

RESEARCH ARTICLE

Open Access



# Adjusting heterogeneous ascertainment bias for genetic association analysis with extended families

Suyeon Park<sup>1,2,3</sup>, Sungyoung Lee<sup>4</sup>, Young Lee<sup>1,2</sup>, Christine Herold<sup>5,6</sup>, Basavaraj Hooli<sup>7</sup>, Kristina Mullin<sup>7</sup>, Taesung Park<sup>8</sup>, Changsoon Park<sup>1</sup>, Lars Bertram<sup>7,9,10</sup>, Christoph Lange<sup>6,11,12,13</sup>, Rudolph Tanzi<sup>7\*</sup> and Sungho Won<sup>14,15,16\*</sup>

## Abstract

**Background:** In family-based association analysis, each family is typically ascertained from a single proband, which renders the effects of ascertainment bias heterogeneous among family members. This is contrary to case-control studies, and may introduce sample or ascertainment bias. Statistical efficiency is affected by ascertainment bias, and careful adjustment can lead to substantial improvements in statistical power. However, genetic association analysis has often been conducted using family-based designs, without addressing the fact that each proband in a family has had a great influence on the probability for each family member to be affected.

**Method:** We propose a powerful and efficient statistic for genetic association analysis that considered the heterogeneity of ascertainment bias among family members, under the assumption that both prevalence and heritability of disease are available. With extensive simulation studies, we showed that the proposed method performed better than the existing methods, particularly for diseases with large heritability.

**Results:** We applied the proposed method to the genome-wide association analysis of Alzheimer's disease. Four significant associations with the proposed method were found.

**Conclusion:** Our significant findings illustrated the practical importance of this new analysis method.

**Keywords:** Family-based association analysis, Ascertainment, Liability model

## Background

Genome-wide association studies (GWASs) have been used to identify many genes involved in human diseases, and during the last decade, many disease-susceptibility variants have been identified. However, despite these successes, we have found that variants discovered from GWASs often explain only a small proportion of the heritability of diseases [1, 2]. For example, SNPs significantly associated with human height explain only about 5 % of phenotypic variance, despite studies of tens of thousands of people [3]. Many reasons, such as rare causal variants and gene/gene interactions, have been attributed

to this so-called "missing heritability". However, the low power induced by the multiple-testing problem is still an intractable issue in GWASs, and further investigations of the most efficient strategies for genetic association analysis are necessary.

Careful selection of samples based on phenotypes can lead to improved power for the discovery of risk variants [4–11]. One such example is the extreme discordant sib-pair design in linkage analysis, which may result in a substantial increase in statistical power when compared to other sib-pair designs [11, 12]. Similarly, ascertaining the extremes of quantitative phenotypes from large population cohorts has also been shown to increase the power to identify associated variants [13–15]. In such a design, the effect of ascertainment conditions are homogeneous between individuals, and existing methods, such as the Cochran-Armitage(CA) trend test [16], can be an efficient choice. However, in association analysis using extended

\* Correspondence: rtanzi@mgh.harvard.edu; won1@snu.ac.kr

<sup>7</sup>Genetics and Aging Research Unit, MassGeneral Institute for Neurodegenerative Diseases, Department of Neurology, Massachusetts General Hospital, Harvard Medical School, Massachusetts, USA

<sup>14</sup>Department of Public Health Science, Seoul National University, Seoul, Korea

Full list of author information is available at the end of the article

families, the effects of ascertainment bias are often heterogeneous among family members, and depending on their relationships with probands, different magnitudes of ascertainment bias may be generated. In particular, the probability of each individual being affected when his or her relatives are affected is similar to the prevalence, if the heritability is small, which indicates that the heterogeneous effect of the ascertainment bias depends on the magnitude of heritability. However, the heterogeneous effects of ascertainment conditions and the influence of heritability on it have not yet been investigated, and should therefore be taken into account for association analysis.

Recently, the CA trend test was extended for association analysis of dichotomous phenotypes with family-based samples [17, 18]. These statistics compares the genotype frequencies between affected and unaffected individuals, and the genetic association with family-based samples is tested by building a genotype correlation matrix with either kinship coefficients or an empirical correlation matrix estimated from large-scale genetic data. This approach has been extended to include family members with known phenotypes and missing genotypes or *vice versa*. By the nature of these statistics, it performs well for ascertained family-based samples and it can be an efficient choice, even for a case-control design, if the relatives' phenotype information is available. However, their statistical efficiency is affected by the heterogeneous effect of the ascertainment bias on family members, and for extended families, its effects on statistical efficiency can be substantial.

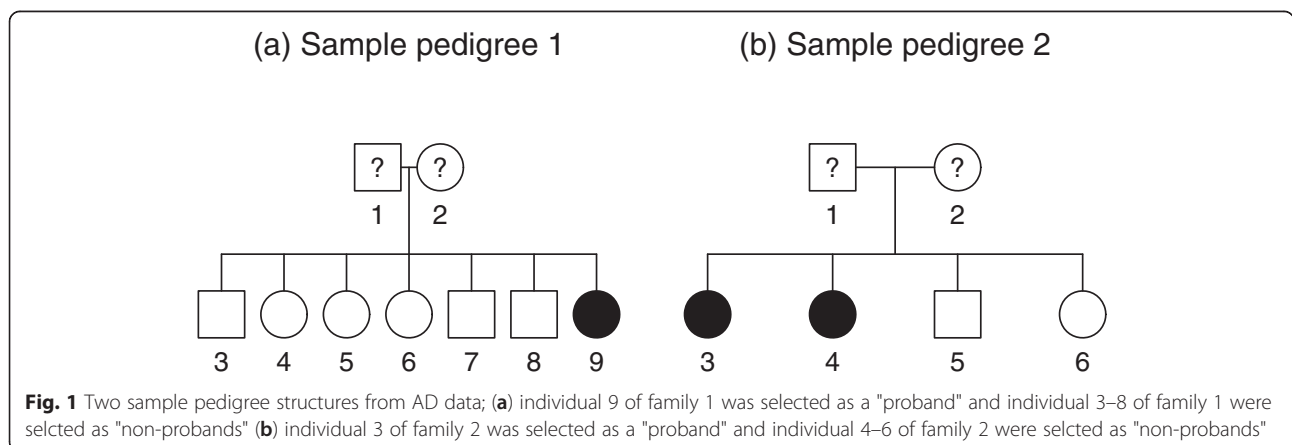
In this report, we consider the heterogeneous effects of the ascertainment bias on family members for dichotomous phenotypes. By the nature of the proposed methods, individuals with missing genotypes and non-missing phenotypes can be utilized, and incorporation of the estimated kinship matrix to the proposed statistic provided robustness against the population substructure. The proposed method consists of two steps; the probability for

each family member to be affected was calculated using a latent continuous liability [19], and then this probability is incorporated into a quasi-likelihood score test. With an extensive simulation, we showed that the proposed method performed better than the existing methods, particularly for a disease with large heritability. Application of our method to Alzheimer's disease (AD) demonstrated its practical use in the detection of genetic associations in ascertained family-based samples.

### Methods

#### Notations and statistic

We assumed that there were  $n$  families and  $n_i$  family members in each family. We considered the situation where the family of size  $n_i$  was ascertained because it contained a particular set of  $p_i$  members, and we let  $q_i = n_i - p_i$ . We called the members of the set of  $p_i$  family members "probands", and the remaining  $q_i$  individuals "non-probands". To provide a clearer motivation on this concept, we randomly selected two families, family 1 and 2, from our AD data (see Fig. 1). In family 1 (Fig. 1-(a)), individual 9 was diagnosed as AD and individuals 3–8 were selected as her relatives for genetic analysis. In family 2 (Fig. 1-(b)), individual 3 was diagnosed as AD, and individuals 4–6 were selected. Therefore  $p_1 = p_2 = 1$ ,  $q_1 = 6$  and  $q_2 = 3$  in this example. In real data analysis,  $p_i$  is often 1 and  $q_i = n_i - 1$ . We assumed that  $N$  individuals were available and thus  $N = \sum_i n_i$ . The genotypes were coded as 0, 1, or 2, according to the number of disease alleles.  $x_{ij}^P$  and  $x_{i'j'}^N$  were defined as the genotypes of proband  $j$  and non-proband  $j'$  in family  $i$  and family  $i'$ , respectively. Phenotypes were coded as 0 for an unaffected individual and 1 for an affected individual. If we let the prevalence of the disease be  $p$ , a missing phenotype was coded as  $p$ . We denoted the phenotypes of a proband and non-proband by  $y_{ij}^P$  and  $y_{i'j'}^N$ , respectively, and the vectors for genotypes and phenotypes in family  $i$  were defined by



$$\mathbf{X}_i^P = \begin{pmatrix} x_{i1}^P \\ x_{i2}^P \\ \vdots \\ x_{ip_i}^P \end{pmatrix}, \mathbf{X}_i^N = \begin{pmatrix} x_{i1}^N \\ x_{i2}^N \\ \vdots \\ x_{iq_i}^N \end{pmatrix}, \mathbf{X}_i = \begin{pmatrix} \mathbf{X}_i^P \\ \mathbf{X}_i^N \end{pmatrix},$$

$$\mathbf{Y}_i^P = \begin{pmatrix} y_{i1}^P \\ y_{i2}^P \\ \vdots \\ y_{ip_i}^P \end{pmatrix}, \mathbf{Y}_i^N = \begin{pmatrix} y_{i1}^N \\ y_{i2}^N \\ \vdots \\ y_{iq_i}^N \end{pmatrix}, \text{ and } \mathbf{Y}_i = \begin{pmatrix} \mathbf{Y}_i^P \\ \mathbf{Y}_i^N \end{pmatrix}$$

We also denoted the  $w \times w$  identity matrix by  $\mathbf{I}_w$ , and the  $w \times 1$  column vector  $\mathbf{1}_w$  indicated a vector in which all elements were 1. Let  $\pi_{ijj}^P$  and  $\pi_{ijj}^N$  be the kinship coefficient between probands  $j$  and  $j'$  in family  $i$ , and non-proband  $j$  and  $j'$  in family  $i$ , respectively. In addition, we let  $\pi_{ijj}^{PN}$  be the kinship coefficient between proband  $j$  and non-proband  $j'$  in family  $i$ , and let  $d_{ij}^P$  and  $d_{ij}^N$  be the inbreeding coefficient for proband  $j$  and non-proband  $j'$  in family  $i$ , respectively. The inbreeding coefficient is the parameter that quantifies the departure from Hardy-Weinberg equilibrium (HWE) and ranges from 0 to 1. Several approaches [20, 21] that can estimate  $d_{ij}$  have been proposed. We let

$$\mathbf{R}_i^P = \begin{pmatrix} 1 + d_{i1}^P & \cdots & 2\pi_{i1p_i}^P \\ \vdots & \ddots & \vdots \\ 2\pi_{ip_i1}^P & \cdots & 1 + d_{ip_i}^P \end{pmatrix},$$

$$\mathbf{R}_i^N = \begin{pmatrix} 1 + d_{i1}^N & \cdots & 2\pi_{i1q_i}^N \\ \vdots & \ddots & \vdots \\ 2\pi_{iq_i1}^N & \cdots & 1 + d_{iq_i}^N \end{pmatrix},$$

$$\mathbf{R}_i^{PN} = \begin{pmatrix} 2\pi_{i11}^{PN} & \cdots & 2\pi_{i1q_i}^{PN} \\ \vdots & \ddots & \vdots \\ 2\pi_{ip_i1}^{PN} & \cdots & 2\pi_{ip_iq_i}^{PN} \end{pmatrix},$$

and  $\mathbf{R}_i$  is defined by

$$\mathbf{R}_i = \begin{pmatrix} \mathbf{R}_i^P & \mathbf{R}_i^{PN} \\ (\mathbf{R}_i^{PN})^t & \mathbf{R}_i^N \end{pmatrix}$$

If we let  $q_A$  be the disease allele frequency,  $E(\mathbf{X}_i)$  was  $2q_A\mathbf{1}_{n_i}$ , and  $q_A$  is estimated with the best linear unbiased estimator (BLUE).  $\text{var}(\mathbf{X}_i)$  is expressed by  $\sigma^2\mathbf{R}_i$ , and  $\sigma^2$  is equal to  $2q_A(1 - q_A)$  under HWE.

When we analyzes the distribution of genotypes as in the FBAT approach, the statistical efficiency of the test statistic could be improved by adjustments of the phenotype with the so-called offset [22]. If we let  $\mu_{ij}^P$  and  $\mu_{ij}^N$  be offsets for proband  $j$  and non-proband  $j'$  in family  $i$

and family  $i'$ , respectively, the offset vector for family  $i$  is defined as

$$\mu_i^P = \begin{pmatrix} \mu_{i1}^P \\ \mu_{i2}^P \\ \vdots \\ \mu_{ip_i}^P \end{pmatrix}, \mu_i^N = \begin{pmatrix} \mu_{i1}^N \\ \mu_{i2}^N \\ \vdots \\ \mu_{iq_i}^N \end{pmatrix}, \text{ and } \mu_i = \begin{pmatrix} \mu_i^P \\ \mu_i^N \end{pmatrix}$$

Setting  $\mathbf{T}_i = \mathbf{Y}_i - \mu_i$ , we can define

$$\mathbf{X} = \begin{pmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \vdots \end{pmatrix}, \mathbf{Y} = \begin{pmatrix} \mathbf{Y}_1 \\ \mathbf{Y}_2 \\ \vdots \end{pmatrix}, \mathbf{T} = \begin{pmatrix} \mathbf{T}_1 \\ \mathbf{T}_2 \\ \vdots \end{pmatrix}, \text{ and}$$

$$\mathbf{R} = \begin{pmatrix} \mathbf{R}_1 & 0 & \cdots \\ 0 & \mathbf{R}_2 & \cdots \\ \vdots & \vdots & \ddots \end{pmatrix}.$$

We denoted a minor allele frequency (MAF) of a variant in unaffected individuals by  $q$ . We assumed [18] that for a constant  $\gamma$ ,

$$E(\mathbf{X}|\mathbf{T}) = 2p\mathbf{1}_N + \gamma\mathbf{T},$$

where  $0 < 2p + \gamma < 1$ . Then, the score for a variant [18, 23] can be defined by

$$S = \mathbf{T}^t(\mathbf{X} - \hat{E}(\mathbf{X})) \text{ and } \hat{E}(\mathbf{X}) = \mathbf{1}_N(\mathbf{1}_N^t\mathbf{R}^{-1}\mathbf{1}_N)^{-1}\mathbf{1}_N^t\mathbf{R}^{-1}\mathbf{X}.$$

The variance of  $S$  is

$$\text{var}(S) = \sigma^2\mathbf{T}^t\mathbf{V}^{-1}(\mathbf{R} - \mathbf{1}_N(\mathbf{1}_N^t\mathbf{R}^{-1}\mathbf{1}_N)^{-1}\mathbf{1}_N^t)\mathbf{V}^{-1}\mathbf{T},$$

and we considered the following statistic [17, 18]:

$$\frac{\mathbf{T}^t(\mathbf{1}_N - (\mathbf{1}_N^t\mathbf{R}^{-1}\mathbf{1}_N)^{-1}\mathbf{1}_N^t\mathbf{R}^{-1})\mathbf{X}}{\sqrt{\sigma^2\mathbf{T}^t(\mathbf{R} - \mathbf{1}_N(\mathbf{1}_N^t\mathbf{R}^{-1}\mathbf{1}_N)^{-1}\mathbf{1}_N^t)\mathbf{T}}} \sim N(0, 1) \text{ if } \gamma = 0.$$

This statistic will be denoted by  $WL$  in the remainder of this report.

### Adjusting the heterogeneous ascertainment bias

Families are often selected based on some probands, and the probability for family members to be affected depends on their relationship with the probands. Additional file 1 shows that the incorporation of conditional probability of each individual being affected to  $WL$  as offset lead to asymptotically smaller variance and therefore the adjustment of heterogeneous ascertainment bias is required to improve the statistical power of  $WL$ . This probability could be estimated with the liability model if the heritabilities,  $h^2$ , and prevalence,  $p$ , were available. We let  $l_{ij}^P$  and  $l_{ij}^N$  be the liability of proband  $j$  and non-proband  $j'$  in family  $i$  and family  $i'$ , respectively, and let  $\mathbf{L}_i^P = (l_{i1}^P, \dots, l_{ip_i}^P)$  and  $\mathbf{L}_i^N = (l_{i1}^N, \dots, l_{iq_i}^N)$ . We assumed that each liability

followed the standard normal distribution, and their joint distributions were

$$\begin{pmatrix} L_i^P \\ L_i^N \end{pmatrix} \sim MVN\left(0, h^2 \begin{pmatrix} \mathbf{R}_i^P & \mathbf{R}_i^{PN} \\ (\mathbf{R}_i^{PN})^t & \mathbf{R}_i^N \end{pmatrix} + (1-h^2)\mathbf{I}_{n_i}\right).$$

Benchek and Morris [24] reported that significant asymptotic biases are likely to arise when the multivariate normal (MVN) liability assumption is not met and in such a case, different assumptions should be considered. We assume that  $\mathbf{M}_i^{P*}$  and  $\mathbf{V}_i^{P*}$  are the expectation and variances of  $L_i^P$  when their disease statuses are conditioned. If all probands are affected, they becomes

$$\mathbf{M}_i^{P*} \equiv E(L_i^P | l_{i1}^P > c, \dots, l_{ip_i}^P > c)$$

and

$$\mathbf{V}_i^{P*} \equiv \text{var}(L_i^P | l_{i1}^P > c, \dots, l_{ip_i}^P > c).$$

They can be calculated with the numerical algorithms [25]. If  $p_i$  is 1, both can be simply calculated. We denote the cumulative and probability density function of standard normal distribution by  $\Phi(\cdot)$  and  $\phi(\cdot)$ . If we let  $c$  be the  $(1-p)$ th quantile of the standard normal distribution,  $\mathbf{M}_i^{P*}$  and  $\mathbf{V}_i^{P*}$  becomes

$$\mathbf{M}_i^{P*} \begin{cases} \phi(c)/[1-\Phi(x)] & \text{if } y_{i1}^P = 1 \\ -\phi(c)/\Phi(x) & \text{if } y_{i1}^P = 0 \end{cases}, \text{ and}$$

$$\mathbf{V}_i^{P*} = 1 - (\mathbf{M}_i^{P*})^2 + \mathbf{M}_i^{P*} c.$$

With Pearson-Aitken formula [26, 27], we could obtain the conditional mean and variance-covariance matrix of  $\mathbf{L}_i^N$  given  $\mathbf{L}_i^P > 1_{p_i} \cdot c$  as follows:

$$\mathbf{M}_i^{N*} = 0 + (\mathbf{R}_i^{PN})^t (\mathbf{V}_i^P)^{-1} (\mathbf{M}_i^{P*} - 0)$$

and

$$\mathbf{V}_i^{N*} = \mathbf{V}_i^N + (\mathbf{R}_i^{PN})^t ((\mathbf{V}_i^P)^{-1} - (\mathbf{V}_i^P)^{-1} \mathbf{V}_i^{P*} (\mathbf{V}_i^P)^{-1}) \mathbf{R}_i^{PN}.$$

We denoted the  $j$ th element in  $\mathbf{M}_i^{N*}$  by  $m_j^{N*}$  and the  $j$ th diagonal element in  $\mathbf{V}_i^{N*}$  by  $v_j^{N*}$ . Then the probability of being affected for a non-proband under multivariate normality of the liabilities could be calculated as

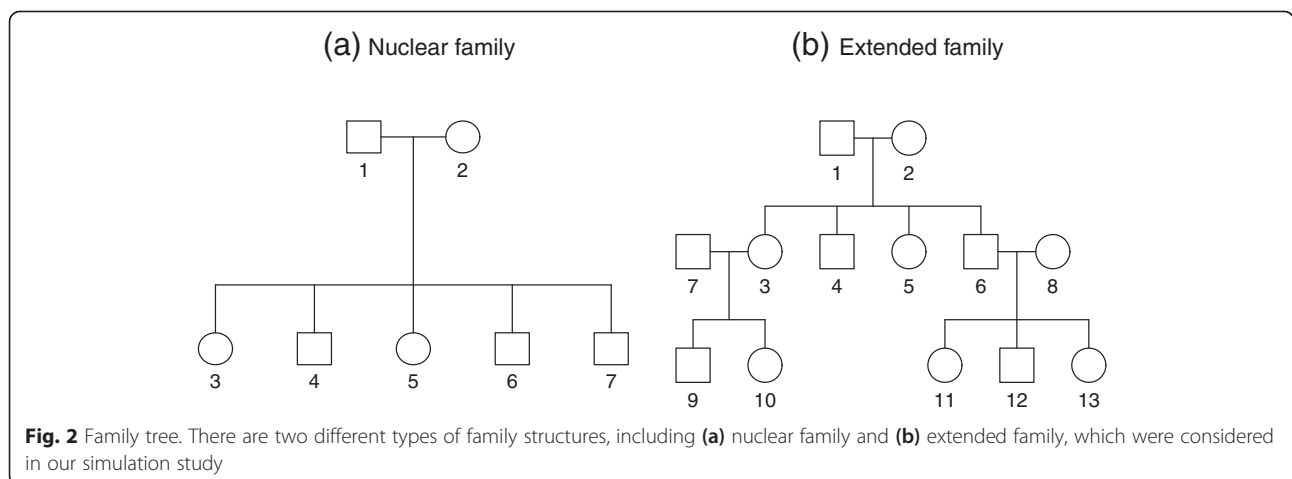
$$\Phi\left(\frac{c - m_j^{N*}}{v_j^{N*}}\right),$$

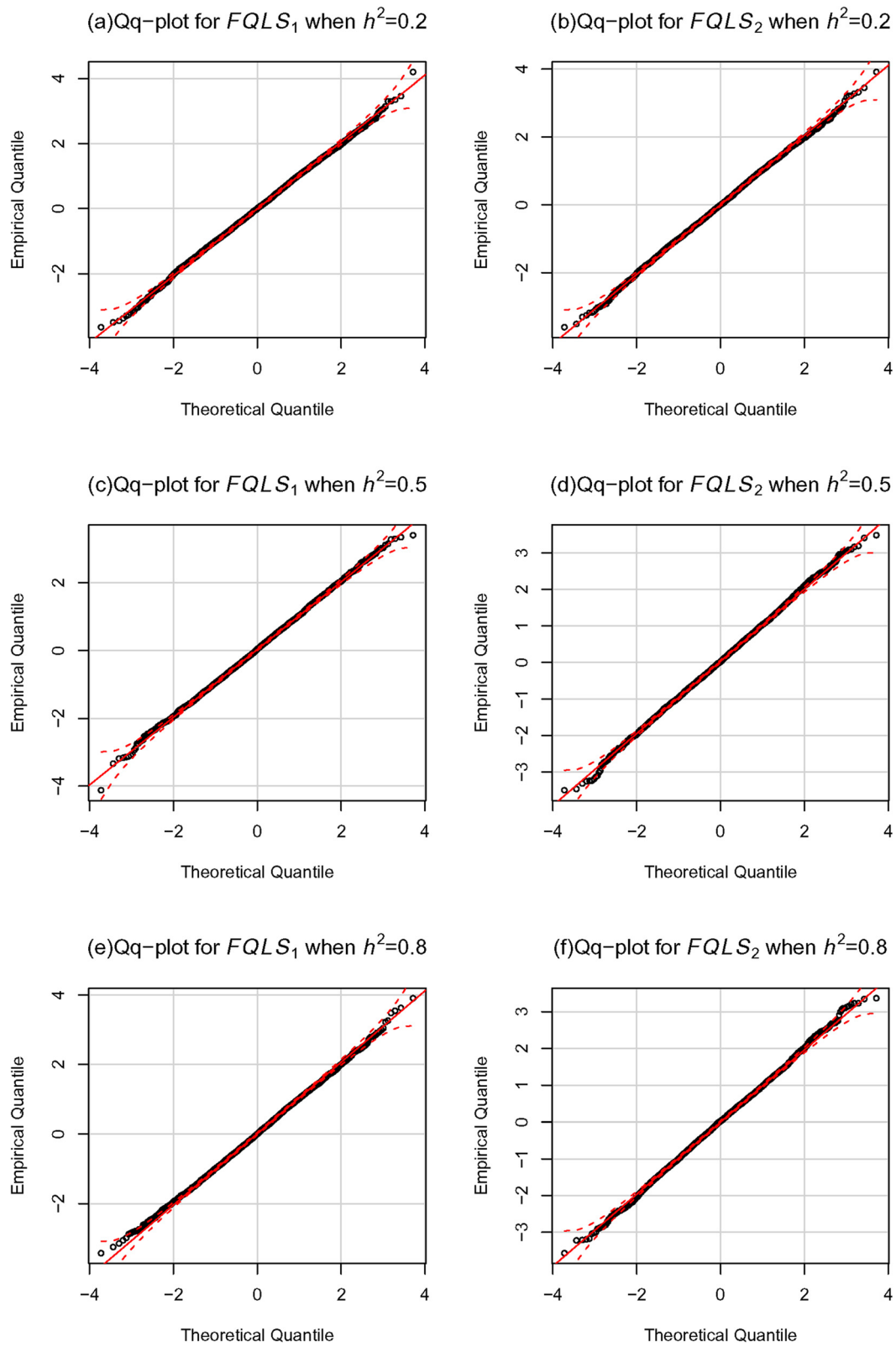
and this will be incorporated into the proposed statistic as offset. Thus far, we have assumed that there was a well-designed set of  $p_i$  individuals who were ‘‘probands’’, and for this situation, we calculated the statistic as indicated and denoted  $FQLS_1$ . However in practice, different ascertainment condition such as sequential sampling frame [28] are often utilized, and the set of  $p_i$  individuals will not be well defined. For this situation, we calculated the probability for each individual to be affected under the assumption that all the other family members were ‘‘probands’’, and thus  $p_i = n_i - 1$  and  $q_i = 1$ . The statistic calculated this way was denoted by  $FQLS_2$ .

## Results

### The simulation model

In our simulation studies, we considered two types of family structures; nuclear families with five offspring and the extended families that consist of 13 individuals along 3 generations (see Fig. 2). The latter will be called extended families in the remainder of this report. The disease allele frequency,  $p$ , was assumed to be 0.2. If we denoted the disease allele frequency by  $q_A$ , the genotype frequencies for  $AA$ ,  $Aa$ , and  $aa$  became  $q_A^2$ ,  $2q_A(1 - q_A)$ , and  $(1 - q_A)^2$  under HWE, respectively, and founders' genotypes were generated under the corresponding multinomial distribution. The genotypes for non-founders were generated with randomly generated Mendelian transmission. The disease status was generated with the liability





**Fig. 3** QQ plots for  $FQLS_1$  and  $FQLS_2$  under the null hypothesis. QQplots for  $FQLS_1$  and  $FQLS_2$  are obtained when  $h^2$  is 0.2((a), (b)), 0.5((c), (d)), or 0.8((e), (f)). P-values were calculated based on 5000 replicates when the number of families was 900. The genetic effect  $\beta$  was assumed to be 0, and the minor allele frequency was 0.2

threshold model. Once continuous liabilities that consisted of polygenic effects and random errors were generated, they were transformed to being affected if they were larger than the threshold; and otherwise, they were considered to be unaffected. The threshold was chosen to preserve the prevalence, and prevalence was assumed to be 0.2. Continuous liability was determined by combining the phenotypic mean, polygenic effect, main genetic effect, and random error. The main genetic effect for each individual was the product of  $\beta$  and the number of disease alleles. If we denoted the relative proportion of the phenotypic variance attributable to the main disease gene by  $h_a^2$ , and  $h^2$  was a heritability for continuous liability,  $\beta$  was calculated by

$$\beta = \sqrt{\frac{h_a^2}{2q_A(1-q_A)(1-h^2)}}$$

For the evaluation of type-1 errors and power,  $h_a^2$  was assumed to be 0 and 0.005, respectively. Phenotypic correlations between family-members were explained by the polygenic effects. Parental polygenic effects were generated from  $N(0, h^2)$ , and  $h^2$  was assumed to be 0.2, 0.5, or 0.8. For non-founders, the average of maternal and paternal polygenic effects was combined with the values independently sampled from  $N(0, 0.5 h^2)$  for the polygenic effects of offspring. Random errors were generated from  $N(0, \sigma_e^2 = 1-h^2)$ . For each replicate, sampling was repeated until a given number of ascertained families was generated. Type-1 error estimates were calculated with 5000 replicates, and empirical power estimates were calculated with 1000 replicates.

**Table 1** Empirical type-1 error estimates. The empirical type-1 error rates and their 95 % confidence intervals were estimated with 5000 replicates at the 0.01 and 0.05 significance level for  $h^2 = 0.2, 0.5,$  and  $0.8$ . The number of families was assumed to be 900, and the disease allele frequency was 0.2

$h^2$	Statistics	Type-1 error estimates	95 % confidence interval		
			Lower	Upper	
0.01	0.2	<i>FQLS</i> <sub>1</sub>	0.011	0.008	0.014
		<i>FQLS</i> <sub>2</sub>	0.011	0.008	0.013
	0.5	<i>FQLS</i> <sub>1</sub>	0.010	0.007	0.013
		<i>FQLS</i> <sub>2</sub>	0.009	0.007	0.012
	0.8	<i>FQLS</i> <sub>1</sub>	0.009	0.006	0.011
		<i>FQLS</i> <sub>2</sub>	0.009	0.007	0.012
0.05	0.2	<i>FQLS</i> <sub>1</sub>	0.049	0.043	0.055
		<i>FQLS</i> <sub>2</sub>	0.049	0.043	0.055
	0.5	<i>FQLS</i> <sub>1</sub>	0.050	0.044	0.056
		<i>FQLS</i> <sub>2</sub>	0.053	0.047	0.059
	0.8	<i>FQLS</i> <sub>1</sub>	0.047	0.042	0.053
		<i>FQLS</i> <sub>2</sub>	0.053	0.047	0.059

**Evaluation of the proposed methods with simulated data**

The empirical type-1 errors for *FQLS*<sub>1</sub> and *FQLS*<sub>2</sub> were evaluated from 5000 replicates under the situation of no association ( $h_a^2 = 0$ ), and 900 nuclear families with five offspring in Fig. 2 were generated for each replicate. Fig. 3 shows the quantile quantile (QQ) plots from 5000 replicates, and the nominal significance levels for both methods were preserved for various significance levels. We also estimated the empirical type-1 error rates at the 0.01 and 0.05 significance levels; the empirical type-1 error estimates of *FQLS*<sub>1</sub> and *FQLS*<sub>2</sub> preserved these nominal significance levels (Table 1). These results verified that the use of the approximation to the standard normal distribution resulted in an accurate assessment of significance for the proposed methods.

**Table 2** Empirical power estimates for scenario 1 when  $h^2$  is 0.2. The empirical power estimates for scenario 1 were calculated with 1000 replicates at the both 0.01 and 0.001 significance levels. The disease allele frequency was assumed to be 0.2, and the prevalence was assumed to be 0.2. The relative phenotypic variance attributable to the main disease gene was assumed to be 0.005

$n_{proband}$	Statistic	<i>N</i>	<i>N</i>					
			100	300	600	900	1200	1400
0.01	1	WL	0.027	0.062	0.129	0.219	0.293	0.369
		<i>FQLS</i> <sub>1</sub>	<b>0.029</b>	<b>0.064</b>	0.122	0.220	0.304	0.372
		<i>FQLS</i> <sub>2</sub>	<b>0.029</b>	0.059	<b>0.129</b>	<b>0.230</b>	<b>0.306</b>	<b>0.385</b>
	2	WL	0.033	<b>0.076</b>	0.165	0.295	0.415	0.456
		<i>FQLS</i> <sub>1</sub>	<b>0.036</b>	0.073	0.166	0.309	0.418	0.461
		<i>FQLS</i> <sub>2</sub>	0.035	<b>0.076</b>	<b>0.178</b>	<b>0.312</b>	<b>0.422</b>	<b>0.465</b>
	3	WL	<b>0.040</b>	0.103	<b>0.257</b>	0.398	0.519	0.609
		<i>FQLS</i> <sub>1</sub>	0.039	<b>0.112</b>	<b>0.257</b>	<b>0.404</b>	0.526	<b>0.623</b>
		<i>FQLS</i> <sub>2</sub>	0.036	0.110	0.253	0.407	<b>0.527</b>	0.621
	4	WL	0.041	0.127	0.297	0.497	0.626	0.720
		<i>FQLS</i> <sub>1</sub>	0.044	<b>0.129</b>	<b>0.307</b>	<b>0.502</b>	0.626	0.719
		<i>FQLS</i> <sub>2</sub>	<b>0.046</b>	0.126	<b>0.307</b>	0.493	<b>0.637</b>	<b>0.726</b>
0.001	1	WL	<b>0.005</b>	<b>0.012</b>	<b>0.038</b>	0.075	0.113	0.153
		<i>FQLS</i> <sub>1</sub>	<b>0.005</b>	<b>0.012</b>	<b>0.038</b>	<b>0.076</b>	<b>0.115</b>	0.157
		<i>FQLS</i> <sub>2</sub>	0.003	0.010	<b>0.038</b>	<b>0.076</b>	0.106	<b>0.163</b>
	2	WL	0.006	<b>0.019</b>	0.060	0.098	0.193	0.224
		<i>FQLS</i> <sub>1</sub>	<b>0.007</b>	0.018	0.059	0.099	<b>0.199</b>	0.217
		<i>FQLS</i> <sub>2</sub>	0.006	<b>0.019</b>	<b>0.064</b>	<b>0.100</b>	0.197	<b>0.219</b>
	3	WL	0.004	0.018	0.086	0.164	0.267	0.337
		<i>FQLS</i> <sub>1</sub>	0.007	0.020	<b>0.091</b>	0.162	<b>0.275</b>	0.333
		<i>FQLS</i> <sub>2</sub>	<b>0.008</b>	<b>0.023</b>	0.087	<b>0.165</b>	<b>0.275</b>	<b>0.348</b>
	4	WL	<b>0.010</b>	0.029	0.116	0.231	0.309	0.451
		<i>FQLS</i> <sub>1</sub>	0.009	0.029	0.116	0.228	<b>0.370</b>	0.449
		<i>FQLS</i> <sub>2</sub>	0.009	<b>0.031</b>	<b>0.118</b>	<b>0.233</b>	0.363	<b>0.459</b>

The bold text indicates the highest empirical estimate of the power for each situation

The empirical powers at the various significance levels were measured based on 1000 replicates at the 0.01 and 0.001 significance levels. The relative proportion,  $h_a^2$ , of phenotypic variance attributable to the main disease gene,  $2p_A(1 - p_A)\beta^2$ , was assumed to be 0.005, and nuclear and extended families in Fig. 2 were considered for the power comparison. In the first simulation setting, the numbers of nuclear families were assumed to be 100, 300, 600, 900, 1200, and 1400, and half of the families were ascertained if the number of affected family members was larger than or equal to  $n_{proband}$ , and the other half of the families were ascertained if the number of unaffected family members was larger than or equal to  $n_{proband}$ . Therefore, if 100 nuclear families were generated, half of nuclear families should have more than or equal to  $n_{proband}$  affected family

**Table 3** Empirical power estimates for scenario 1 when  $h^2$  is 0.5. The empirical power estimates for scenario 1 were calculated with 1000 replicates at the both 0.1 and 0.001 significance levels. The disease allele frequency was assumed to be 0.2, and the prevalence was assumed to be 0.2. The relative phenotypic variance attributable to the main disease gene was assumed to be 0.005

	$n_{proband}$	Statistic	N						
			100	300	600	900	1200	1400	
0.01	1	WL	0.052	0.142	0.369	0.518	0.682	0.765	
		FQLS <sub>1</sub>	0.053	0.145	0.352	0.523	0.681	0.766	
		FQLS <sub>2</sub>	<b>0.056</b>	<b>0.151</b>	<b>0.389</b>	<b>0.543</b>	<b>0.702</b>	<b>0.796</b>	
		2	WL	0.053	0.174	0.396	0.616	0.761	0.829
			FQLS <sub>1</sub>	<b>0.054</b>	0.183	0.400	0.619	0.780	0.834
			FQLS <sub>2</sub>	0.053	<b>0.202</b>	<b>0.422</b>	<b>0.658</b>	<b>0.799</b>	<b>0.859</b>
	3	WL	0.068	0.118	0.470	0.692	0.808	0.818	
		FQLS <sub>1</sub>	<b>0.073</b>	0.216	0.489	0.705	<b>0.826</b>	<b>0.905</b>	
	4	FQLS <sub>2</sub>	0.067	<b>0.232</b>	<b>0.498</b>	<b>0.729</b>	<b>0.826</b>	0.904	
		WL	0.066	0.222	0.528	0.755	0.860	0.938	
		FQLS <sub>1</sub>	0.068	<b>0.250</b>	<b>0.560</b>	<b>0.774</b>	<b>0.882</b>	<b>0.939</b>	
	0.001	1	FQLS <sub>2</sub>	<b>0.072</b>	0.244	0.544	0.766	0.870	0.935
WL			0.012	0.034	0.143	0.247	0.376	0.505	
FQLS <sub>1</sub>			<b>0.015</b>	<b>0.046</b>	0.142	0.245	0.387	0.496	
FQLS <sub>2</sub>			0.013	<b>0.046</b>	<b>0.169</b>	<b>0.305</b>	<b>0.442</b>	<b>0.567</b>	
2			WL	0.008	0.060	0.151	0.342	0.494	0.612
			FQLS <sub>1</sub>	<b>0.014</b>	<b>0.055</b>	0.163	0.357	0.512	0.640
3		FQLS <sub>2</sub>	0.012	<b>0.055</b>	<b>0.184</b>	<b>0.378</b>	<b>0.538</b>	<b>0.648</b>	
		WL	0.005	0.033	0.223	0.404	0.579	0.595	
4		FQLS <sub>1</sub>	0.008	0.076	0.235	0.432	0.604	<b>0.717</b>	
		FQLS <sub>2</sub>	<b>0.010</b>	<b>0.078</b>	<b>0.236</b>	<b>0.438</b>	<b>0.610</b>	0.699	
4		WL	0.010	0.079	0.274	0.484	0.655	0.763	
		FQLS <sub>1</sub>	0.008	<b>0.088</b>	<b>0.296</b>	<b>0.490</b>	<b>0.677</b>	<b>0.783</b>	
	FQLS <sub>2</sub>	<b>0.014</b>	0.084	0.280	0.474	0.671	0.769		

The bold text indicates the highest empirical estimate of the power for each situation

**Table 4** Empirical power estimates for scenario 1 when  $h^2$  is 0.8. The empirical power estimates for scenario 1 were calculated with 1000 replicates at the both 0.01 and 0.001 significance levels. The disease allele frequency was assumed to be 0.2, and the prevalence was assumed to be 0.2. The relative phenotypic variance attributable to the main disease gene was assumed to be 0.005

	$n_{proband}$	Statistic	N						
			100	300	600	900	1200	1400	
0.01	1	WL	0.068	0.233	0.470	0.699	0.819	0.903	
		FQLS <sub>1</sub>	0.071	0.238	0.471	0.717	0.841	0.896	
		FQLS <sub>2</sub>	<b>0.078</b>	<b>0.298</b>	<b>0.559</b>	<b>0.817</b>	<b>0.899</b>	<b>0.943</b>	
		2	WL	0.071	0.222	0.521	0.708	0.881	0.937
			FQLS <sub>1</sub>	0.080	0.304	0.568	0.788	0.907	0.942
		3	FQLS <sub>2</sub>	<b>0.085</b>	<b>0.311</b>	<b>0.605</b>	<b>0.830</b>	<b>0.930</b>	<b>0.948</b>
	WL		0.075	0.253	0.555	0.786	0.911	0.931	
	4	FQLS <sub>1</sub>	<b>0.079</b>	0.298	0.592	0.813	0.921	0.972	
		FQLS <sub>2</sub>	0.078	<b>0.311</b>	<b>0.629</b>	<b>0.823</b>	<b>0.941</b>	<b>0.982</b>	
	4	WL	0.081	0.307	0.585	0.800	0.917	0.951	
		FQLS <sub>1</sub>	0.088	<b>0.325</b>	<b>0.622</b>	0.828	0.928	0.957	
		FQLS <sub>2</sub>	<b>0.092</b>	0.318	0.614	<b>0.832</b>	<b>0.929</b>	<b>0.962</b>	
0.001	1	WL	0.016	0.074	0.229	0.444	0.602	0.710	
		FQLS <sub>1</sub>	0.017	0.072	0.221	0.436	0.643	0.693	
		FQLS <sub>2</sub>	<b>0.021</b>	<b>0.120</b>	<b>0.309</b>	<b>0.573</b>	<b>0.739</b>	<b>0.820</b>	
		2	WL	0.020	0.074	0.251	0.436	0.676	0.778
			FQLS <sub>1</sub>	0.017	0.103	0.313	0.540	0.740	0.798
		3	FQLS <sub>2</sub>	<b>0.024</b>	<b>0.116</b>	<b>0.350</b>	<b>0.589</b>	<b>0.782</b>	<b>0.845</b>
	WL		0.024	0.081	0.278	0.513	0.734	0.810	
	4	FQLS <sub>1</sub>	0.016	0.124	0.341	0.581	0.769	0.864	
		FQLS <sub>2</sub>	<b>0.017</b>	<b>0.129</b>	<b>0.365</b>	<b>0.604</b>	<b>0.802</b>	<b>0.881</b>	
	4	WL	0.015	0.109	0.310	0.542	0.756	0.830	
		FQLS <sub>1</sub>	<b>0.021</b>	0.118	0.335	0.588	0.783	0.867	
		FQLS <sub>2</sub>	<b>0.021</b>	<b>0.133</b>	<b>0.345</b>	<b>0.597</b>	<b>0.790</b>	<b>0.874</b>	

The bold text indicates the highest empirical estimate of the power for each situation

members, and the other half should have at least  $n_{proband}$  unaffected family members. We assumed that the heritabilities were 0.2, 0.5, and 0.8, and results are shown in Tables 2, 3, 4, respectively. In the second simulation setting, the numbers of extended families were assumed to be 100, 300, 600, and 900, and all families were ascertained if the number of affected family members was larger than or equal to  $n_{proband}$ . Empirical power estimates for scenario 2 were calculated when  $h^2 = 0.2, 0.5, \text{ and } 0.8$ , and the data are shown in Tables 5, 6, 7, respectively. Our results showed that either FQLS<sub>1</sub> or FQLS<sub>2</sub> was usually the most efficient statistic, and the least efficiency was provided from WL. In particular, the power gap between the proposed methods and WL was largest if  $h^2$  was 0.8, which indicates that power improvement may be proportional to the heritability.

**Table 5** Empirical power estimates for scenario 2 when  $h^2$  is 0.2. The empirical power estimates for scenario 2 were calculated with 1000 replicates at the both 0.01, and 0.001 significance levels. The disease allele frequency was assumed to be 0.2, and the prevalence was assumed to be 0.2. The relative phenotypic variance attributable to the main disease gene was assumed to be 0.005

	$n_{proband}$	Statistic	N			
			100	300	600	900
0.01	1	WL	0.072	0.149	0.304	0.409
		FQLS <sub>1</sub>	0.072	0.147	0.305	0.415
		FQLS <sub>2</sub>	<b>0.075</b>	<b>0.158</b>	<b>0.312</b>	<b>0.434</b>
	2	WL	0.042	0.137	0.300	0.448
		FQLS <sub>1</sub>	0.039	0.136	<b>0.303</b>	0.455
		FQLS <sub>2</sub>	<b>0.041</b>	<b>0.139</b>	0.295	<b>0.471</b>
	3	WL	<b>0.058</b>	0.188	0.410	0.608
		FQLS <sub>1</sub>	0.054	<b>0.191</b>	<b>0.424</b>	<b>0.620</b>
		FQLS <sub>2</sub>	0.055	0.190	0.423	0.615
0.001	1	WL	<b>0.025</b>	0.059	0.147	0.211
		FQLS <sub>1</sub>	0.023	<b>0.062</b>	<b>0.152</b>	0.197
		FQLS <sub>2</sub>	0.022	0.055	<b>0.152</b>	<b>0.212</b>
	2	WL	0.010	<b>0.038</b>	0.123	0.229
		FQLS <sub>1</sub>	<b>0.012</b>	0.036	0.123	<b>0.232</b>
		FQLS <sub>2</sub>	0.010	0.036	<b>0.127</b>	0.227
	3	WL	0.006	0.055	<b>0.182</b>	0.342
		FQLS <sub>1</sub>	0.007	0.055	<b>0.182</b>	0.355
		FQLS <sub>2</sub>	<b>0.009</b>	<b>0.057</b>	<b>0.182</b>	<b>0.356</b>

The bold text indicates the highest empirical estimate of the power for each situation

If  $h^2$  was 0.2, the proposed methods were only slightly better than WL. While all methods in our power comparison focused on the distribution of genotypes to calculate statistics, the proposed methods uniquely considered the heterogeneous effects of ascertainment bias among family members which were proportional to the magnitude of heritability; this explained the power improvement of the proposed methods. Furthermore the differences of empirical power estimates from WL and the proposed methods are larger for Tables 5, 6, 7 than Tables 2, 3, 4, which indicates that the heterogeneity of ascertainment condition may be positively related with family size and the proposed methods become more efficient for large families. Last our simulation results show that FQLS<sub>2</sub> was slightly better than FQLS<sub>1</sub>, and this may be induced by the uncertainty of probands in our simulation studies. Therefore, we concluded that the incorporation of a sampling scheme to the offset could make a substantial difference, and test statistic should be carefully selected depending on type of sampling scheme. The bold text indicates the highest empirical estimate of the power for each situation. The bold text indicates the highest empirical estimate of the power

**Table 6** Empirical power estimates for scenario 2 when  $h^2$  is 0.5. The empirical power estimates for scenario 2 were calculated with 1000 replicates at the both 0.1 and 0.001 significance levels. The disease allele frequency was assumed to be 0.2, and the prevalence was assumed to be 0.2. The relative phenotypic variance attributable to the main disease gene was assumed to be 0.005

	$n_{proband}$	Statistic	N			
			100	300	600	900
0.01	1	WL	0.130	0.293	0.567	0.787
		FQLS <sub>1</sub>	0.130	0.295	0.568	0.773
		FQLS <sub>2</sub>	<b>0.140</b>	<b>0.323</b>	<b>0.603</b>	<b>0.823</b>
	2	WL	0.093	0.332	0.645	0.864
		FQLS <sub>1</sub>	0.094	<b>0.357</b>	0.654	0.871
		FQLS <sub>2</sub>	<b>0.108</b>	0.354	<b>0.694</b>	<b>0.902</b>
	3	WL	0.100	0.382	0.735	0.904
		FQLS <sub>1</sub>	0.108	0.406	0.751	0.915
		FQLS <sub>2</sub>	<b>0.130</b>	<b>0.408</b>	<b>0.772</b>	<b>0.932</b>
0.001	1	WL	0.046	0.148	0.341	0.560
		FQLS <sub>1</sub>	<b>0.050</b>	0.149	0.353	0.559
		FQLS <sub>2</sub>	0.047	<b>0.166</b>	<b>0.386</b>	<b>0.617</b>
	2	WL	0.019	0.127	0.387	0.634
		FQLS <sub>1</sub>	<b>0.021</b>	0.119	0.394	0.648
		FQLS <sub>2</sub>	0.017	<b>0.144</b>	<b>0.432</b>	<b>0.695</b>
	3	WL	0.023	0.166	0.481	0.749
		FQLS <sub>1</sub>	0.026	0.183	0.511	0.772
		FQLS <sub>2</sub>	<b>0.028</b>	<b>0.196</b>	<b>0.532</b>	<b>0.782</b>

The bold text indicates the highest empirical estimate of the power for each situation

for each situation. The bold text indicates the highest empirical estimate of the power for each situation. The bold text indicates the highest empirical estimate of the power for each situation. The bold text indicates the highest empirical estimate of the power for each situation. The bold text indicates the highest empirical estimate of the power for each situation.

**Robustness of the proposed methods against the misspecification of prevalence**

The statistical powers of the proposed methods may depend on the accuracy of the prevalence and we evaluated the sensitivity of the proposed method to the misspecified prevalence with simulated data.  $h_a^2$  and  $h^2$  were assumed to be 0.05 and 0.8, and nuclear families (Fig. 2-(a)) were considered in this simulation. The number of nuclear families was assumed to be 900 and  $n_{proband}$  was assumed to be 1, 2, or 3. Prevalence was assumed to be 0.2 for phenotype generation, and the offset for which we recommended prevalence was set to be 0.1, 0.2 or 0.3 for calculation of FQLS<sub>1</sub> and FQLS<sub>2</sub>. In particular,  $y_{ij}$  for individuals with missing phenotypes are coded by the



**Table 7** Empirical power estimates for scenario 2 when  $h^2$  is 0.8. The empirical power estimates for scenario 2 were calculated with 1000 replicates at the both 0.01 and 0.001 significance levels. The disease allele frequency was assumed to be 0.2, and the prevalence was assumed to be 0.2. The relative phenotypic variance attributable to the main disease gene was assumed to be 0.005

	$n_{proband}$	Statistic	N			
			100	300	600	900
0.01	1	WL	0.164	0.445	0.749	0.906
		FQLS <sub>1</sub>	0.156	0.441	0.751	0.905
		FQLS <sub>2</sub>	<b>0.194</b>	<b>0.508</b>	<b>0.817</b>	<b>0.944</b>
	2	WL	0.132	0.473	0.823	0.970
		FQLS <sub>1</sub>	0.131	0.505	0.861	0.969
		FQLS <sub>2</sub>	<b>0.150</b>	<b>0.564</b>	<b>0.905</b>	<b>0.978</b>
	3	WL	0.140	0.475	0.835	0.958
		FQLS <sub>1</sub>	0.134	0.520	0.867	0.970
		FQLS <sub>2</sub>	<b>0.167</b>	<b>0.556</b>	<b>0.886</b>	<b>0.981</b>
0.001	1	WL	0.059	0.230	0.519	0.759
		FQLS <sub>1</sub>	0.053	0.236	0.519	0.757
		FQLS <sub>2</sub>	<b>0.070</b>	<b>0.311</b>	<b>0.632</b>	<b>0.841</b>
	2	WL	0.039	0.239	0.561	0.858
		FQLS <sub>1</sub>	0.033	0.250	0.594	0.884
		FQLS <sub>2</sub>	<b>0.053</b>	<b>0.300</b>	<b>0.702</b>	<b>0.924</b>
	3	WL	0.033	0.215	0.629	0.865
		FQLS <sub>1</sub>	<b>0.046</b>	0.247	0.671	0.900
		FQLS <sub>2</sub>	0.044	<b>0.287</b>	<b>0.713</b>	<b>0.925</b>

The bold text indicates the highest empirical estimate of the power for each situation

assumed prevalence, and sensitivity of the proposed methods can be substantial when there are individuals with missing phenotypes. Therefore, individuals were randomly selected from non-probands, and their phenotypes were assumed to be unknown for calculation of the proposed statistics. The number of family members with missing phenotypes in each family was denoted by  $n_{missing}$ . The empirical powers were calculated at the 0.01 significance level with 1000 replicates. Table 8 shows that the results obtained by setting prevalence to be 0.1 and 0.3 are similar to the results when the prevalence was set to be 0.2, which indicates that the power loss attributable to the misspecified prevalence is not substantial. Furthermore the empirical power estimates are positively related with  $n_{proband}$  and inversely related with  $n_{missing}$ . If  $n_{missing}$  is larger than 3, the power loss may be more substantial. The bold text indicates the reference for the proposed statistics to compare with the misspecified prevalence.

**Application of the proposed method to AD**

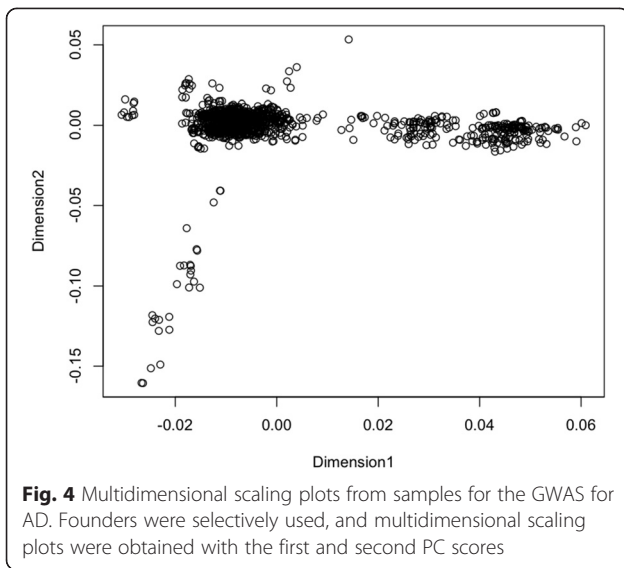
AD is an irreversible, progressive brain disorder characterized by genetic heterogeneity. However, the genetic

**Table 8** Empirical power estimates for three situations when  $h^2$  is 0.8. The empirical power estimates for three situations were calculated with 1000 replicates at the 0.01 significance levels. Phenotypes were generated under the assumption that the prevalence was assumed to be 0.2. Prevalence was set to be 0.1, 0.2 or 0.3 to calculate the proposed statistics. The relative phenotypic variance attributable to the main disease gene was assumed to be 0.005

$n_{missing}$	$n_{proband}$	Statistic	Prevalence to be set for statistics		
			0.1	0.2	0.3
0	1	FQLS <sub>1</sub>	0.486	<b>0.479</b>	<b>0.464</b>
		FQLS <sub>2</sub>	<b>0.572</b>	<b>0.539</b>	0.525
	2	FQLS <sub>1</sub>	0.562	<b>0.557</b>	0.532
		FQLS <sub>2</sub>	<b>0.606</b>	<b>0.621</b>	0.624
	3	FQLS <sub>1</sub>	0.632	<b>0.642</b>	0.634
		FQLS <sub>2</sub>	0.644	<b>0.647</b>	0.654
1	1	FQLS <sub>1</sub>	0.432	<b>0.434</b>	<b>0.418</b>
		FQLS <sub>2</sub>	0.468	<b>0.469</b>	0.462
	2	FQLS <sub>1</sub>	0.522	<b>0.530</b>	0.519
		FQLS <sub>2</sub>	0.547	<b>0.554</b>	0.559
	3	FQLS <sub>1</sub>	0.580	<b>0.584</b>	0.576
		FQLS <sub>2</sub>	0.575	<b>0.584</b>	0.585
2	1	FQLS <sub>1</sub>	0.415	<b>0.412</b>	0.399
		FQLS <sub>2</sub>	0.428	<b>0.420</b>	0.412
	2	FQLS <sub>1</sub>	0.482	<b>0.487</b>	<b>0.468</b>
		FQLS <sub>2</sub>	0.483	<b>0.495</b>	0.492
	3	FQLS <sub>1</sub>	0.487	<b>0.497</b>	0.490
		FQLS <sub>2</sub>	<b>0.500</b>	<b>0.515</b>	0.511

The bold text indicates the reference for the proposed statistics to compare with the misspecified prevalence

variations that contribute to AD still remain elusive. Thus, we applied the proposed method for identification of the disease susceptibility loci for AD. The heritability and prevalence of AD are approximately 0.8 [29, 30] and 0.1, respectively; therefore, we chose heritabilities of 0.8, and a prevalence of 0.1 for the calculation of proposed methods. Samples were collected as part of the National Institute of Mental Health Genetics Initiative (NIMH). The NIMH Alzheimer’s Disease Genetics Family Sample was used along with the information about the genotype platform (Affy 6.0) [20, 31], and ethical approach and participant approval were obtained through the NIMH IRB panel. Families were selected based on the disease status of a certain family member. However the proband for each family was not clear, and FQLS<sub>2</sub> was uniquely applied for the proposed method. 1376 individuals from 410 families were available, and all families were nuclear. All individuals were of self-reported European ancestry. HWE for each single nucleotide polymorphism (SNP) was tested, and MAFs were estimated. SNPs for which  $p$ -values for HWE were less than  $10^{-6}$  or MAFs less than



0.05 were excluded, and therefore, 417,680 SNPs were analyzed for genetic association analysis.

The choice of kinship coefficient matrix for the proposed method depends on the presence of population substructure. To confirm the presence of population substructure, multidimensional-scaling [32] analysis was performed with PLINK1.07 [20], and we constructed a multidimensional-scaling plot (Fig. 4) to provide evidence concerning the presence of population substructure. Therefore, the genomic control method [31, 33, 34] was used for explicit detection and correction of population stratification with common SNPs, and they were incorporated into both *FQLS<sub>2</sub>* and *WL*. Fig. 5 shows

**Table 9** Top significant results of GWAS for AD. For the genome-wide significant SNPs from *FQLS<sub>2</sub>* and *WL*, their *p*-values are given

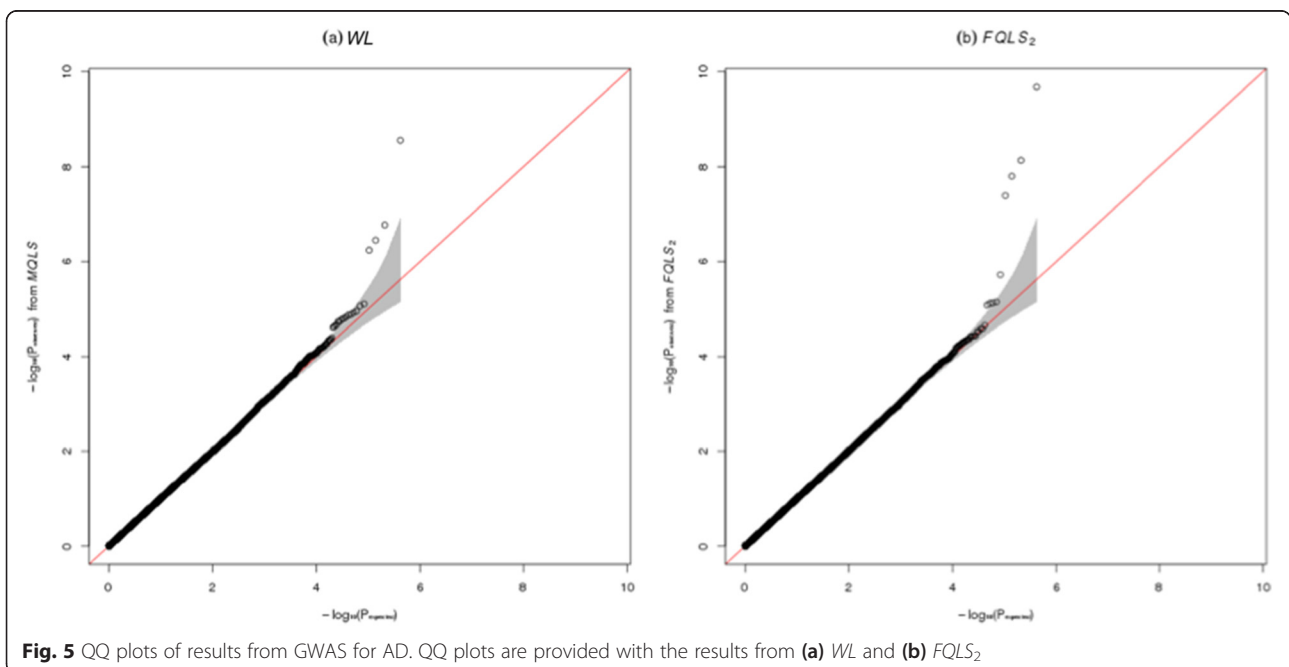
SNP	<i>WL</i>	<i>FQLS<sub>2</sub></i>	<i>FBAT</i>
SNP1	$3.53 \times 10^{-7}$	$1.94 \times 10^{-8}$	$4.23 \times 10^{-4}$
SNP2	$5.74 \times 10^{-7}$	$4.45 \times 10^{-8}$	$4.9 \times 10^{-5}$
SNP3	$1.69 \times 10^{-7}$	$8.36 \times 10^{-9}$	$8.6 \times 10^{-5}$
SNP4	$2.79 \times 10^{-9}$	$2.86 \times 10^{-10}$	$6.94 \times 10^{-12}$

QQ plots for *FQLS<sub>2</sub>* and *WL*; these plots revealed that the presence of population substructure was appropriately adjusted. Results showed four genome-wide significant results for *FQLS<sub>2</sub>* and one significant result for *WL*.

Detailed information for these significant results is provided in Table 9. We also considered FBAT statistics [35]. *FQLS<sub>2</sub>* identified four genome-wide significant SNPs, while *WL* and FBAT identified one genome-wide significant SNP. In addition, one SNP that was significant according to *WL* was more significant according to *FQLS<sub>2</sub>*. The most significant result acquired using *FQLS<sub>2</sub>* was SNP4 ( $p = 2.10 \times 10^{-10}$ ). The other three SNPs, i.e., SNP1, SNP2 and SNP3, reached the genome-wide significance level.

**Discussion**

Although major advances in high-density genome scans have enabled the genetic association analysis of more than 10,000 individuals, disappointing results in the mapping of many common diseases have illustrated the need for more powerful methods for detecting disease susceptibility loci. Statistical efficiency is known to being affected by the ascertainment bias, and its careful adjustment can



lead to substantial improvement of statistical power [36–38]. In particular, genetic association analysis has often been conducted using family-based designs, but without addressing the fact that the probability for each family member to be affected is inversely related with the familial relationship with affected probands. In this report, we proposed new methods to adjust this heterogeneity with known prevalence and heritability. Our simulation studies showed that the proposed methods provided substantial power improvement. In particular, the mis-specified heritability and prevalence can lead to the statistical power loss for the proposed methods, but it was found to be not substantial at least in our simulation studies.  $FQLS_1$  and  $FQLS_2$  were suggested, and  $FQLS_1$  is an efficient choice if probands for each family are clearly defined and all remaining family members are incorporated to the genetic analysis. However, these conditions are often not satisfied, and different methods such as sequential sampling frame [28] are usually utilized. Simulation studies showed that  $FQLS_2$  is usually better than  $FQLS_1$  if the ascertaining condition is not clearly defined and thus we recommend  $FQLS_2$  unless probands are clearly defined. However we considered the limited ascertainment conditions and comprehensive simulation studies are still necessary.

Furthermore, the proposed method was conceptually simple and can be applied to the large families. Our methods require only a single calculation of offset for all markers, and the real data analysis could be completed with a single CPU in a few hours. For  $M$  markers and  $N$  individuals, the time complexity is  $O(N^3 + MN^2)$  for the proposed method. The proposed method was implemented with C++, and can be downloaded from <http://healthstat.snu.ac.kr/mfqls/>.

Heterogeneity between samples is an important issue in large-scale genetic analysis, and the proposed method can likely be applied to various additional scenarios with some modifications. For instance, the disease status of relatives reveals the importance of genetic components for each individual, and for this reason, such information has been used, albeit only on occasion, in genetic association analysis. The effect of relatives' disease statuses is dependent on prevalence and heritability, and the probability for each individual to be affected could be calculated with the proposed method. This probability can be used to improve the statistical efficiency of genetic association analysis. In addition, the heterogeneity of the ascertainment bias is often an important issue for genome-wide meta-analysis because samples are collected from multiple medical centers [39, 40], and different sampling schemes among studies need to be adjusted to improve statistical efficiency. Therefore, we believe that proposed method can be extended to provide a statistical framework that adjusts the heterogeneity between samples.

## Conclusions

We proposed  $FQLS$  method to adjust this heterogeneity with known prevalence and heritability and the software was implemented with C++. We identified several significant associations between AD and SNPs, and their potential functional information will provide the better understanding of the pathogenesis of AD. Although this study has some limitations, our proposed methods illustrated important features required for genetic analysis with family-based samples, and an extension of the proposed method to rare variant association analysis such as *FARVAT* [41] will be investigated in future studies.

## Additional file

**Additional file 1: Web-based Supporting Materials for “Adjusting heterogeneous ascertainment bias for genetic association analysis with extended families” by Suyeon Park, Sungyoung Lee, Young Lee, Christine Herold, Basavaraj Hooli, Kristina Mullin, Lars Bertram, Taesung Park, Changsoon Park, Christoph Lange, Rudolph Tanzi, and Sungho Won.**

## Abbreviations

GWAS: Genome-wide association studies; CA: Cochran-Armitage; AD: Alzheimer's disease; PD: Parkinson's disease; HWE: Hardy-Weinberg equilibrium; BLUE: Best linear unbiased estimator; QQ: Quantile quantile; MVN: Multivariate normal; MAF: Minor allele frequency; SNP: Single nucleotide polymorphism; FBAT: Family Based Association Testing software; FQLS: Family based quasi-likelihood score test.

## Competing interests

The authors declare that they have no competing interests.

## Authors' contributions

SP participated in the design of the study, and performed the simulation studies and statistical analysis. SL and TP developed the program. YL conducted the statistical analysis. CH, BH, KM, CP, LB, CL and RT conceived of the study. SW conceived of the study, and participated in data analysis. All authors read and approved the final manuscript.

## Acknowledgements

SP, SL, YL, and SW were supported by the Industrial Core Technology Development Program (10040176, Development of Various Bioinformatics Software Using Next Generation Bio-Data) funded by the Ministry of Trade, Industry and Energy (MOTIE, Korea), and by Basic Science Research Program through the National Research Foundation of Korea (NRF) funded by the Ministry of Education (2013R1A1A2010437). CP was supported by the National Research Foundation of Korea Grant funded by the Korean Government (NRF-2011-220-C00004). BH, KM, CL and RT were supported by NIMH RO1 grant (2R01MH060009) and the Cure Alzheimer's Fund.

## Author details

<sup>1</sup>Department of Applied Statistics, Chung-Ang University, Seoul, Korea. <sup>2</sup>Center for Genome Science, National Institute of Health, Osong Health Technology Administration complex, Chungcheongbuk-do, Seoul, Korea. <sup>3</sup>Department of Biostatistics, Soonchunhyang University, College of Medicine, Seoul, Korea. <sup>4</sup>Interdisciplinary Program in Bioinformatics, Seoul National University, Seoul, Korea. <sup>5</sup>German Center for Neurodegenerative Diseases (DZNE), Sigmund-Freud-Str. 25, Bonn 53127, Germany. <sup>6</sup>Department of Biostatistics, Harvard School of Public Health, Boston, MA, USA. <sup>7</sup>Genetics and Aging Research Unit, MassGeneral Institute for Neurodegenerative Diseases, Department of Neurology, Massachusetts General Hospital, Harvard Medical School, Massachusetts, USA. <sup>8</sup>Department of Statistics, Seoul National University, Seoul, Korea. <sup>9</sup>Department of Vertebrate Genomics, Neuropsychiatric Genetics Group, Max Planck Institute for Molecular Genetics, Berlin, Germany. <sup>10</sup>Department of Medicine, School of Public Health, Imperial

College London, London, UK. <sup>11</sup>Harvard Medical School, Boston, MA, USA.

<sup>12</sup>Institute for Genomic Mathematics, University of Bonn, Bonn, Germany.

<sup>13</sup>German Center for Neurodegenerative Diseases, Bonn, Germany.

<sup>14</sup>Department of Public Health Science, Seoul National University, Seoul,

Korea. <sup>15</sup>Institute of Health and Environment, Seoul National University,

Seoul, Korea. <sup>16</sup>National Cancer Center, Seoul, Korea.

Received: 6 December 2014 Accepted: 6 July 2015

Published online: 19 August 2015

## References

- Maher B. Personal genomes: the case of the missing heritability. *Nature*. 2008;456(7218):18–21.
- Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorf LA, Hunter DJ, et al. Finding the missing heritability of complex diseases. *Nature*. 2009;461(7265):747–53.
- Visscher PM. Sizing up human height variation. *Nat Genet*. 2008;40(5):489–90.
- Clement K, Vaisse C, Manning BS, Basdevant A, Guygrand B, Ruiz J, et al. Genetic-variation in the beta(3)-adrenergic receptor and an increased capacity to gain weight in patients with morbid-obesity. *New Engl J Med*. 1995;333(6):352–4.
- Gu C, Todorov AA, Rao DC. Genome screening using extremely discordant and extremely concordant pairs. *Genet Epidemiol*. 1997;14(6):791–6.
- Hu SW, Zhong YF, Hao YT, Luo MQ, Zhou Y, Guo H, et al. Novel rare alleles of ABCA1 are exclusively associated with extreme high-density lipoprotein-cholesterol levels among the Han Chinese. *Clin Chem Lab Med*. 2009;47(10):1239–45.
- Khor CC, Goh DLM. Strategies for identifying the genetic basis of dyslipidemia: genome-wide association studies vs. the resequencing of extremes. *Curr Opin Lipidol*. 2010;21(2):123–7.
- Li BS, Leal SM. Methods for detecting associations with rare variants for common diseases: application to analysis of sequence data. *Am J Hum Genet*. 2008;83(3):311–21.
- Liang KY, Huang CY, Beatty TH. A unified sampling approach for multipoint analysis of qualitative and quantitative traits in sib pairs. *Am J Hum Genet*. 2000;66(5):1631–41.
- Price RA, Li WD, Zhao H. FTO gene SNPs associated with extreme obesity in cases, controls and extremely discordant sister pairs. *BMC Med Genet*. 2008;9:4.
- Risch N, Zhang H. Extreme discordant sib pairs - the design of choice for mapping quantitative trait loci in humans. *Am J Hum Genet*. 1995;57(4):1159–9.
- Risch NJ, Zhang HP. Mapping quantitative trait loci with extreme discordant sib pairs: Sampling considerations. *Am J Hum Genet*. 1996;58(4):836–43.
- Kryukov GV, Shpunt A, Stamatoyannopoulos JA, Sunyaev SR. Power of deep, all-exon resequencing for discovery of human trait genes. *Proc Natl Acad Sci U S A*. 2009;106(10):3871–6.
- Lander ES, Botstein D. Mapping Mendelian factors underlying quantitative traits using Rflp linkage maps. *Genetics*. 1989;121(1):185–99.
- Van Gestel S, Houwing-Duistermaat JJ, Adolfsson R, van Duijn CM, Van Broeckhoven C. Power of selective genotyping in genetic association analyses of quantitative traits. *Behav Genet*. 2000;30(2):141–6.
- Sasieni PD. From genotypes to genes: doubling the sample size. *Biometrics*. 1997;53(4):1253–61.
- Won S, Lange C. A general framework for robust and efficient association analysis in family-based designs: quantitative and dichotomous phenotypes. *Stat Med*. 2013;32(25):4482–98.
- Thornton T, McPeck MS. Case-control association testing with related individuals: a more powerful quasi-likelihood score test. *Am J Hum Genet*. 2007;81(2):321–37.
- Falconer DS. The inheritance of liability to diseases with variable age of onset, with particular reference to diabetes mellitus. *Ann Hum Genet*. 1967;31(1):1–20.
- Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet*. 2007;81(3):559–75.
- Yang JA, Lee SH, Goddard ME, Visscher PM. GCTA: a tool for genome-wide complex trait analysis. *Am J Hum Genet*. 2011;88(1):76–82.
- Lange C, DeMeo DL, Laird NM. Power and design considerations for a general class of family-based association tests: Quantitative traits. *Am J Hum Genet*. 2002;71(6):1330–41.
- McPeck MS, Wu X, Ober C. Best linear unbiased allele-frequency estimation in complex pedigrees. *Biometrics*. 2004;60(2):359–67.
- Benckeh PH, Morris NJ. How meaningful are heritability estimates of liability? *Hum Genet*. 2013;132(12):1351–60.
- Leppard P, Tallis GM. Evaluation of the mean and covariance of the truncated multinormal distribution. *J R Stat Soc C*. 1989;38(3):543–53.
- Aitken AC. Note on selection from a multivariate normal population. *Proc Edinburgh Mathematical Society B*. 1934;4:106–10.
- Pearson K. Mathematical contributions to the theory of evolution. XI. On the influence of natural selection on the variability and correlation of organs. *Philos Transact A Math Phys Eng Sci*. 1903;200:1–66.
- Elston RC, Sobel E. Sampling considerations in the gathering and analysis of pedigree data. *Am J Hum Genet*. 1979;31(1):62–9.
- Gatz M, Reynolds CA, Fratiglioni L, Johansson B, Mortimer JA, Berg S, et al. Role of genes and environments for explaining Alzheimer disease. *Arch Gen Psychiatry*. 2006;63(2):168–74.
- Gatz M, Pedersen NL, Berg S, Johansson B, Johansson K, Mortimer JA, et al. Heritability for Alzheimer's disease: the study of dementia in Swedish twins. *J Gerontol A Biol Sci Med Sci*. 1997;52(2):M117–25.
- Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet*. 2006;38(8):904–9.
- Li Q, Yu K. Improved correction for population stratification in genome-wide association studies by identifying hidden population structures. *Genet Epidemiol*. 2008;32(3):215–26.
- Devlin B, Roeder K. Genomic control for association studies. *Biometrics*. 1999;55(4):997–1004.
- Reich DE, Goldstein DB. Detecting association in a case-control study while correcting for population stratification. *Genet Epidemiol*. 2001;20(1):4–16.
- Laird NM, Horvath S, Xu X. Implementing a unified approach to family-based tests of association. *Genet Epidemiol*. 2000;19 Suppl 1:S36–42.
- Vieland VJ, Hodge SE. Inherent intractability of the ascertainment problem for pedigree data: a general likelihood framework. *Am J Hum Genet*. 1995;56(1):33–43.
- Nielsen R, Hubisz MJ, Clark AG. Reconstituting the frequency spectrum of ascertained single-nucleotide polymorphism data. *Genetics*. 2004;168(4):2373–82.
- Clark AG, Hubisz MJ, Bustamante CD, Williamson SH, Nielsen R. Ascertainment bias in studies of human genome-wide polymorphism. *Genome Res*. 2005;15(11):1496–502.
- Cantor RM, Lange K, Sinsheimer JS. Prioritizing GWAS results: A review of statistical methods and recommendations for their application. *Am J Hum Genet*. 2010;86(1):6–22.
- Ripke S, Wray NR, Lewis CM, Hamilton SP, Weissman MM, Breen G, et al. A mega-analysis of genome-wide association studies for major depressive disorder. *Mol Psychiatry*. 2013;18(4):497–511.
- Choi S, Lee S, Cichon S, Nothen MM, Lange C, Park T, Won S: FARVAT: a family-based rare variant association test. *Bioinformatics*. 2014;30(22):3197–205.

**Submit your next manuscript to BioMed Central and take full advantage of:**

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at  
www.biomedcentral.com/submit

