1.5,0.02 | 3.13 | 71.07 | 166590.60 |
| 0.1,500,I,1.5,0.0001 | 3.81 | 55.52 | 166590.84 |
| 0.1,500,I,1.1,0.02 | 3.22 | 67.56 | 166590.64 |
| 0.1,500,I,1.1,0.0001 | 3.84 | 54.26 | 166590.52 |
| 0.1,2000,I,2.0,0.02 | 9.91 | 72.77 | 233527.88 |
| 0.1,2000,I,2.0,0.0001 | 10.95 | 73.18 | 233788.92 |
| 0.1,2000,I,1.5,0.02 | 9.94 | 71.34 | 233538.28 |
| 0.1,2000,I,1.5,0.0001 | 10.97 | 71.25 | 233803.52 |

| | | | |
|---|---|---|---|
| 0.1,2000,I,1.1,0.02 | 9.82 | 69.45 | 231225.76 |
| 0.1,2000,I,1.1,0.0001 | 10.76 | 73.08 | 233841.76 |
| 0.1,1000,I,2.0,0.02 | 5.46 | 66.40 | 181210.92 |
| 0.1,1000,I,2.0,0.0001 | 6.10 | 65.14 | 181210.64 |
| 0.1,1000,I,1.5,0.02 | 5.41 | 67.52 | 181210.80 |
| 0.1,1000,I,1.5,0.0001 | 6.17 | 63.74 | 181210.80 |
| 0.1,1000,I,1.1,0.02 | 5.49 | 65.42 | 181210.52 |
| 0.1,1000,I,1.1,0.0001 | 6.25 | 57.76 | 181210.68 |
| 0.05,500,I,2.0,0.02 | 3.06 | 74.66 | 166590.52 |
| 0.05,500,I,2.0,0.0001 | 3.67 | 63.00 | 166590.84 |
| 0.05,500,I,1.5,0.02 | 3.10 | 73.32 | 166590.60 |
| 0.05,500,I,1.5,0.0001 | 3.70 | 60.99 | 166590.68 |
| 0.05,500,I,1.1,0.02 | 3.09 | 74.07 | 166590.96 |
| 0.05,500,I,1.1,0.0001 | 3.74 | 60.54 | 166590.84 |
| 0.05,2000,I,2.0,0.02 | 10.87 | 75.38 | 233762.72 |
| 0.05,2000,I,2.0,0.0001 | 10.88 | 76.33 | 233830.32 |
| 0.05,2000,I,1.5,0.02 | 9.76 | 75.67 | 233551.64 |
| 0.05,2000,I,1.5,0.0001 | 10.84 | 77.87 | 233818.16 |
| 0.05,2000,I,1.1,0.02 | 9.76 | 76.10 | 233559.48 |
| 0.05,2000,I,1.1,0.0001 | 10.89 | 78.26 | 233821.16 |
| 0.05,1000,I,2.0,0.02 | 5.45 | 69.61 | 181210.40 |
| 0.05,1000,I,2.0,0.0001 | 6.01 | 69.13 | 181210.88 |
| 0.05,1000,I,1.5,0.02 | 5.24 | 74.24 | 181210.68 |
| 0.05,1000,I,1.5,0.0001 | 6.04 | 68.74 | 181211.00 |
| 0.05,1000,I,1.1,0.02 | 5.34 | 71.82 | 181210.52 |
| 0.05,1000,I,1.1,0.0001 | 5.99 | 68.72 | 181210.64 |
| 0.01,500,I,2.0,0.02 | 3.13 | 72.69 | 166590.40 |
| 0.01,500,I,2.0,0.0001 | 3.69 | 60.83 | 166590.96 |
| 0.01,500,I,1.5,0.02 | 3.02 | 76.70 | 166590.72 |
| 0.01,500,I,1.5,0.0001 | 3.77 | 59.05 | 166590.52 |
| 0.01,500,I,1.1,0.02 | 3.20 | 69.81 | 166590.60 |
| 0.01,500,I,1.1,0.0001 | 3.72 | 60.54 | 166590.64 |
| 0.01,2000,I,2.0,0.02 | 9.70 | 77.18 | 233557.00 |
| 0.01,2000,I,2.0,0.0001 | 10.87 | 76.75 | 233813.88 |
| 0.01,2000,I,1.5,0.02 | 9.76 | 76.94 | 233554.60 |
| 0.01,2000,I,1.5,0.0001 | 10.83 | 77.23 | 233817.52 |
| 0.01,2000,I,1.1,0.02 | 9.81 | 81.09 | 233562.32 |
| 0.01,2000,I,1.1,0.0001 | 10.86 | 76.11 | 233839.32 |

| | | | |
|---|---|---|---|
| 0.01,1000,I,2.0,0.02 | 5.28 | 73.35 | 181210.76 |
| 0.01,1000,I,2.0,0.0001 | 6.00 | 69.20 | 181210.72 |
| 0.01,1000,I,1.5,0.02 | 5.35 | 71.68 | 181210.80 |
| 0.01,1000,I,1.5,0.0001 | 6.05 | 68.01 | 181210.36 |
| 0.01,1000,I,1.1,0.02 | 5.35 | 71.71 | 181210.84 |
| 0.01,1000,I,1.1,0.0001 | 5.99 | 68.05 | 181211.00 |

*MAF,POP,MOD,OR,PREV where MAF represents the minor allele frequency, POP is the number of individuals, MOD is the used model (with or without main effect and with or without epistasis effect), OR is the odds ratio and PREV is the prevalence of the disease.