

3-4-2014

# Comparing Partial Least Square Approaches in Gene-or Region-based Association Study for Multiple Quantitative Phenotypes

Zhongshang Yuan

*Department of Epidemiology and Biostatistics, School of Public Health, Shandong University*

Xiaoshuai Zhang

*Department of Epidemiology and Biostatistics, School of Public Health, Shandong University*

Fangyu Li

*Department of Epidemiology and Biostatistics, School of Public Health, Shandong University*

Jinghua Zhao

*MRC Epidemiology Unit and Institute of Metabolic Science, Addenbrooke's Hospital, Cambridge, UK*

Fuzhong Xue

*Department of Epidemiology and Biostatistics, School of Public Health, Shandong University, xuefzh@sdu.edu.cn*

---

## Recommended Citation

Yuan, Zhongshang; Zhang, Xiaoshuai; Li, Fangyu; Zhao, Jinghua; and Xue, Fuzhong, "Comparing Partial Least Square Approaches in Gene-or Region-based Association Study for Multiple Quantitative Phenotypes" (2014). *Human Biology Open Access Pre-Prints*. Paper 50.

[http://digitalcommons.wayne.edu/humbiol\\_preprints/50](http://digitalcommons.wayne.edu/humbiol_preprints/50)

**Comparing Partial Least Square Approaches in Gene-or Region-based  
Association Study for Multiple Quantitative Phenotypes**

Zhongshang Yuan,<sup>1</sup> Xiaoshuai Zhang,<sup>1</sup> Fangyu Li,<sup>1</sup> Jinghua Zhao,<sup>2</sup> and

Fuzhong Xue<sup>1</sup>

<sup>1</sup>Department of Epidemiology and Biostatistics, School of Public Health,  
Shandong University.

<sup>2</sup>MRC Epidemiology Unit and Institute of Metabolic Science, Addenbrooke's  
Hospital, Cambridge, UK.

\*Correspondence to: Fuzhong Xue, Department of Epidemiology and  
Biostatistics, School of Public Health, Shandong University. E-mail:  
xuefzh@sdu.edu.cn.

**Key words: gene- or region-based association study, multiple quantitative  
traits, partial least square, thinking quantitatively for complex disease.**

**Abstract.** On thinking quantitatively of complex diseases, there are at least  
three statistical strategies for association study: single SNP on single trait,  
gene-or region (with multiple SNPs) on single trait and on multiple traits. The

third of which is the most general in dissecting the genetic mechanism underlying complex diseases underpinning multiple quantitative traits. Gene-or region association methods based on partial least square (PLS) approaches have been shown to have apparent power advantage. However, few attempts are developed for multiple quantitative phenotypes or traits underlying a condition or disease, and the performance of various PLS approaches used in association study for multiple quantitative traits had not been assessed. We, from regression perspective, exploit association between multiple SNPs and multiple phenotypes or traits through exhaustive scan statistics (sliding window) using PLS and sparse PLS (SPLS) regression. Simulations are conducted to assess the performance of the proposed scan statistics and compare them with the existed method. The proposed methods are applied to 12 regions of GWAS data from the European Prospective Investigation of Cancer (EPIC)-Norfolk study.

The well-documented successes (Stranger et al. 2011) in genome-wide association studies (GWASs) have greatly improved our understanding of the genetic architecture of complex traits. However, a notable issue with GWAS using case-control design, though discussed extensively in the literature, still needs to be addressed (Zhang and Liu 2007; Gayan et al. 2008). Such a design

is usually furnished through division of a particular quantitative phenotype into case and control groups with a cutoff that may or may not relate to the underlying genetic variation. The issue stems from incomplete penetrance commonly seen in complex traits such that individuals not showing any symptoms selected as controls are not far from being cases. A number of examples can be given. The first of these is diabetes, where an individual can be defined as being diabetic (case) when the fasting plasma glucose is  $\geq 7$  mmol/L or non-diabetic (control) otherwise (Rowe et al. 2000). In this case, many of the so-called 'controls' can almost be cases phenotypically and genetically (say fast plasma glucose=6.8mmol/L). It is apparent that assigning an artificial cutoff to a continuous trait runs a risk of information loss, for the phenotypically similar individuals on either side of the cutoff are treated as being phenotypically different. In general, phenotypes of most complex diseases (obesity, hypertension, diabetes, etc.) are effectively quantitative. It is important to appreciate this point for which we will dwell on a bit more. First, it should be noted that quantitative traits sometimes can only be acquired through collection of a large sample from the population and that a well-designed case-control study can achieve comparable statistical power with smaller sample size when variances in the cases and controls differ greatly. Second, in the framework of treating common disorders as

quantitative traits (Plomin et al. 2009), a complex disease is due to multiple genes with small effects, gene-gene and gene-environment interactions, Third, there may be obvious quantitative traits for some diseases as in the cases of body mass index (BMI)-obesity, blood pressure- hypertension, and mood-depression relationships, but the relevant quantitative traits may not be entirely clear for others in the cases of arthritis, autism, cancers, dementia and heart disease due to limited availability of biomarkers. As for obesity, BMI is understood to be merely a proxy since it crudely measures the mean body weight under given body surface area and varies with the amount of body fat but not its distribution. It has been shown that people with abdominal fat (with more weight around the waist) face more risks of cardiovascular diseases (Donahue and Abbott 1987; Ducimetiere et al. 1986) and other related diseases (hypertension, type 2 diabetes, and high cholesterol) than those with more fat around the hip (Bjorntorp 1988; Yusuf et al. 2005; Wells 2007), suggesting that a single measurement (BMI) may be ineffective marker of disease (obesity).

The consideration for quantitative trait(s) can lead to three statistical strategies for SNP-trait association: namely single SNP-single trait, gene/region (with multiple SNPs)-single trait, and gene/region-multiple traits. Among these, the first strategy is most susceptible to high false positive rate

and low power in detecting modest effects owing to the ignorance of the linkage disequilibrium (LD) (Beyene et al. 2009; Buil et al. 2009). The second strategy, though may alleviate the problems of multiple testing and facilitate stable results and interpretation (Lo et al. 2008; Qiao et al. 2009), fails to account for multiple quantitative traits. The last strategy is most general and deserves more attention.

Many gene- and region-based association tests have been developed, which include but do not limit to haplotype-based (Tregouet et al. 2009), P-value or odds ratio combination (Yang et al. 2009; Li et al. 2009), PCA-based (Peng et al. 2010), scan statistic (sliding window) (Sun et al. 2009; Hoh et al. 2000) and partial least squares (PLS)-based methods (Zhang et al. 2011; Turkmen et al. 2011; Chun et al. 2011; Xue et al. 2012). Among these, PLS-based methods have been shown to have apparent advantage in statistical power over others. Zhang et al. (2011) proposed a PLS regression-based multilocus association study for single quantitative traits, considering only the first component in their simulations. Turkmen and Lin (2011) considered methods to aggregate the signals of many SNPs within a gene for possible genetic effects from rare variants. As association may be contaminated by irrelevant markers in a gene or a region, Chun et al. (2011) considered sparse PLS (SPLS) regression for dimension reduction and derivation of components

that are linear combinations of only relevant markers. Furthermore, Xue et al. (2012) introduced a partial least squares path modeling (PLSPM) framework for association between single or multiple SNPs and a latent trait that can involve single or multiple correlated phenotype(s), which naturally provides estimators of polygenic effect by appropriately weighting trait-attributing alleles. However, few attempts were made for multiple quantitative phenotypes or traits underlying a disease, and performances of various PLS approaches used in association study for multiple quantitative traits have not been assessed. Much work needs to be done in this respect since the final result may vary greatly with the way to summarize the association on a gene level involving multiple markers and the approach to consider multiple phenotypes and their correlation structure to achieve efficiency and validity.

We therefore exploited association between multiple SNPs and multiple phenotypes via exhaustive scan statistics (sliding window) using PLS and SPLS regressions. Simulations were conducted to assess the performance of the newly proposed scan statistics and compare them with existing PLSPM method (Xue. et al. 2012). The methods were then applied to 12 regions of GWAS data from the European Prospective Investigation of Cancer (EPIC)-Norfolk study (Loos et al. 2008).

## **Material and Methods**

**Study Samples.** The EPIC-Norfolk study was a population-based, ethnically homogeneous, white Europe cohort study of 25,631 residents living in the city of Norwich, United Kingdom, and its surrounding area. Participants were 39-79 years old during the baseline health check between 1993 and 1997. A case-cohort study was conducted in which 3867 individuals were assayed with Affymetrix 500K genechips among whom subcohort (N=2,566) was a random sample of the study cohort at baseline and cases were part of the remaining individuals with BMI bigger than  $30 \text{ kg/m}^2$  (N=1,301). We analyzed the 2,417 individuals in the subcohort who passed the quality control and had complete genotype data for 446,861 SNPs on the whole genome.

**The Modeling Framework.** We used the exhaustive scan statistics based on PLS and SPLS regressions for multiple phenotypes (Figure 1). A major concern over the sliding window approach related to how to determine the optimal window size. A large window may include too many non-informative markers while a small window may miss those which are informative, both reducing statistical power. We here employed a brute-force search strategy with variable window sizes to alleviate this problem, which is likely to be feasible with a multiprocessor and multithreading computing environment.



Specifically, the strategy started with a pre-set largest window size  $L$  and was followed by exhausting search of the candidate region from the first SNP with sliding-window of all possible sizes  $s$  ranging from 1 to  $L$ . For a given window size (rectangle in Figure 1), a PLS or SPLS regression was used to detect association.

Cross validation was used to tune the number of latent components in PLS regression and the two parameters in SPLS regression (the number of latent components and the sparsity). Unlike the traditional multivariate hypothesis testing, the statistic obtained from PLS and SPLS regression does not follow Wilks' lambda distribution since the latent component was derived from multiple phenotype information and therefore permutation test was used to assess the statistical significance by comparing the observed statistic to its empirical distribution generated from 5,000 permutations with permuted phenotypes. For each permuted data, we used the same number of components from the original data for both PLS and SPLS regressions, but different sparsity parameter tuned from the permuted data for SPLS (Chun et al. 2011).

## **Simulation**

**Design.** Simulation was conducted as follows: (1) HapMap phase II CEU data at the brain-derived neurotrophic factor (*BDNF*) region (Chr

11:27633610..27692970 with 31 SNPs) was chosen to generate simulated genotypes. The pair-wise  $r^2$  are shown in Figure 2. (2) From (1), 500,000 individuals were generated via software gs 2.0 (Li and Chen 2008) with the 10th SNP (minor allele frequency, MAF=36.7%) being the causal variant; (3) Multiple quantitative traits data was generated from a trivariate  $Y = (Y_1, Y_2, Y_3)$  normal distribution  $Y \sim N(\mu, \Sigma)$ , where the three variables formed the random vector (waist, hip, BMI) for ‘‘apple-shaped’’ types ( $N = 355$ ) in EPIC-Norfolk GWAS subcohort with sample mean  $\bar{Y} = (105.2746, 106.0051, 29.2172)$  and sample covariance

$$S = \begin{pmatrix} 52.1991 & 36.8688 & 16.9545 \\ 36.8688 & 37.1419 & 13.7969 \\ 16.9545 & 13.7969 & 8.3859 \end{pmatrix}$$

For the case that the causal SNP had no effect ( $H_0$ ), let

$\mu = (105.2746, 106.0051, 29.2172)$  to be invariant with all three genotypes ( $GG$ ,  $GA$ , and  $AA$ ), while for  $H_1$ , the causal SNP was assumed to have effects on waist but not hip with the single allele effect size on BMI  $\delta$ ,

$\mu = (105.2746, 106.0051, 29.2172 + i\delta)$ , where  $i = 0, 1, 2$  for  $GG$ ,  $GA$  and  $AA$ ,

respectively. The range of  $\delta = (0.10, 0.15, 0.20, 0.25, 0.30)$  was estimated by

published data on genetic predisposition score (Li et al. 2010). Using the same

‘‘apple-shaped’’ data in the EPIC-Norfolk GWAS, estimation of waist under

fixed hip was obtained by

$$\widehat{waist} = 10.20345 + 0.62138 \cdot hip + 0.99947 \cdot BMI (p < 0.0001, R^2 = 0.7635); (4)$$

Genotypic data were simulated under variable sample sizes ( $N = 500, 1000, 1500, 2000, 2500$ ) from the simulated CEU population (500000 individuals), and quantitative genetics models with a given  $\delta$  generated by the R *mvtnorm* package. The window size was set to 10 SNPs from the 7th to the 16th SNP. Meanwhile, simulations with the 10th SNP and the 19th SNP (minor allele frequency, MAF=35%) being the causal variants have also been conducted, where one SNP have effect on waist and another have same effect on BMI. 1000 simulations were conducted under various sample sizes and various single allele effect sizes  $\delta$  to assess the type I error and power.

To further investigate the performance of the proposed method on different LD structures in a gene or region, we also chose another two regions from *FTO* gene in the simulation, one in high LD (Figure 3a) and the other in low LD (Figure 3b) while keeping the window size and the MAF of the causal SNP (30% for high LD, 35.8% for low LD) nearly the same as above.

Moreover, we removed the causal SNP to be in line with the more practically indirect association in all simulations. All procedures were implemented under Linux and involved R *plspm*, *pls*, *spls* and *mvtnorm* packages all available from CRAN (<http://cran.r-project.org/>).

**Type I Error.** Results from the simulation under the null hypothesis are shown in Table 1, indicating that type I error rates of PLSPM, PLS1 (only capturing the first component in PLS regression), PLS and SPLS were close to given nominal values ( $\alpha = 0.05$ ) given different sample sizes.

**Power.** Shown in Figure 4 are the statistical power under four different scenarios, regarding under different sample sizes at a given effect size  $\delta = 0.25$  (4a), under different effect sizes but a fixed sample size 2,000 (4b), with two causal SNPs (4c,d). As expected, the power was monotonically increasing functions of sample size and effect size for all the models. For a given sample size and a given effect size, PLS and SPLS regressions had comparable power while power for PLS1 regression, though higher than PLSPM, was not as high since it only extracted the first component. Table 2 gives results involving both high and low LD under a fixed sample size of 1,000 and an effect size of 0.25. SPLS and PLS generally had comparable power which were higher than other methods when there was high LD but all methods lost power when there was low LD. Besides, simulations given different correlation structures of phenotypes and causal SNP with lower MAF (8th SNP, MAF=17%) were considered and shown to have similar results

(data not shown), though the power for each method was not quite high owing to the lower MAF of causal SNP.

**Computational Efficiency.** The success of a GWAS requires computational efficiency and feasibility. Different PLS-based methods had different theoretical basis, and the computing time can be influenced by many factors such as numerical algorithm, sample size, window size, coding language and number of permutations required. Taking our multiprocessor and multithreading computational cluster as an example, 100 jobs can be submitted at a time (10 nodes and each with 10 concurrent processes). For a typical 500,000 SNP GWAS and a fixed window size, nearly 500,000 scan statistics only required 5,000 times. Table 3 provides the estimated computing time in hours for a typical whole genome scan under a window size of 10 and 1,000 permutations, with each node as Intel Xeon 5620 with 2.4GHz CPU and 16 GB RAM.

## **Application**

We employed both the proposed exhaustive scan statistics and the existing PLSPM to the EPIC-Norfolk data involving 12 genomic regions (*NEGR1*, *SEC16B*, *TMEM18*, *ETV15*, *GNPDA2*, *BDNF*, *MTCH2*, *SH2B1*, *FAIM2*, *FTO*,

*MC4R, KCTD15*) and three obesity-related phenotypes waist, hip and BMI with sliding windows sized between 1 and 12 SNPs, 1000000 permutations were conducted to obtain the empirical  $p$  value. The proposed method for the single trait was also used for comparisons.

One region rs7204609~rs9939881 at the first intron of *FTO* gene was found to have strong association ( $p = 2.86 \times 10^{-6}$ , window size=10) by the PLS regression with the three correlated traits (waist, hip, BMI) as response variables, compared to  $p = 3.8 \times 10^{-4}$  for waist,  $p = 1.1 \times 10^{-3}$  for hip and  $p = 7.6 \times 10^{-4}$  for BMI as single traits. For this region, the results for SPLS were  $p = 1.2 \times 10^{-5}$  for three multiple traits,  $p = 0.55, 0.276, 0.214$  for waist, hip and BMI respectively. In contrast the results by PLSPM was  $p = 8 \times 10^{-5}$  for three multiple traits,  $p = 0.1, 0.283, 0.5$  for waist, hip and BMI respectively.

## **Discussion**

Under the hypothesis of thinking quantitatively (Plomin et al. 2009), we have considered a general framework for association study on quantitative phenotype, which includes single SNP on single trait, gene or region (each can involve multiple SNPs) on single trait, and gene or region on multiple traits as the most reasonable in underlying genetic mechanism involving multiple quantitative phenotypes for complex diseases. We exploited the association

through exhaustive scan statistics using PLS and SPLS, and compared the performances of various PLS-based approaches. Simulations based on real data from the EPIC-Norfolk study indicated that, under a variety of scenarios SPLS and PLS had comparable power higher than other methods when there was high or moderate LD, while all methods lost power when the LD was low. Furthermore, the power for PLS1 regression was not so high because it only extracted the first component and it was still more powerful than PLSPM. Xue et al. (2012) showed in their simulation that the PLSPM had good power when the causal SNP had not been removed, while the causal SNP had been removed in our simulation which is a more practical scenario involving indirect association between the traits and a genomic region. Similar PLSPM results were obtained if the causal SNP was kept (results not shown).

The computational efficiency and feasibility of an implementation were function of the numerical algorithm, sample size, window size, the coding language and the number of permutations required. We only provided estimates for computing time under a fixed window size and 1,000 permutations (the threshold being  $p = 1 \times 10^{-3}$ ) in R. In a real-world GWAS, given that only a few genes or regions may be related to the traits in question, one can first use a somewhat lower threshold (e.g.  $1 \times 10^{-3}$ ) to screen potential trait-related genes or regions, to be followed by a higher threshold (e.g.  $1 \times 10^{-5}$ )

for further analysis on all genes or regions screened. SPLS essentially embeds variable selection into the PLS regression and would be preferred in order to alleviate the influence from unrelated variants when the region contains many variants, notwithstanding its highest computational burden among all the methods.

Analysis of the EPIC-Norfolk data suggested that the scan statistics for multiple quantitative traits were more powerful than those for single trait with the window size 10 providing the strongest evidence and in agreement with the literature regarding the optimality of 10-SNP window (Tregouet et al. 2009). In particular, the region (rs7204609~rs9939811) detected within the first intron 1 of *FTO* gene was also identified earlier (Xue et al. 2012).

In conclusion, PLS and SPLS are valid and powerful gene-or region-based association method for multiple quantitative phenotypes. However, how to obtain its theoretical distribution rigorously remains a challenge for which further work is warranted.

**Acknowledgements.** The EPIC-Norfolk study is supported by research programme grant funding from Cancer Research UK and the Medical Research Council. This work was supported by grants from National Natural Science Foundation of China (31071155), grants from the China Postdoctoral



Science Foundation (2011M501147) and Young Talents Innovation

Foundation of School of Public Health, Shandong University. We are very grateful of participants of the EPIC-Norfolk GWAS and other colleagues for their support. We wish to thank the associate Editor and anonymous referees for their comments which help to improve the presentation greatly.

### **Literature Cited**

Beyene, J., D. Tritchler, J. Asimit et al. 2009. Gene- or region-based analysis of genome-wide association studies. *Genetic Epidemiology* 33:S105–S110.

Bjorntorp, P. 1988. Abdominal obesity and the development of noninsulin-dependent diabetes mellitus. *Diabetes/Metabolism Reviews* 4:615–622.

Buil, A., A. Martinez-Perez, A. Perera-Lluna et al. 2009. A new gene-based association test for genome-wide association studies. *BMC Proceedings* 3 (Suppl. 7):S130–S130.

Chun, H., D. Ballard, J. Cho et al. 2011. Identification of association between disease and multiple markers via sparse partial least-squares regression. *Genet. Epidemiol.* 35:479–486.

Donahue, R., and R. Abbott. 1987. Central obesity and coronary heart disease

in men. *Lancet* 2:1,215.

Ducimetiere, P., J. Richard, and F. Cambien. 1986. The pattern of subcutaneous fat distribution in middle-aged men and the risk of coronary heart disease: The Paris Prospective Study. *International Journal of Obesity* 10:229–240.

Gayán, J., A. Gonzalez-Perez, F. Bermudo et al. 2008. A method for detecting epistasis in genome-wide studies using case-control multi-locus association analysis. *BMC Genomics* 9:360.

Hoh, J., and J. Ott. 2000. Scan statistics to scan markers for susceptibility genes. *Proc. Natl. Acad. Sci. USA* 97:9,615–9,617.

Li, J., and Y. Chen. 2008. Generating samples for association studies based on HapMap data. *BMC Bioinformatics* 9: 44.

Li, M., K. Wang, S. Grant et al. 2009. ATOM: A powerful genebased association test by combining optimally weighted markers. *Bioinformatics* 25:497–503.

Li, S., J. Zhao, J. Luan et al. 2010. Cumulative effects and predictive value of common obesity-susceptibility variants identified by genome-wide association studies. *Am. J. Clin. Nutr.* 91:184–190.

Lo, S., H. Chernoff, L. Cong et al. 2008. Discovering interactions among BRCA1 and other candidate genes associated with sporadic breast cancer. *Proceedings of the National Academy of Sciences.* 105:12,387–12,392.

- Loos, R., C. Lindgren, S. Li et al. 2008. Common variants near MC4R are associated with fat mass, weight and risk of obesity. *Nat Genet.* 40:768–775.
- Peng, Q., J. Zhao, and F. Xue. 2010. PCA-based bootstrap confidence interval tests for gene-disease association involving multiple SNPs. *BMC Genet.* 11:6.
- Plomin, R., C. Haworth, and O. Davis, O. 2009. Common disorders are quantitative traits. *Nat. Rev. Genet.* 10:872–878.
- Qiao, B., C. Huang, L. Cong et al. 2009. Genome-wide gene based analysis of rheumatoid arthritis-associated interaction with PTPN22 and HLA-DRB1. *BMC Proc.* 3(Suppl. 7):S132.
- Rowe, N., P. Mitchell, R. Cumming et al. 2000. Diabetes, fasting blood glucose and age-related cataract: The Blue Mountains Eye Study. *Ophthalmic Epidemiology* 7:103–114.
- Stranger, B., E. Stahl, and T. Raj. 2011. Progress and promise of genome-wide association studies for human complex trait genetics. *Genetics* 187:367–383.
- Sun, Y., D. Jacobsen, S. Turner et al. 2009. A fast implementation of a scan statistic for identifying chromosomal patterns of genome wide association studies. *Computational Statistics & Data Analysis* 53:1,794–1,801.
- Tregouet, D., I. Konig, J. Erdmann et al. 2009. Genome-wide haplotype

- association study identifies the SLC22A3-LPAL2-LPA gene cluster as a risk locus for coronary artery disease. *Nat. Genet.* 41:283–285.
- Turkmen, A., S. Lin. 2011. Gene-based partial least-squares approaches for detecting rare variant associations with complex traits. *BMC Proc.* 5 (Suppl. 9):S19.
- Wells, J. 2007. BMI compared with 3-dimensional bodyshape: The UK National Sizing Survey. *Am. J. Clin. Nutr.* 85:7.
- Willer, C., E. Speliotes, R. Loos et al. 2009. Six new loci associated with body mass index highlight a neuronal influence on body weight regulation. *Nat. Genet.* 41:25–34.
- Xue, F., S. Li, J. Luan et al. 2012. A latent variable partial least squares path modeling approach to regional association and polygenic effect with applications to a human obesity study. *PLoS ONE* 7:e31927.
- Yang, H., Y. Liang, C. Chung et al. 2009. Genome-wide gene-based association study. *BMC Proc.* 3 (Suppl. 7):S135.
- Yusuf, S., S. Hawken, S. Ounpuu et al. 2005. Obesity and the risk of myocardial infarction in 27,000 participants from 52 countries: A case-control study. *Lancet* 366:1,640–1,649.
- Zhang, F., X. Guo, and H.-W. Deng. 2011. Multilocus association testing of quantitative traits based on partial least-squares analysis. *PLoS ONE* 6:e16739.

Zhang, Y., and J. Liu. 2007. Bayesian inference of epistatic interactions in case-control studies. *Nature Genetics* 39:1,167–1,173.

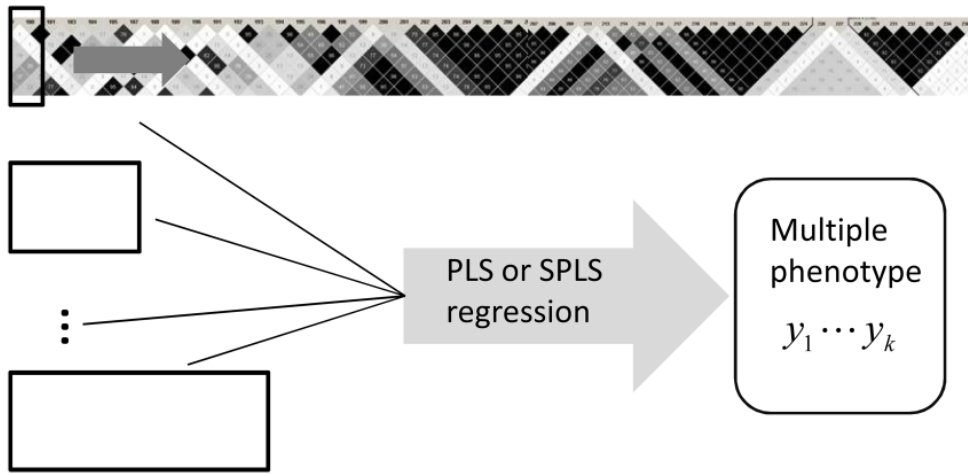


Figure 1. Modeling framework of PLS or SPLS regression-based exhaustive scan statistics. Rectangle represents a fixed window size.

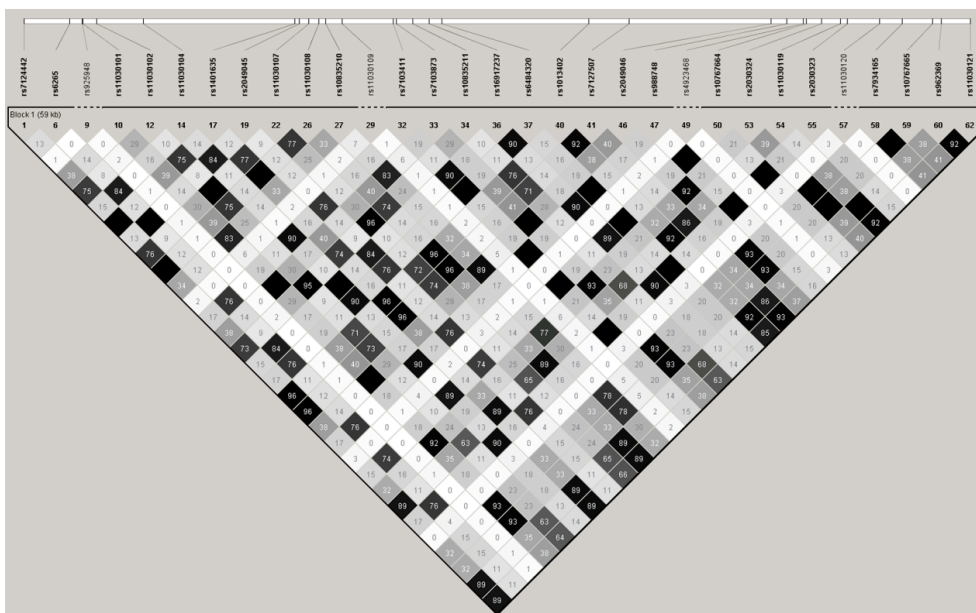


Figure 2. Pair-wise  $r^2$  pattern among the 31 SNPs in the selected region.

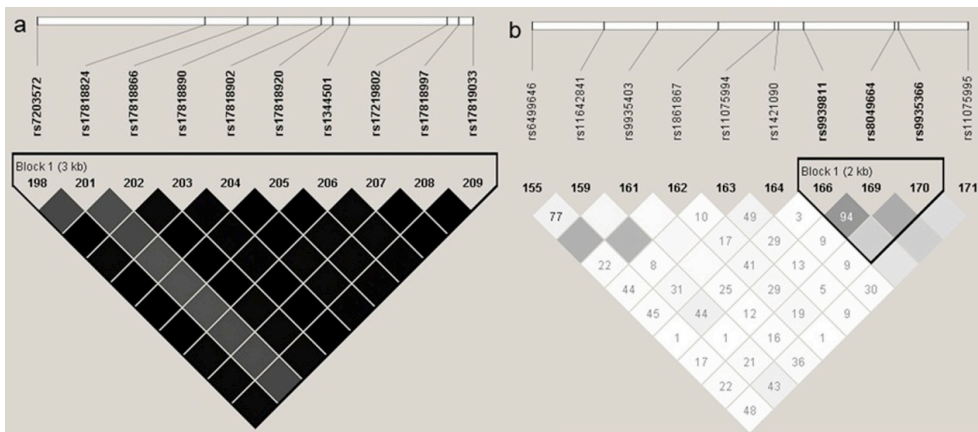


Figure 3. LD patterns for two regions in the *FTO* gene. (a) region in high LD.

(b) region in low LD.

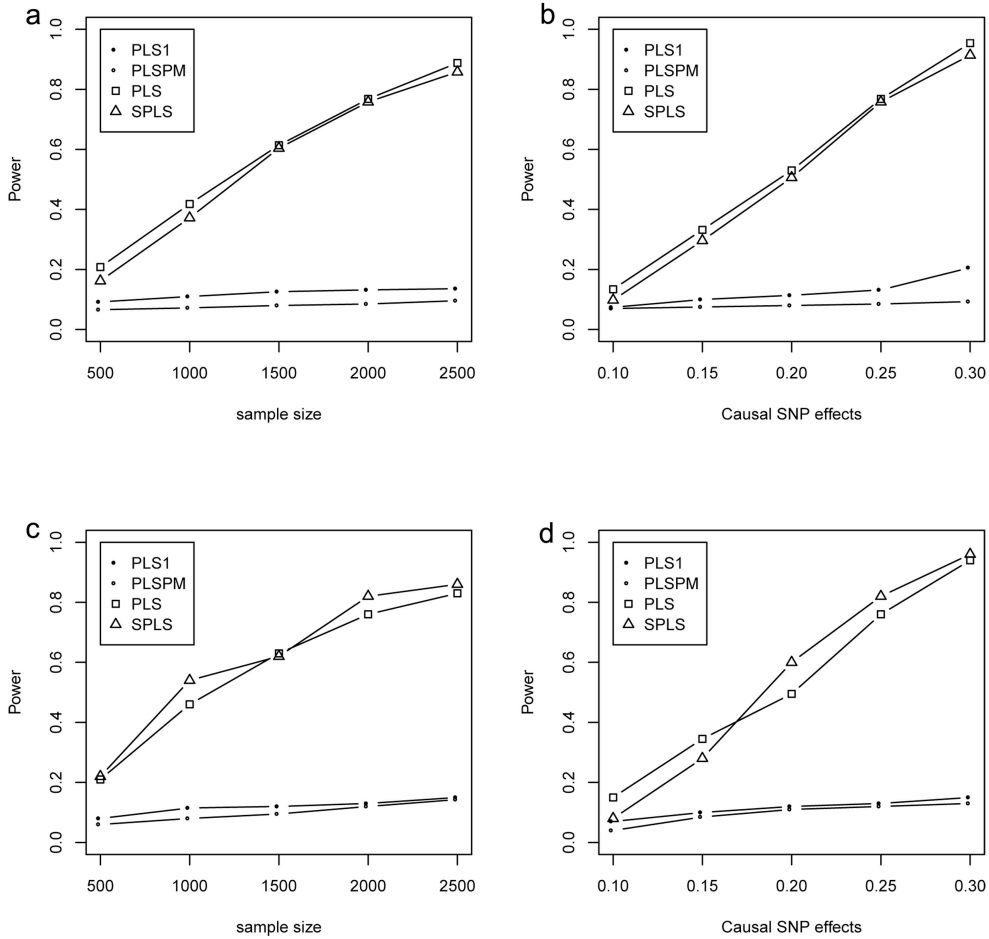


Figure 4. Power of PLSPM, PLS1, PLS and SPLS under four scenarios. (a) power under different sample sizes but a given effect size of  $\delta = 0.25$  for one causal SNP. (b) power under different effect sizes but a fixed sample size of 2,000 for one causal SNP. (c) power under different sample sizes but a given effect size of  $\delta = 0.25$  for two causal SNPs. (d) power under different effect sizes but a fixed sample size of 2,000 for two causal SNPs



Table 1. Type I error as a function of sample size ( $\alpha = 0.05$ )

<b>Sample size</b>	<b>PLSPM</b>	<b>PLS1</b>	<b>PLS</b>	<b>SPLS</b>
500	0.058	0.058	0.046	0.061
1000	0.046	0.060	0.044	0.046
1500	0.046	0.062	0.048	0.054
2000	0.050	0.048	0.054	0.058
2500	0.044	0.060	0.054	0.060

Table 2. Power in relation to LD for given a sample size of 1,000 and an effect size of 0.25

	<b>PLSPM</b>	<b>PLS1</b>	<b>PLS</b>	<b>SPLS</b>
<b>High LD region</b>	0.178	0.112	0.532	0.558
<b>Low LD region</b>	0.070	0.058	0.074	0.064

Table 3. Estimated computing time (hours) for a whole genome scan for a fixed window size of 10 and 1,000 permutations

<b>Sample size</b>	<b>PLSPM</b>	<b>PLS1</b>	<b>PLS</b>	<b>SPLS</b>
500	20	12.5	15	240
1,000	40	30	55	421

1,500	81	70	90	680
2,000	165	160	176	840
2,500	260	251	255	1,000

---