# Analysis of population specific pharmacogenomic variants using next generation sequencing data

Eunyong Ahn[1,2], Taesung Park[1,3*]

[1]Interdisciplinary Program of Bioinformatics, Seoul National University, Seoul, 151-747, Korea.

[2]Department of Computer Science, Technion − Israel Institute of Technology, Haifa, 3200003, Israel.

[3]Department of Statistics, Seoul National University, Seoul, 151-747, Korea.

[*]Corresponding Authors:

Taesung Park, Department of Statistics, Seoul National University, 1 Gwanak-Ro Gwanak-Gu, Seoul, 151-747, Korea, Tel: +82-2-880-8924; E-mail: tspark@stats.snu.ac.kr.

**S1 Text. The dependence of $F_{st}$ on MAF**

Wu, *et al*. reported the PD among human genes using the HapMap data according to $F_{st}$ [1]. They called the genes of PD if they contained at least one SNP with an $F_{st}$ more than or equal to 0.6. However, it is not feasible that the $F_{st}$ from rare variants exceed 0.6 and this method could not detect the population differentiation in rare variants. S2 figure shows the $F_{st}$ in our WES data according to weighted average of MAF from each ancestry groups, and MAF from total population. The red line describes the maximum $F_{st}$ when the MAF is given. The maximum $F_{st}$ for MAF is defined when all the minor allele only exist in the ancestry group with the smallest sample size. If we assume there is no genotyping error, then we know how many loci are genotyped for each ancestry groups. In case of our data, the max $F_{st}$ was identified when American Hispanics (1938) only have minor allele and other ancestry groups such as African Americans (2025), East Asians (2164), South Asians (2199), and Europeans (4518) do not have minor allele at all. The maximum $F_{st}$ for the rare or less common variants (MAF < 0.05) is less than 0.36 and especially the maximum of $F_{st}$ for the rare variants (MAF < 0.01) is less than 0.073. Therefore, for these variants, very high divergence ($0.25 < F_{st}$) or high divergence ($0.15 < F_{st} < 0.25$) in the wright's criteria need to be modified to find PD in rare variants.

The upper and lower bound of $F_{st}$ when the MAF is given are already reported [2] and here we also showed that the maximum of $F_{st}$ is bounded according to MAF when MAF is small enough through the equation. We used the initial definition of $F_{st}$ which is developed by Wright for the simplification of proof. Since Wright developed this measure, many estimators has been proposed to estimate the $F_{st}$ correctly under various situations. However, we only choose the Wright's $F_{st}$ for our proof because other estimates are originated to estimate this parameter. Since Wright's $F_{st}$ assume the ideal

42    situation with infinite allele and balanced sample sizes, this ideal condition would be

43    different from the real world. However, we can assume the ideal condition theoretically,

44    and our proof will be able to be extended to other estimators.

45        $n_i$ denotes minor allele count of ancestry group i; N denotes the total genotyped

46    allele counts; k denotes the total minor allele counts in population; m denotes the

47    number of ancestry groups. Under Hardy Weinberg Equilibrium,

48        $H_S$ = mean expected heterozygosity within random mating subpopulations = $2\overline{p_i q_i}$

49        $H_T$ = expected heterozygosity in random mating total population = $2\overline{p}\,\overline{q}$

50    Wright's $F_{st} = \dfrac{H_T - H_s}{H_T}$

51        $= 1 - \dfrac{H_T - H_s}{H_T}$

52        $= 1 - \dfrac{\dfrac{2}{m}\sum\limits_{i=1}^{m}\left(\dfrac{n_i}{N/m}\left(1 - \dfrac{n_i}{N/m}\right)\right)}{2\left(\dfrac{k}{N}\right)\left(1 - \dfrac{k}{N}\right)}$

53        $= 1 - \dfrac{\sum\limits_{i=1}^{m} m\left(n_i\left(\dfrac{N}{m} - n_i\right)\right)}{k(N - k)}$

54    When k is less than N/m, $\exists n_i = k$ such that $i = 1 \cdots m$ and the minimum of nominator is

55    $m\left(k\left(\dfrac{N}{m} - k\right)\right)$. Since $\sum\limits_{i=1}^{m} n_i = k$, the minimum of nominator can be easily proved by

56    Jensen's inequality [3]. Therefore, when k is less than N/m, the maximum of above

57    equation is

58
$$= 1 - \frac{m\left(\dfrac{N}{m} - k\right)}{N - k}$$

59
$$= \frac{(m-1)k}{N-k}$$
$$= \frac{(1-m)(N-k) + (m-1)N}{N-k}$$
$$= (1-m) + \frac{(m-1)N}{N-k}$$

60 As k is increasing; the denominator is decreasing; the maximum of $F_{st}$ is increasing.

61 Therefore, when we focus on the variants with small MAF, then k is less than N/m and

62 the maximum of $F_{st}$ is bounded according to their MAF

63

64 **S2 Text. The permutation to confirm the distribution of PDRC**

65        Considering small p-values from real data analysis, the distribution of PDRC

66 may be claimed not to follow chi-square distribution. As an attempt to answer this issue,

67 we permuted the ancestral allele information of each 48 VIP genes for 100 thousand

68 times and calculated the PDRC statistics with three different weighting schemes. All

69 the variants are used for this permutation regardless of their MAFs. In Supplementary

70 S3 Fig, from several permuted data sets, the PDRC statistics without weight did not

71 follow the chi-square distribution, but the PDRC statistics seem to be controlled by

72 using weights as inverse of MAF, or inverse of $MAF^2$ (S4 Fig, and S5 Fig).

73

74 **S3 Text. The variance of common odds ratio**

$$\text{var}(\log(\widehat{\theta}_{MH_i})) = \frac{\sum_k w_k^2 \left(n_{i1k} + n_{52k}\right)\left(n_{i1k} n_{52k}\right) / \left(n_{i.k} + n_{5.k}\right)^2}{2\left(\sum_k w_k \left(n_{i2k} \cdot n_{51k}\right) / \left(n_{i.k} + n_{5.k}\right)\right)^2}$$

$$+ \frac{\sum_k w_k^2 \left(\left(n_{i1k} + n_{52k}\right)\left(n_{i2k} n_{51k}\right) + \left(n_{i2k} + n_{51k}\right)\left(n_{i1k} n_{52k}\right)\right) / \left(n_{i.k} + n_{5.k}\right)^2}{2\left(\sum_k w_k \left(n_{i1k} \cdot n_{52k}\right) / \left(n_{i.k} + n_{5.k}\right)\right)\left(\sum_k w_k \left(n_{i2k} \cdot n_{51k}\right) / \left(n_{i.k} + n_{5.k}\right)\right)}$$

$$+ \frac{\sum_k w_k^2 \left(n_{i2k} + n_{51k}\right)\left(n_{i2k} n_{51k}\right) / \left(n_{i.k} + n_{5.k}\right)^2}{2\left(\sum_k w_k \left(n_{i2k} \cdot n_{51k}\right) / \left(n_{i.k} + n_{5.k}\right)\right)^2}$$

76

## S4 Text. The simulation under the assumed null distribution

For the simulation to evaluate the type-1 error rate, we generated the data set without PD. We designed the data with four ancestral groups with 500 individuals and one ancestral group with 1000 individuals to assume a similar setting of sample sizes like our WES dataset. We also specified the number of SNPs in genes and assumed the distributions of MAF for a range of scenarios to investigate the potential effect of the number of variants in a gene and the distribution of MAF. In Scenarios 1 to 5, the rare or less common variants were generated; in Scenarios 6 to 10, the common variants were generated; in Scenarios 11 to 15, the variants were generated with the same MAF distribution as our WES data. The ancestral group information was randomly assigned to follow the null hypothesis, and $10^5$ genes are simulated under fifteen different scenarios. For each setting of MAF distribution, 5 different numbers of SNP in a gene are assumed as following, 5, 10, 20, 50, and 100. Let $\text{Gene}_{i,s}$ represent the $i$th Gene of scenario $s$ for $i \in \{1, 2, \cdots, 99999, 100000\}$ and $SNP_{i,j,s}$ the $j$th SNP in $\text{Gene}_{i,s}$ for

91    $j \in \{1, 2, \cdots, n_s\}$ . The MAF of $SNP_{i,j,s}$ , $p_{i,j,s}$ , is sampled from $unif(0, 0.05)$ for

92    scenarios 1 to 5, and from $unif(0.05, 0.5)$ for scenarios 6 to 10. For scenarios 11 to 15,

93    $p_{i,j,s}$ is sampled from the MAF distribution of all 3,130,381 variants in our WES data

94    sets (S1 Fig). Also, it is known that the *p*-values from the  null distribution follow the

95    uniform distribution and the distribution of simulated *p*-values can be investigated by

96    QQ-plot. Supplementary Figures S6 to S8 show that the p-value of the PDRC test from

97    the simulated data sets follows the uniform distribution when using three types of

98    weights, 1, inverse of MAF, and inverse of MAF$^2$. According to these results, the type-

99    1 error rate of PDRC tests seems to be reasonably controlled regardless the MAF values

100   and the number of variants in a gene.

101

102   1      Wu, D. D. & Zhang, Y. P. Different level of population differentiation among human genes.
103          *Bmc Evol Biol* **11**, 16, doi:10.1186/1471-2148-11-16 (2011).
104   2      Jakobsson, M., Edge, M. D. & Rosenberg, N. A. The relationship between F(ST) and the
105          frequency    of    the    most    frequent    allele.    *Genetics*    **193**,    515-528,
106          doi:10.1534/genetics.112.144758 (2013).
107   3      Hölder, O. Ueber einen Mittelwertsatz. *Göttinger Nachr*, 38–47 (1889).

108

109   **Legends to Supplementary Figures**

110

111   **S1 Fig.** Histogram of $\log_{10}(MAF)$ from our WES data sets

112

113   **S2 Fig.** The maximum of $F_{st}$ is bounded according to the MAF from total population

114   The red line represents theoretical maximum of $F_{st}$.

115

**S3 Fig.** QQ-plot results of PDRC without weight from the ancestral allele information permutation of VIP gene datasets

**S4 Fig.** QQ-plot results of PDRC with the weight as inverse of MAF from the ancestral allele information permutation of VIP gene datasets

**S5 Fig.** QQ-plot results of PDRC with the weight as inverse of MAF$^2$ from the ancestral allele information permutation of VIP gene datasets

**S6 Fig.** QQ-plot results from the simulation under null hypothesis 1

MAF of $SNP_{i,j,s}$, $p_{i,j,s}$, is sampled from $unif(0, 0.05)$ for scenarios 1 to 5

**1:** Number of SNPs in a Gene is 5. **2**: Number of SNPs in a Gene is 10. **3**: Number of SNPs in a Gene is 20. **4**: Number of SNPs in a Gene is 50. **5**: Number of SNPs in a Gene is 100. **A**: No weight, **B**: Weight is 1/MAF, **C**: Weight is 1/MAF$^2$

**S7 Fig.** QQ-plot results from the simulation under null hypothesis 2

MAF of $SNP_{i,j,s}$, $p_{i,j,s}$, is sampled from $unif(0.05, 0.5)$ for scenarios 6 to 10

**6**: Number of SNPs in a Gene is 5. **7**: Number of SNPs in a Gene is 10. **8**: Number of SNPs in a Gene is 20. **9**: Number of SNPs in a Gene is 50. **10**: Number of SNPs in a Gene is 100. **A**: No weight, **B**: Weight is 1/MAF, **C**: Weight is 1/MAF$^2$

137    **S8 Fig.** QQ-plot results from the simulation under null hypothesis 3

138    MAF of $SNP_{i,j,s}$, $p_{i,j,s}$, is sampled from the real MAF distribution of our WES data for
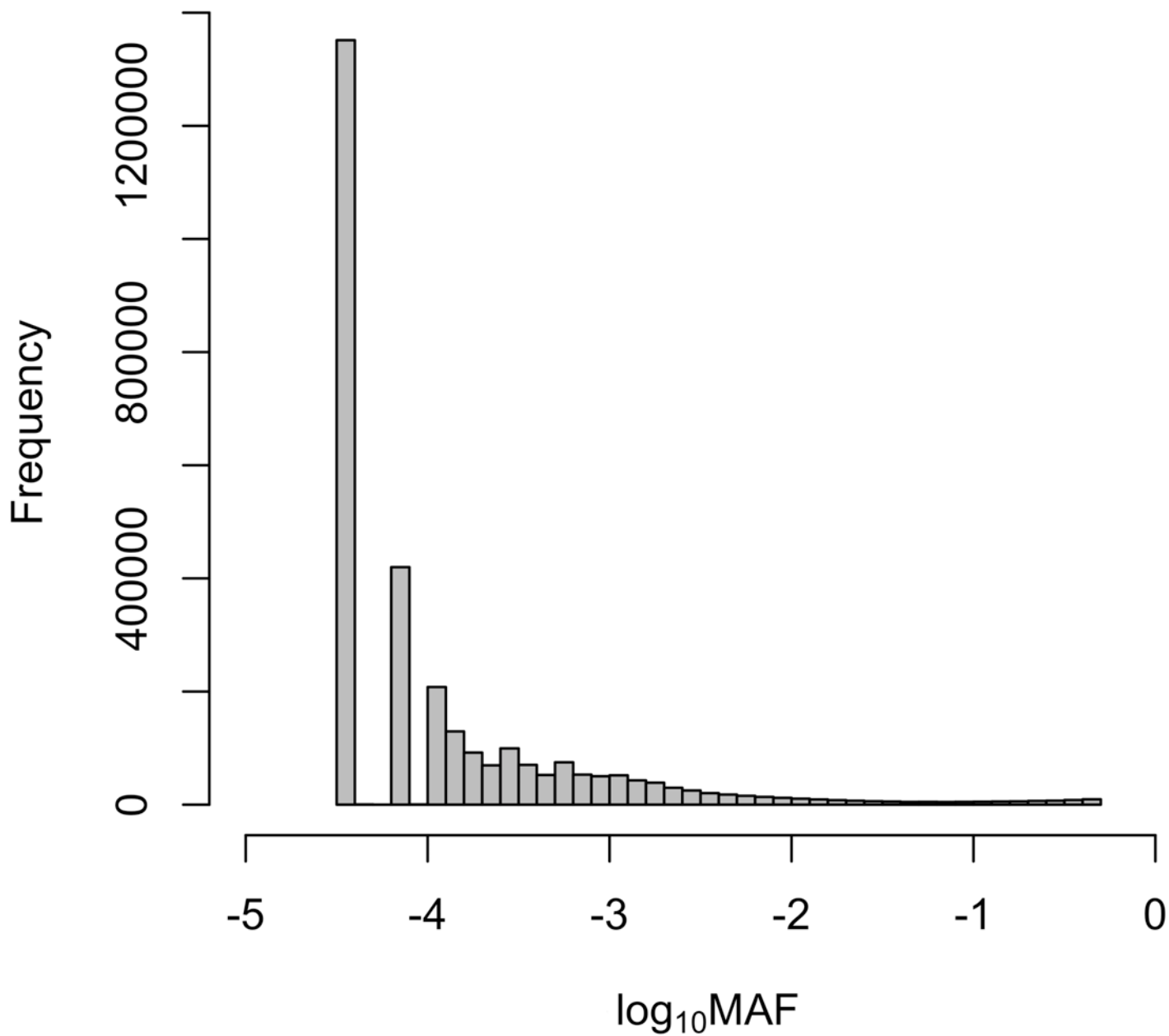
139    scenarios 11 to 15

140    **11**: Number of SNPs in a Gene is 5. **12**: Number of SNPs in a Gene is 10. **13**: Number

141    of SNPs in a Gene is 20. **14**: Number of SNPs in a Gene is 50. **15**: Number of SNPs in
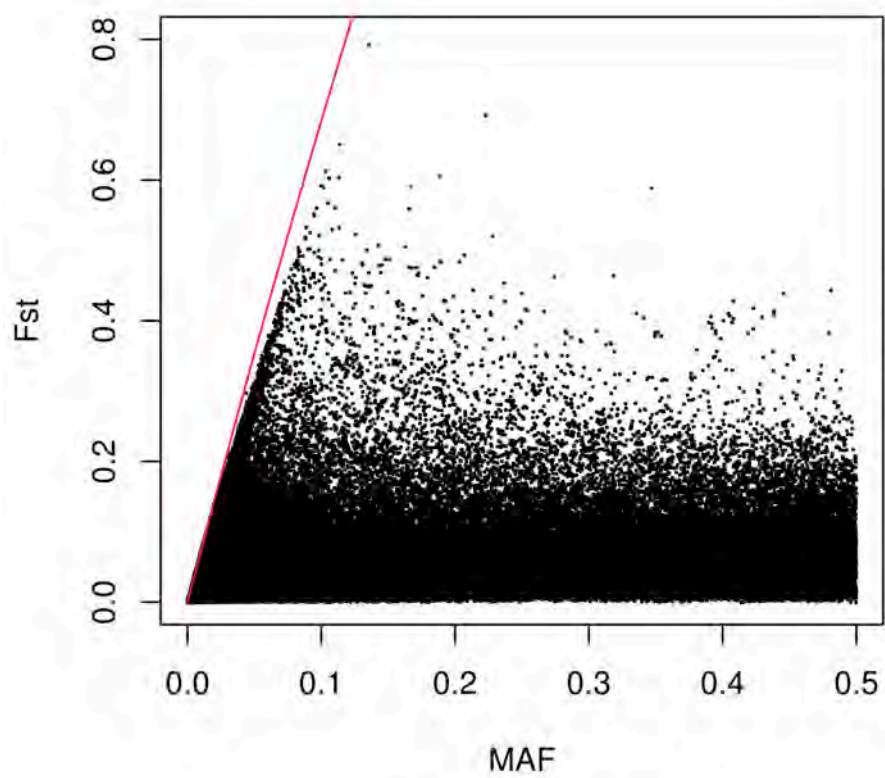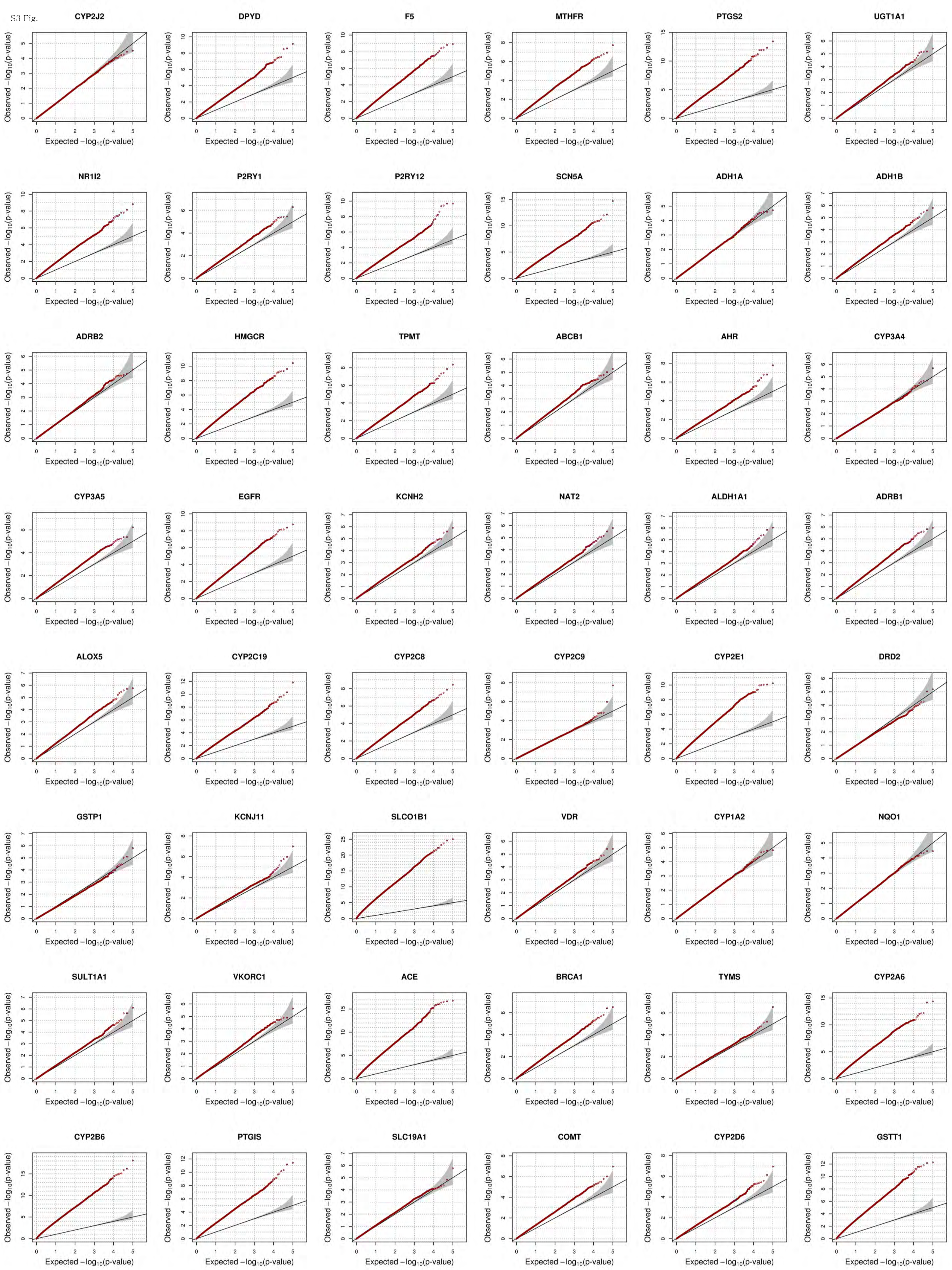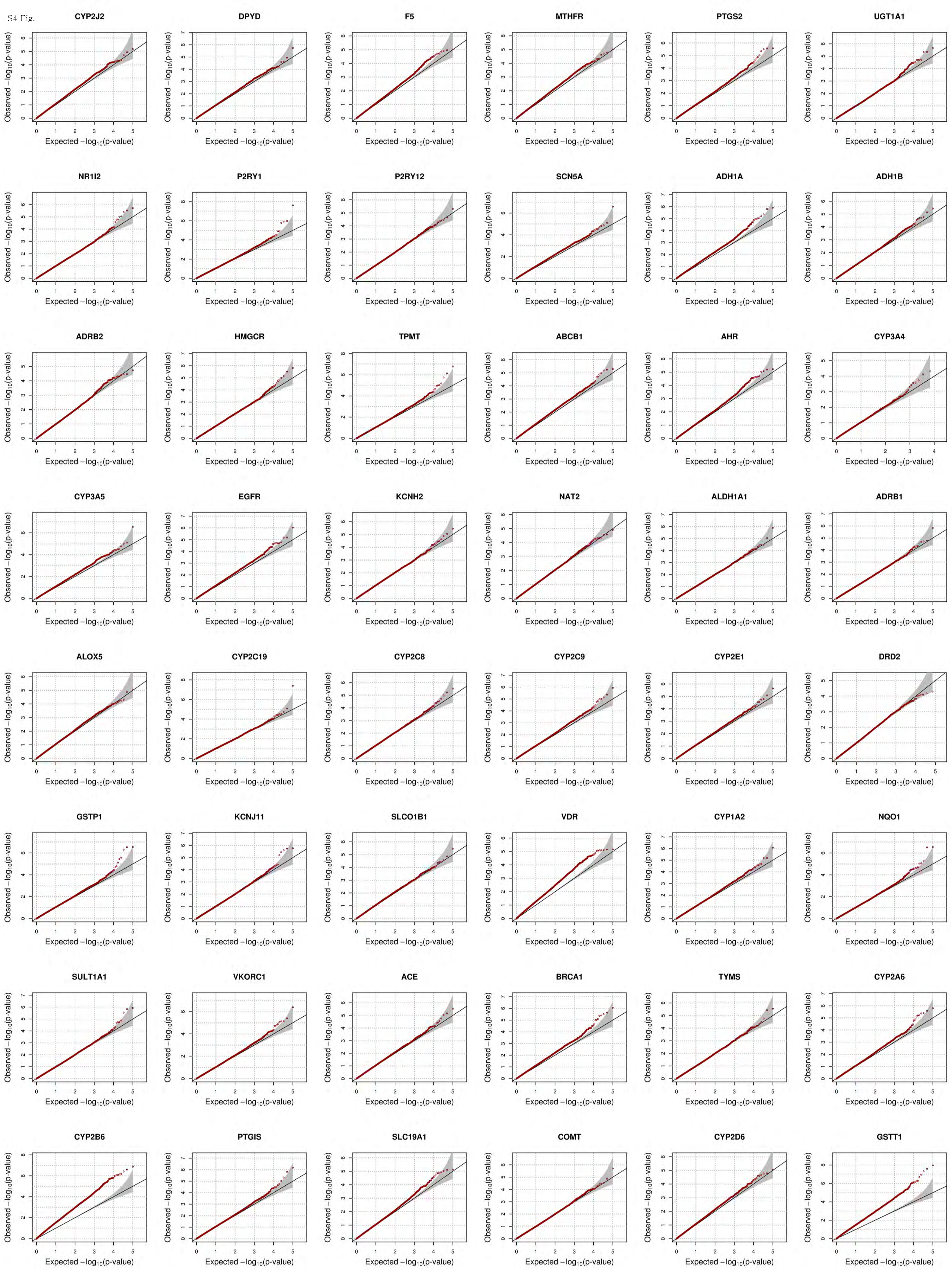
142    a Gene is 100. **A**: No weight, **B**: Weight is 1/MAF, **C**: Weight is 1/MAF$^2$

S1 Fig.



Histogram of $\log_{10}$MAF

S2 Fig.

S6 Fig.
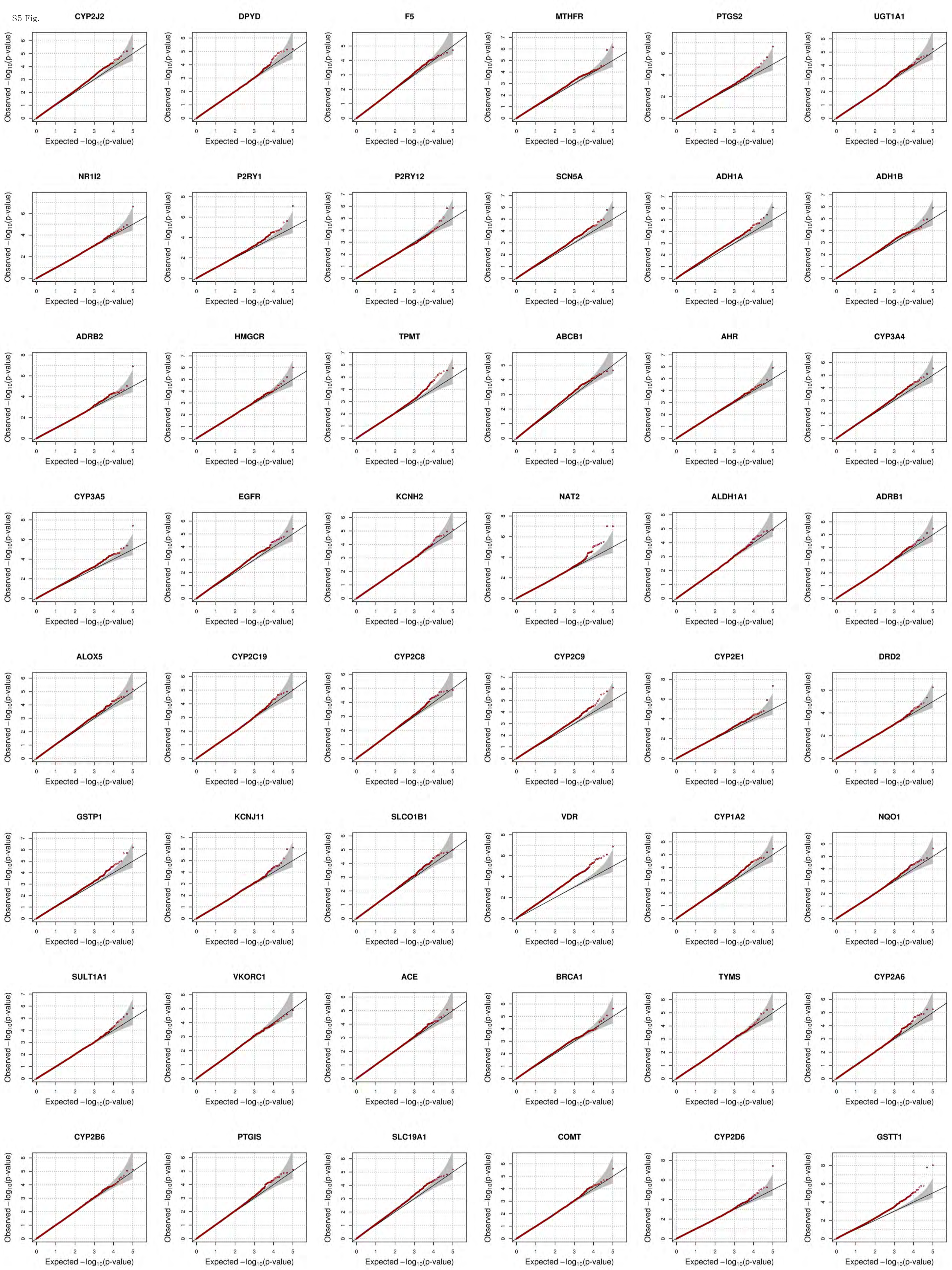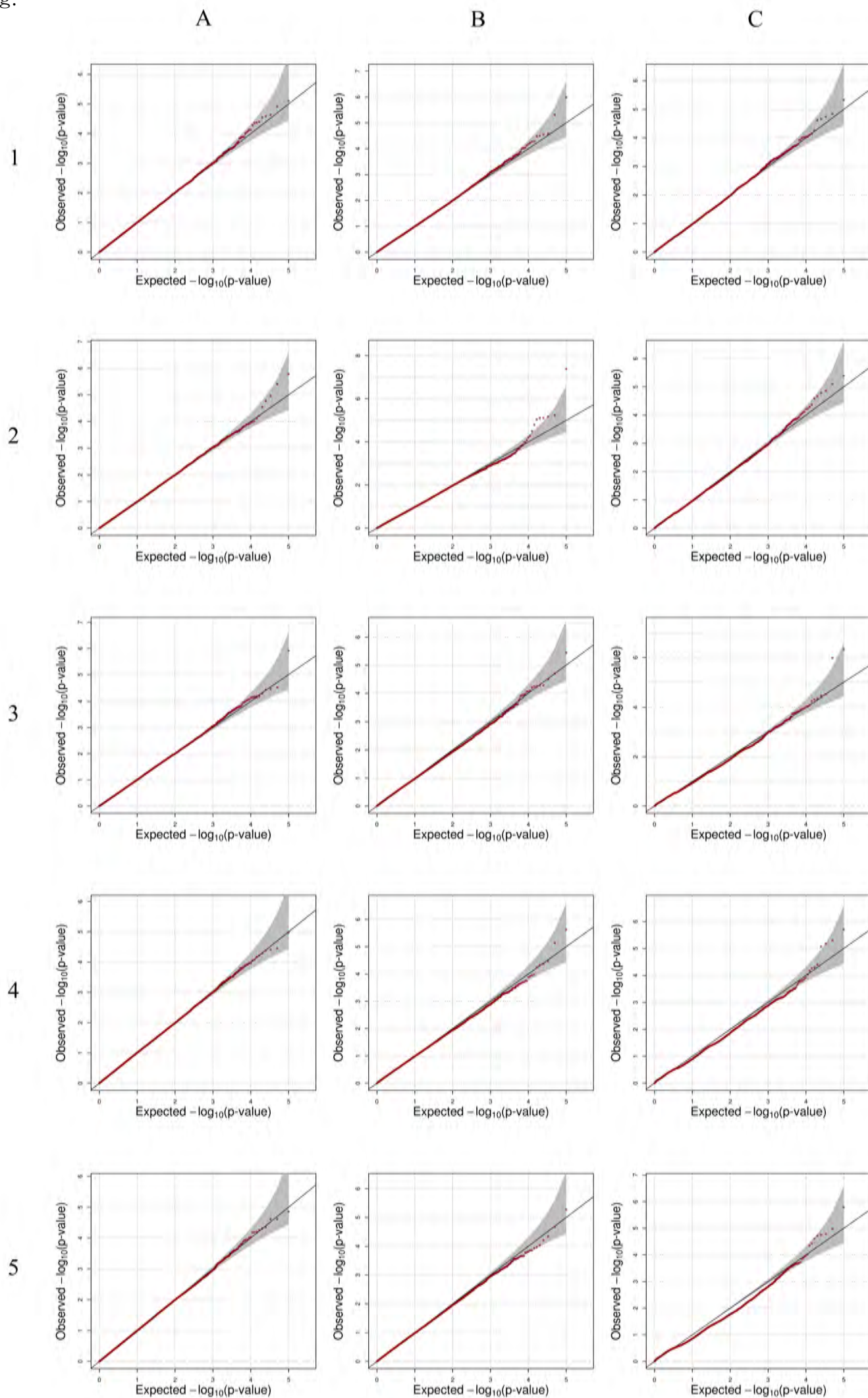
S7 Fig.

A               B               C

6

7

8

9

10

S8 Fig.