



Published in final edited form as:

Genet Epidemiol. 2015 September ; 39(6): 399–405. doi:10.1002/gepi.21913.

Sequence kernel association analysis of rare variant set based on the marginal regression model for binary traits

Baolin Wu^{1,*}, James S. Pankow², and Weihua Guan^{1,*}

¹Division of Biostatistics, School of Public Health, University of Minnesota

²Division of Epidemiology and Community Health, School of Public Health, University of Minnesota

Abstract

Recent sequencing efforts have focused on exploring the influence of rare variants on the complex diseases. Gene-level based tests by aggregating information across rare variants within a gene have become attractive to enrich the rare variant association signal. Among them, the sequence kernel association test has proved to be a very powerful method for jointly testing multiple rare variants within a gene. In this article, we explore an alternative sequence kernel association test. We propose to use the univariate likelihood ratio statistics from the marginal model for individual variants as input into the kernel association test. We show how to compute its significance p-value efficiently based on the asymptotic chi-square mixture distribution. We demonstrate through extensive numerical studies that the proposed method has competitive performance. Its usefulness is further illustrated with application to associations between rare exonic variants and type 2 diabetes in the Atherosclerosis Risk in Communities (ARIC) Study. We identified an exome-wide significant rare variant set in the gene *ZZZ3* worthy of further investigations.

Keywords

GWAS; SKAT; Score statistic; Sequencing data

Introduction

In GWAS, observed effect sizes for common variants have typically been quite small. In combination they explain a small proportion of the phenotypic variance. Manolio *et al.* (2009) have suggested that rare variants could have substantial effect sizes without demonstrating clear Mendelian segregation, and could contribute substantially to missing heritability. Individual rare variant based tests typically lack power due to low minor allele frequencies, and gene-level based association tests implemented by aggregating information across rare variants within a gene have become attractive to enrich the association signal. An intuitive and simple approach to aggregating signals across rare variants collapses the rare variants into a burden score to be linked to the phenotype (Morgenthaler and Thilly, 2007;

Correspondence to: Baolin Wu, Telephone: (612) 624-0647, Fax: (612) 626-0660, baolin@umn.edu, Address: A460 Mayo Building, MMC 303, Division of Biostatistics, School of Public Health, University of Minnesota, Minneapolis, Minnesota 55455-0392, USA.
*Co-correspondence authors

Madsen and Browning, 2009; Morris and Zeggini, 2010; Price *et al.*, 2010; Lin and Tang, 2011). The combined multivariate and collapsing (CMC) method is an extension of the burden test by collapsing rare variants in a region within subgroups defined according to their minor allele frequencies (MAFs) (Li and Leal, 2008). The variable threshold (VT) method is a data adaptive burden test by choosing an optimal MAF threshold (Price *et al.*, 2010; Lin and Tang, 2011). The burden test works well for variants with similar effects and could lose substantial power with both protective and deleterious variants, or in the presence of many non-causal variants. The sequence kernel association test (SKAT) is based on the variance component score test and works well under various combinations of protective and deleterious variants (Wu *et al.*, 2010; Neale *et al.*, 2011; Wu *et al.*, 2011). A more flexible approach is SKAT-O, which adaptively combines the burden and the SKAT statistics (Lee *et al.*, 2012). The SKAT based approach performs well and is widely used in rare variant based association test.

Rare variants have been postulated to have large effect sizes (Manolio *et al.*, 2009). It is likely that typical GWAS only have sufficient power to detect variants with large effects. This is indeed the case for most rare disease-causing variants identified to date (Bonneton *et al.*, 2012; Zhan *et al.*, 2013; Steinthorsdottir *et al.*, 2014; Wang *et al.*, 2014; Estrada *et al.*, 2014). The SKAT is based on the score test, thus is computationally very efficient. The score test performs well when parameter is close to the null value, but could have suboptimal performance with large deviation from the null (e.g., when testing those rare variants with large effect sizes).

Recently Chen *et al.* (2014) developed a Cox SKAT for survival outcomes and adopted the likelihood ratio test for its better performance compared to the score test in the Cox proportional hazard model. In this article, we explore an alternative sequence kernel association test for binary trait in the same spirit as Chen *et al.* (2014). We use the univariate likelihood ratio statistics from the marginal model for individual variants as input into the sequence kernel association test and its adaptive test. Their significance p-values can be computed efficiently based on the asymptotic chi-square mixture distribution. We demonstrate through extensive numerical studies that the proposed method has competitive performance. We illustrate the usefulness of the proposed method through an application to associations between rare variants and type 2 diabetes in the ARIC Study.

Materials and Methods

Consider a GWAS with genotype scores G , coded as (0,1,2) for the copies of minor allele, disease status indicator Y , and additional covariates X , which could include ancestry covariate (e.g., ancestry indicator or principal components).

Consider n subjects sequenced in a region with m genotyped rare variants. For the i -th subject, let y_i denote the case-control status, $\mathbf{G}_i = (g_{i1}, \dots, g_{im})$ the genotypes for the m variants, $\mathbf{X}_i = (x_{i1}, \dots, x_{ip})$ the covariates to be adjusted. We study the disease association of rare variants based on the following logistic regression model

$$\Pr(y_i=1|\mathbf{X}_i, \mathbf{G}_i)=\text{expit}(\beta_0 + \mathbf{X}_i\boldsymbol{\alpha} + \mathbf{G}_i\boldsymbol{\beta}), \quad (1)$$

where $\boldsymbol{\alpha}$ and $\boldsymbol{\beta} = (\beta_1, \dots, \beta_m)'$ are the vector of regression coefficients for the covariates and rare variants. Here $\text{expit}(x) = 1/(1 + \exp(-x))$ is the inverse-logit function. The disease association of the m rare variants can be tested by evaluating the null hypothesis $H_0 : \boldsymbol{\beta} = 0$.

Sequence kernel association test

The sequence kernel association test (SKAT; Wu *et al.*, 2011) is derived as a variance-component score statistic by assuming that each β_j follows an arbitrary zero-mean distribution with variance $w_j^2 \psi$, where weight w_j is fixed and typically computed based on MAF, e.g., the Wu weights $w_j = \text{Beta}(f_j; 1, 25)$ (Wu *et al.*, 2011). Here f_j is the MAF of G_j and Beta is the beta distribution density function. Under this assumption, the null hypothesis $H_0 : \boldsymbol{\beta} = 0$ is equivalent to $H_0 : \psi = 0$.

Let $\mathbf{y} = (y_1, \dots, y_n)'$ denote the response vector, \mathbf{X} the $n \times p$ covariates matrix, $\mathbf{G} = (\mathbf{G}'_1, \dots, \mathbf{G}'_n)'$ the $n \times m$ genotype matrix, $\mathbf{W} = \text{diag}(w_1, \dots, w_m)$ the diagonal matrix of weights. The SKAT statistic can be computed as

$$Q = (\mathbf{y} - \hat{\boldsymbol{\pi}}_0)' \mathbf{G} \mathbf{W} \mathbf{W} \mathbf{G}' (\mathbf{y} - \hat{\boldsymbol{\pi}}_0),$$

where $\hat{\boldsymbol{\pi}}_0 = (\hat{\pi}_1, \dots, \hat{\pi}_n)'$ with $\hat{\pi}_i = \hat{\text{Pr}}(y_i = 1 | \mathbf{X}_i, \mathbf{G}_i)$ derived under the null model ($\boldsymbol{\beta} = 0$). Let $\mathbf{V}_0 = \text{diag}\{\hat{\pi}_0(1 - \hat{\pi}_0)\}$ denote the $n \times n$ diagonal matrix of marginal variances, and $\mathbf{X}_0 = (\mathbf{1}, \mathbf{X})$ the $n \times (p + 1)$ null model design matrix. Define $\mathbf{P} = \mathbf{V}_0 - \mathbf{V}_0 \mathbf{X}_0 (\mathbf{X}'_0 \mathbf{V}_0 \mathbf{X}_0)^{-1} \mathbf{X}'_0 \mathbf{V}_0$, which is the asymptotic covariance matrix $\text{Cov}(\mathbf{y} - \hat{\boldsymbol{\pi}}_0)$. Under null, Q follows a mixture of 1-DF chi-square distributions (Liu *et al.*, 2007; Tzeng and Zhang, 2007), with the mixture coefficients being the eigen values of $\mathbf{P}^{1/2} \mathbf{G} \mathbf{W} \mathbf{W} \mathbf{G}' \mathbf{P}^{1/2}$, which is of dimension $n \times n$. The p-value can be obtained by matching moments (Liu *et al.*, 2009) or by inverting the characteristic function (Davies, 1980).

The SKAT statistic can be equivalently derived based on the score vector \mathbf{U} for $\boldsymbol{\beta}$ (Pan, 2009). We can check that $\mathbf{U} = \mathbf{G}'(\mathbf{y} - \hat{\boldsymbol{\pi}}_0)$. Under null, the score vector \mathbf{U} are asymptotically zero-mean multivariate normal with covariance that can be consistently estimated by (Cox and Hinkley, 1979)

$$\boldsymbol{\Sigma} = \mathbf{G}' \mathbf{V}_0 \mathbf{G} - \mathbf{G}' \mathbf{V}_0 \mathbf{X}_0 (\mathbf{X}'_0 \mathbf{V}_0 \mathbf{X}_0)^{-1} \mathbf{X}'_0 \mathbf{V}_0 \mathbf{G} = \mathbf{G}' \mathbf{P} \mathbf{G}, \quad (2)$$

which accounts for the linkage disequilibrium among variants. The SKAT statistic can be equivalently written as $Q = \mathbf{U}' \mathbf{W} \mathbf{W} \mathbf{U}$. Hence the mixture coefficients can be equivalently computed based on the eigen values of $\boldsymbol{\Sigma}^{1/2} \mathbf{W} \mathbf{W} \boldsymbol{\Sigma}^{1/2}$, which is an $m \times m$ matrix. Note that m is typically much smaller than n , and the eigen values can be very efficiently solved.

Likelihood ratio test based kernel association test

For the score vector $\mathbf{U} = \mathbf{G}'(\mathbf{y} - \hat{\boldsymbol{\pi}}_0)$, consider its j -th element $U_j = \mathbf{G}'_j(\mathbf{y} - \hat{\boldsymbol{\pi}}_0)$, where $\mathbf{G}_j = (g_{1j}, \dots, g_{nj})'$ is the j -th column of \mathbf{G} . Here U_j can be checked equal to the score statistic for testing the significance of the j -th SNP based on the following marginal logistic model

$$\Pr(y_i=1|\mathbf{X}_i, g_{ij})=\text{expit}(\beta_{0j}+\mathbf{X}_i\boldsymbol{\alpha}_j+g_{ij}\beta_{1j}). \quad (3)$$

Alternatively we can employ the likelihood ratio test (LRT) to assess the marginal significance of the j -th SNP. Under null, the score test is asymptotically equivalent to the LRT. However the LRT could be more powerful than the score test if the j -th rare variant has potentially large effect size, when it is either a risk variant or in linkage disequilibrium with other risk variants.

We propose to develop a marginal LRT based sequence kernel association test (denoted as SKAT_L) as follows. Denote χ_j as the LRT chi-square statistic for testing β_{1j} under model (3). Let $S_j=\text{sign}(\hat{\beta}_{1j})\sqrt{\chi_j}$ and $\mathbf{S}=(S_1, \dots, S_m)'$, where $\hat{\beta}_{1j}$ is the maximum likelihood estimator (MLE). Define the SKAT_L statistic

$$L=\mathbf{S}'\mathbf{W}\mathbf{W}\mathbf{S}=\sum_{j=1}^m w_j^2 \chi_j.$$

Under the null of no rare variant effects (all $\beta_j=0$), we have $\beta_{1j}=0$, and S_j is asymptotically equivalent to the standardized U_j . Let $\mathbf{R}=\text{diag}(\boldsymbol{\Sigma})^{-1/2}\boldsymbol{\Sigma}\text{diag}(\boldsymbol{\Sigma})^{-1/2}$, which is the corresponding correlation matrix of $\boldsymbol{\Sigma}$ in (2). The null distribution of L is a mixture of 1-DF chi-square distributions with mixture coefficients being the eigen values of $\mathbf{R}^{1/2}\mathbf{W}\mathbf{W}\mathbf{R}^{1/2}$.

Note that the SKAT_L only depends on the LRT chi-square statistic, and in principle we do not need the MLE $\hat{\beta}_{1j}$, which could have convergence issues and aberrant testing behavior (Hauck and Donner, 1977). When computing the SKAT_L in our numerical studies, we set χ_j equal to the squared standardized score statistics for extremely rare variants (specifically with minor allele count less than ten).

Data adaptive kernel association test

An alternative approach to aggregating signals across rare variants is the burden test (Li and Leal, 2008; Madsen and Browning, 2009). The burden test is typically computed as the weighted sum of score statistics. the burden test works well for variants with similar effects and could lose substantial power in the presence of large number of non-causal variants, or with both protective and deleterious variants. A more flexible approach is to data adaptively combine the burden test and the kernel association test following the SKAT-O approach of Lee *et al.* (2012), which tested the rare variant effects using the minimum p-value of weighted SKAT statistic, $(\mathbf{y}-\boldsymbol{\pi}_0)'\mathbf{K}_\rho(\mathbf{y}-\boldsymbol{\pi}_0)$, where $\mathbf{K}_\rho=\mathbf{G}\mathbf{W}[(1-\rho)\mathbf{I}+\rho\mathbf{J}]\mathbf{W}\mathbf{G}'$, $\rho\in[0, 1]$. Here \mathbf{I} is an $m\times m$ identity matrix and \mathbf{J} $m\times m$ matrix with all elements equal to one.

Similarly we consider the following weighted SKAT_L statistic

$$L_\rho=\mathbf{S}'\mathbf{W}[(1-\rho)\mathbf{I}+\rho\mathbf{J}]\mathbf{W}\mathbf{S}, \rho\in[0, 1].$$

Given ρ , the significance p-value of L_ρ , $P\text{-val}(L_\rho)$, can be similarly computed based on the 1-DF chi-square mixture distribution with coefficients being the eigen values of $\mathbf{R}^{1/2}\mathbf{W}[(1-\rho)\mathbf{I}+\rho\mathbf{J}]\mathbf{W}\mathbf{R}^{1/2}$.

$\rho)\mathbf{I} + \rho\mathbf{J}\mathbf{W}\mathbf{R}^{1/2}$. Data adaptive SKAT_L statistic (denoted as SKAT-O_L) is defined as the minimum p-value, $T = \min_{0 \leq \rho \leq 1} P\text{-val}(L_\rho)$, where the minimum is often taken over a finite grid of ρ : $0 = \rho_1 < \dots < \rho_b = 1$, and the significance of T can be efficiently computed using an one-dimensional numerical integration (Lee *et al.*, 2012). We discuss computational details in the following section.

P-value computation for kernel association tests

We offer some insights into the efficient p-value computation for SKAT, SKAT_L and data adaptive kernel association tests. First note that the non-zero eigen values of $\mathbf{A}\mathbf{A}'$ are the same as $\mathbf{A}'\mathbf{A}$ for any matrix \mathbf{A} , which can be verified from the singular value decomposition of matrix \mathbf{A} : $\mathbf{A}=\mathbf{U}_A \mathbf{D}_A \mathbf{V}_A'$, where \mathbf{U}_A and \mathbf{V}_A are orthogonal and \mathbf{D}_A diagonal matrix.

Therefore $\mathbf{A}\mathbf{A}'=\mathbf{U}_A \mathbf{D}_A^2 \mathbf{U}_A'$ and $\mathbf{A}'\mathbf{A}=\mathbf{V}_A \mathbf{D}_A^2 \mathbf{V}_A'$ and hence their eigen values equal to the squared singular values of \mathbf{A} . So for computing the p-values of proposed SKAT_L, the eigen values of $\mathbf{R}^{1/2}\mathbf{W}\mathbf{W}\mathbf{R}^{1/2}$ can be equivalently computed from $\mathbf{W}\mathbf{R}\mathbf{W}$. For SKAT, the eigen values of $\mathbf{P}^{1/2}\mathbf{G}\mathbf{W}\mathbf{W}\mathbf{G}'\mathbf{P}^{1/2}$ are the same as $\mathbf{W}\mathbf{G}'\mathbf{P}\mathbf{G}\mathbf{W} = \mathbf{W}\Sigma\mathbf{W}$.

For matrix $\mathbf{B} = (1 - \rho)\mathbf{I} + \rho\mathbf{J}$, $\rho \in [0, 1]$, we can check that $\mathbf{B}=\mathbf{B}_h^2$, where

$\mathbf{B}_h = \sqrt{1 - \rho}\mathbf{I} + \frac{\sqrt{1+(m-1)\rho} - \sqrt{1-\rho}}{m} \mathbf{J}$. Therefore for computing p-values of weighted SKAT_L, the eigen values of $\mathbf{R}^{1/2}\mathbf{W}[(1 - \rho)\mathbf{I} + \rho\mathbf{J}]\mathbf{W}\mathbf{R}^{1/2}$ can be equivalently computed from $\mathbf{B}_h\mathbf{W}\mathbf{R}\mathbf{W}\mathbf{B}_h$.

Null distribution of SKAT-O_L

The significance of SKAT-O_L can be computed as (see Appendix for technical details)

$$1 - \int_0^{\tilde{q}_1} M(\delta(x))f(x|\chi_1^2)dx,$$

where

$$\delta(x) = \left(\min_{v < b} \frac{q_{\rho v} - \tau_{\rho v} x}{1 - \rho v} - \mu \right) \frac{\sigma}{\sigma_0} + \mu, \tilde{q}_1 = F^{-1}(1 - T|\chi_1^2),$$

$f(\cdot|\chi_1^2)$ and $F(\cdot|\chi_1^2)$ are the 1-DF chi-square density/distribution functions, and $M(\cdot)$ is the distribution function of 1-DF chi-square mixture with coefficients $(\lambda_1, \dots, \lambda_m)$, which are the eigen values of $(\mathbf{I} - H_1) \mathbf{R}(\tilde{\mathbf{I}} - H_1)$, where $\mathbf{R} = \mathbf{W}\mathbf{R}\mathbf{W}$. Here

$$\mu = \sum_{j=1}^m \lambda_j, \sigma^2 = 2 \sum_{j=1}^m \lambda_j^2, \sigma_0^2 = \sigma^2 + 4tr[\tilde{\mathbf{R}}\mathbf{H}_1\tilde{\mathbf{R}}(\mathbf{I} - H_1)], \tau_\rho = \rho\|\mathbf{R}_1\|^2 + (1-\rho)\mathbf{R}'_1\tilde{\mathbf{R}}\mathbf{R}_1/\|\mathbf{R}_1\|^2, H_1 = \mathbf{R}_h\mathbf{J}\mathbf{R}_h/(\mathbf{R}'_1\mathbf{I}\mathbf{R}_1),$$

where $\mathbf{R}_h\mathbf{R}_h = \mathbf{R}$, and $\mathbf{R}_1 = \mathbf{R}_h(1, \dots, 1)'$.

Results

Simulation studies

We conducted extensive simulation studies to evaluate the performance of the proposed and existing methods. Following Lee *et al.* (2012), we generated 10,000 European-like haplotypes of length 1000 kb under a calibrated coalescent model (Schaffner *et al.*, 2005). We randomly pair the haplotypes to simulate a total population of 10^6 individuals. We randomly select a gene region of length 10 kb and study those rare variants with MAF 0.01. We consider two covariates $Z = (Z_1, Z_2)'$: $Z_1 \in \{0, 1\}$ follows Bernoulli(0.5), and $Z_2 \sim N(0, 1)$. We model the logit disease risk as $\text{expit}(\beta_0 + Z' \beta_Z + \sum_{j=1}^m \beta_j G_j)$. We set $\beta_0 = -3.4$, $\beta_Z = (0.5, 0.5)'$ (corresponding to 5% population disease rate). We randomly select 2500 cases and 2500 controls from the simulated population of 10^6 samples. We compared five rare variant set analysis methods: SKAT, SKAT-O, SKAT_L, SKAT-O_L and burden test. In the burden and SKAT tests, we assign weight $\text{Beta}(f_j; a_0, b_0)$ to the j th variant G_j . And for the proposed method we assign weight $\text{Beta}(f_j; a_1, b_1)$. Here f_j is the MAF of G_j . For a given variant, the likelihood ratio test statistic is inherently standardized and roughly corresponds to the standardized score statistics, which is the score statistics used in SKAT scaled by its standard error, which is roughly proportional to $\sqrt{f_j(1-f_j)}$. Therefore for the proposed method, we set $a_1 = a_0 + 0.5$ and $b_1 = b_0 + 0.5$. Following Wu *et al.* (2011), we set $a_0 = 1$, $b_0 = 25$ for the following simulation studies. We have investigated three sets of weights for (a_0, b_0) : (0.5, 24.5), (1, 25), and (1.5, 25.5). The overall conclusions remain the same (see supplementary material for complete results). As shown in Ma *et al.* (2013), the performance of single rare variant LRT depends on the case-control ratio. We have investigated different case-control ratios for (n_e, n_c) . Here we reported the results for $n_e = n_c = 2500$, and $n_e = 1700$, $n_c = 3300$. The supplementary material provided simulation results for more unbalanced case-control ratios (1:6 and 1:10).

We use 2.5×10^6 experiments to evaluate the type I error at the nominal significance level $\alpha = 10^{-5}$, 10^{-4} , and 10^{-3} by setting all $\beta_j = 0$. The results are summarized in Table 1 and 2. All methods appropriately control the Type I errors. We also verify that the Type I errors are appropriately controlled at the 10^{-6} significance level by conducting 10^8 experiments (please see the supplementary material for detailed results including the QQ plots).

We use 10^4 experiments to evaluate the power under various combinations of β_j at $\alpha = 10^{-6}$, 10^{-5} , 10^{-4} , and 10^{-3} . The rare variant effects β_j are set as follows. Each time we randomly select θ proportion of rare variants and set their $|\beta_j| = d \log_{10}(f_j)$. The other null rare variants have zero coefficients. We have assumed that rarer variants have larger effect sizes. We conducted simulations for (1) $\theta = 0.05$, $d = -0.6$, (2) $\theta = 0.1$, $d = -0.5$, (3) $\theta = 0.2$, $d = -0.4$, (4) $\theta = 0.5$, $d = -0.25$. They correspond to odds ratio of 3.32, 2.72, 2.23 and 1.65 for MAF=0.01 respectively. We have investigated two scenarios for the direction of causal variant effects. First, we assume a mix of equal proportions of protective and deleterious variants, which will in general favor the kernel association test. Second, we assume a mix of unequal proportions of protective and deleterious variants. Specially we randomly set signs of β_j as negative or positive with probability 0.9 and 0.1 respectively.

Table 3 summarized the results assuming equal proportions of protective and deleterious variants and equal case-control ratio ($n_e = n_c = 2500$). Overall the proposed SKAT_L has the best performance. As expected the burden test suffers a dramatic power loss since the burden sum score cancels out those causal variants, and as a result the adaptive SKAT-O and SKAT-O_L have reduced performance compared to the SKAT and SKAT_L. The proposed SKAT_L has the largest power gain over SKAT with relatively large rare variant effect sizes.

Table 4 summarized the results assuming unequal proportions of protective and deleterious variants and equal case-control ratio ($n_e = n_c = 2500$). The adaptive SKAT-O and SKAT-O_L now perform better than the SKAT and SKAT_L under relatively more causal variants with $\theta = 0.5$. With small proportion of causal variants, the burden test suffered much power loss, and as a result the SKAT and SKAT_L performed better than the adaptive SKAT-O and SKAT-O_L.

Table 5 and 6 summarized the corresponding power results under unequal case-control ratio: $n_e = 1700$, $n_c = 3300$. When there are equal proportion of protective and deleterious variants, the proposed LRT based SKAT offered more improvement compared to the score test based SKAT with equal case-control ratio (Table 3 versus 5). While under unequal proportion of protective and deleterious variants, the proposed LRT based SKAT offered more improvement compared to the score test based SKAT with unequal case-control ratio (Table 4 versus 6). The results are in agreement with the observations of Ma *et al.* (2013), who showed that the performance of single rare variant LRT depends on the case-control ratio.

Diabetes study

The Atherosclerosis Risk in Communities (ARIC) study (The ARIC Investigators, 1989) is a multi-center prospective investigation of atherosclerotic disease in a predominantly bi-racial population. Men and women aged 45–64 years at baseline were recruited from four U.S. communities: Forsyth County, North Carolina; Jackson, Mississippi; suburban areas of Minneapolis, Minnesota; and Washington County, Maryland. A total of 15,792 individuals participated in the baseline examination in 1987–1989. The vast majority of ARIC participants are of European (73%) or African ancestry (26%).

We applied the proposed SKAT_L and other competing methods in ARIC to test for association between type 2 diabetes (T2D) and rare variants in each gene. Genotypes were obtained from the Illumina HumanExome BeadChip (Grove *et al.*, 2013), which has information on 247,870 variants. Prevalent T2D diabetes was defined as in previous GWAS analyses using phenotypic information collected at the baseline examination (Morris *et al.*, 2012). Exome chip data were analyzed for 1048 white T2D cases and 6598 white non-cases.

We conducted two different analyses of T2D and adjusted for age, gender and center. First, we analyzed the rare variants (with MAF ≤ 0.01 and at least five copies in the total sample) in the gene *PAM*, which has been recently identified to contain a rare missense variant that contributes to the risk of T2D (Steinthorsdottir *et al.*, 2014). Second, we ran a genome-wide scan and tested the association of rare variants located in each gene.

For the eight rare variants located in the gene *PAM* and available on the exome chip, the proposed $SKAT_L$ has a p-value of 0.039, and SKAT's p-value is 0.115. The burden test has a p-value of 0.894. For the data adaptive tests, the proposed $SKAT-O_L$ has a p-value of 0.072, and SKAT-O's p-value is 0.196.

In total we analyzed 11426 rare variant sets in the genome-wide scan for T2D. $SKAT_L$ identified a significant set with three rare variants in the gene *ZZZ3* (p-value= 1.4×10^{-6}) that passed genome-wide significance after a Bonferonni correction for the total number of sets (4.4×10^{-6}). And the other tests did not identify any significant rare variant set. SKAT reported a p-value of 2.7×10^{-5} for the gene *ZZZ3*, and did not identify any genome-wide significant rare variant set. *ZZZ3* is a protein-coding gene which is a component of the ATAC complex, a complex with histone acetyltransferase activity on histones H3 and H4. A common variant in *ZZZ3* was recently found to be associated with obesity and body mass index in a genome-wide meta-analysis of 263,407 European individuals (Berndt *et al.*, 2013). Obesity is a major risk factor for T2D. This suggests that *ZZZ3* is likely involved in T2D. Further research is needed on the possible role of identified rare variants in the gene *ZZZ3*.

Discussion

To enrich association signals for rare variants, it is attractive and often customary to combine multiple rare variants in a gene. The widely used SKAT is powerful and computationally efficient by combining rare variants based on the variance component score test. The proposed $SKAT_L$ is based on the observation that the score statistics used in SKAT are asymptotically equivalent to the LRT statistics in the marginal regression modeling of individual rare variants, and that the score test performs well when parameter is close to the null value, but could have suboptimal performance with large deviation from the null (e.g., when testing those rare variants with large effect sizes). We developed efficient algorithms to compute p-values based on the asymptotic distribution of the proposed $SKAT_L$ and $SKAT-O_L$. In our extensive numerical studies, the proposed $SKAT_L$ and $SKAT-O_L$ have well controlled type I errors and shown very competitive performance.

Our approach is in the same spirit as Xing *et al.* (2012), Ma *et al.* (2013) and Chen *et al.* (2014), who have shown that the likelihood ratio test often has better performance than the score test for either single rare variant or rare variant set analysis. In practice, the score test has the computational advantage in that we only need to fit one null model. For the ARIC diabetes data, when analyzing 1415 rare variant sets on chromosome 1 on a single Linux workstation, SKAT takes 42 sec CPU time, SKAT-O takes 674 sec CPU time, $SKAT_L$ takes 151 sec CPU time, and $SKAT-O_L$ takes 630 sec CPU time on the same machine. In the supplementary material, we provide more time comparison of score test versus LRT based SKAT in numerical studies. The proposed approach can be readily extended to handle across study meta analyses of gene-level tests, and the analysis of multiple traits. In summary, we advocate using the proposed method as a complementary approach to enhancing the power of detecting association for rare variants in case-control genome-wide association studies.

We have implemented the proposed methods in R programs posted at http://www.umn.edu/~baolin/research/skatl_Rcode.html

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgements

This research was supported in part by NIH grant GM083345 and CA134848. We are grateful to the University of Minnesota Supercomputing Institute for assistance with the computations. We want to thank the editor and reviewers for their constructive comments which have greatly improved the presentation of the paper.

The ARIC Study is carried out as a collaborative study supported by National Heart, Lung, and Blood Institute (NHLBI) contracts (HHSN268201100005C, HHSN268201100006C, HHSN268201100007C, HHSN268201100008C, HHSN268201100009C, HHSN268201100010C, HHSN268201100011C, and HHSN268201100012C). The authors thank the staff and participants of the ARIC study for their important contributions. Support for exome chip genotyping in the ARIC Study was provided by the National Institutes of Health (NIH) American Recovery and Reinvestment Act of 2009 (ARRA) (5RC2HL102419).

References

- Berndt SI, Gustafsson S, Magi R, Ganna A, Wheeler E, Feitosa MF, et al. Genome-wide meta-analysis identifies 11 new loci for anthropometric traits and provides insights into genetic architecture. *Nature genetics*. 2013; 45(5):501–512. [PubMed: 23563607]
- Bonnefond A, Clement N, Fawcett K, Yengo L, Vaillant E, Guillaume JL, et al. Rare MTNR1B variants impairing melatonin receptor 1B function contribute to type 2 diabetes. *Nature genetics*. 2012; 44(3):297–301. [PubMed: 22286214]
- Chen H, Lumley T, Brody J, Heard-Costa NL, Fox CS, Cupples LA, Dupuis J. Sequence kernel association test for survival traits. *Genetic Epidemiology*. 2014; 38(3):191–197. [PubMed: 24464521]
- Cox, DR.; Hinkley, DV. *Theoretical Statistics*. CRC Press; 1979.
- Davies RB. Algorithm AS 155: the distribution of a linear combination of χ^2 random variables. *Applied Statistics*. 1980; 29(3):323.
- Estrada K, Aukrust I, Bjorkhaug L, Burt NP, Mercader JM, et al. Association of a low-frequency variant in *HNF1A* with type 2 diabetes in a latino population. *JAMA*. 2014; 311(22):2305–2314. [PubMed: 24915262]
- Grove ML, Yu B, Cochran BJ, Haritunians T, Bis JC, Taylor KD, et al. Best practices and joint calling of the HumanExome BeadChip: the CHARGE consortium. *PloS One*. 2013; 8(7):e68095. [PubMed: 23874508]
- Hauck WW, Donner A. Wald's test as applied to hypotheses in logit analysis. *Journal of the American Statistical Association*. 1977; 72(360):851.
- Lee S, Wu MC, Lin X. Optimal tests for rare variant effects in sequencing association studies. *Biostatistics*. 2012; 13(4):762–775. [PubMed: 22699862]
- Li B, Leal SM. Methods for detecting associations with rare variants for common diseases: application to analysis of sequence data. *The American Journal of Human Genetics*. 2008; 83(3):311–321. [PubMed: 18691683]
- Lin DY, Tang ZZ. A general framework for detecting disease associations with rare variants in sequencing studies. *The American Journal of Human Genetics*. 2011; 89(3):354–367. [PubMed: 21885029]
- Liu D, Lin X, Ghosh D. Semiparametric regression of multidimensional genetic pathway data: least-squares kernel machines and linear mixed models. *Biometrics*. 2007; 63(4):1079–1088. [PubMed: 18078480]

- Liu H, Tang Y, Zhang HH. A new chi-square approximation to the distribution of non-negative definite quadratic forms in non-central normal variables. *Computational Statistics & Data Analysis*. 2009; 53(4):853–856.
- Ma C, Blackwell T, Boehnke M, Scott LJ. GoT2D Investigators. Recommended joint and meta-analysis strategies for case-control association testing of single low-count variants. *Genetic Epidemiology*. 2013; 37(6):539–550. [PubMed: 23788246]
- Madsen BE, Browning SR. A groupwise association test for rare mutations using a weighted sum statistic. *PLOS Genetics*. 2009; 5(2):e1000384. [PubMed: 19214210]
- Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorf LA, et al. Finding the missing heritability of complex diseases. *Nature*. 2009; 461(7265):747–753. [PubMed: 19812666]
- Morgenthaler S, Thilly WG. A strategy to discover genes that carry multi-allelic or mono-allelic risk for common diseases: a cohort allelic sums test (CAST). *Mutation research*. 2007; 615(1–2):28–56. [PubMed: 17101154]
- Morris AP, Voight BF, Teslovich TM, Ferreira T, Segre AV, et al. Large-scale association analysis provides insights into the genetic architecture and pathophysiology of type 2 diabetes. *Nature Genetics*. 2012; 44(9):981–990. [PubMed: 22885922]
- Morris AP, Zeggini E. An evaluation of statistical approaches to rare variant analysis in genetic association studies. *Genetic epidemiology*. 2010; 34(2):188–193. [PubMed: 19810025]
- Neale BM, Rivas MA, Voight BF, Altshuler D, Devlin B, Orho-Melander M, Kathiresan S, Purcell SM, Roeder K, Daly MJ. Testing for an unusual distribution of rare variants. *PLoS Genet*. 2011; 7(3):e1001322. [PubMed: 21408211]
- Pan W. Asymptotic tests of association with multiple SNPs in linkage disequilibrium. *Genetic Epidemiology*. 2009; 33(6):497–507. [PubMed: 19170135]
- Price AL, Kryukov GV, de Bakker PIW, Purcell SM, Staples J, Wei LJ, Sunyaev SR. Pooled association tests for rare variants in exon-resequencing studies. *American journal of human genetics*. 2010; 86(6):832–838. [PubMed: 20471002]
- Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, Maller J, Sklar P, Bakker PIWd, Daly MJ, Sham PC. PLINK: a tool set for whole-genome association and population-based linkage analyses. *The American Journal of Human Genetics*. 2007; 81(3):559–575. [PubMed: 17701901]
- Schaffner SF, Foo C, Gabriel S, Reich D, Daly MJ, Altshuler D. Calibrating a coalescent simulation of human genome sequence variation. *Genome Research*. 2005; 15(11):1576–1583. [PubMed: 16251467]
- Steinthorsdottir V, Thorleifsson G, Sulem P, Helgason H, Grarup N, et al. Identification of low-frequency and rare sequence variants associated with elevated or reduced risk of type 2 diabetes. *Nature Genetics*. 2014; 46(3):294–298. [PubMed: 24464100]
- The ARIC Investigators. The atherosclerosis risk in communities (aric) study: design and objectives. *American Journal of Epidemiology*. 1989; 129(4):687–702. [PubMed: 2646917]
- Tzeng JY, Zhang D. Haplotype-based association analysis via variance-components score test. *American Journal of Human Genetics*. 2007; 81(5):927–938. [PubMed: 17924336]
- Wang Y, McKay JD, Rafnar T, Wang Z, Timofeeva MN, Broderick P, et al. Rare variants of large effect in BRCA2 and CHEK2 affect risk of lung cancer. *Nature Genetics*. 2014; 46(7):736–741. [PubMed: 24880342]
- Wu MC, Kraft P, Epstein MP, Taylor DM, Chanock SJ, Hunter DJ, Lin X. Powerful SNP-Set analysis for case-control genome-wide association studies. *The American Journal of Human Genetics*. 2010; 86(6):929–942. [PubMed: 20560208]
- Wu M, Lee S, Cai T, Li Y, Boehnke M, Lin X. Rare-variant association testing for sequencing data with the sequence kernel association test. *The American Journal of Human Genetics*. 2011; 89(1):82–93. [PubMed: 21737059]
- Xing G, Lin CY, Wooding SP, Xing C. Blindly using wald's test can miss rare disease-causal variants in case-control association studies. *Annals of human genetics*. 2012; 76(2):168–177. [PubMed: 22256951]

Zhan X, Larson DE, Wang C, Koboldt DC, Sergeev YV, Fulton RS, et al. Identification of a rare coding variant in complement 3 associated with age-related macular degeneration. *Nature Genetics*. 2013; 45(11):1375–1379. [PubMed: 24036949]

APPENDIX

Null distribution of SKAT- O_L

The significance of SKAT- O_L can be computed following the approach of Lee *et al.* (2012). Denote $\tilde{\mathbf{R}} = \mathbf{W}\mathbf{R}\mathbf{W}$. Define a symmetric matrix R_h such that $R_h R_h = \tilde{\mathbf{R}}$. Let $\mathbf{Z} = (z_1, \dots, z_m)'$ be independent standard normal random variables. Then the null distribution of L_ρ is the same as $L_\rho = \mathbf{Z}' R_h [(1 - \rho)\mathbf{I} + \rho\mathbf{J}] R_h \mathbf{Z}$. Denote $R_1 = R_h \mathbf{1}$, where $\mathbf{1} = (1, \dots, 1)'$ is a column vector of ones. Note that $H_1 = R_h \mathbf{J} R_h / (R_1' R_1)$ is a projection matrix into a space spanned by R_1 . Therefore $Z_1 = H_1 \mathbf{Z}$ and $Z_2 = (\mathbf{I} - H_1) \mathbf{Z}$ are independent. Define

$\eta_2 = Z_2' \tilde{\mathbf{R}} Z_2$, $\eta_1 = Z_1' \tilde{\mathbf{R}} Z_1$, and $\eta_0 = Z_1' Z_1$. Here η_2 follows a mixture of 1-DF chi-square distributions with coefficients being the eigen values of $(\mathbf{I} - H_1) \tilde{\mathbf{R}} (\mathbf{I} - H_1)$, denoted as $(\lambda_1, \dots, \lambda_m)$. Note $\text{Cov}(\eta_1, \eta_2) = \text{Cov}(\eta_1, \eta_0) = 0$, and $E(\eta_1) = 0$, $\text{Var}(\eta_1) = \text{tr}[\tilde{\mathbf{R}} H_1 \tilde{\mathbf{R}} (\mathbf{I} - H_1)]$. We can check that

$$L_\rho = (1 - \rho)(\eta_2 + 2\eta_1) + \tau_\rho \eta_0, \tau_\rho = \rho \|R_1\|^2 + (1 - \rho) R_1' \tilde{\mathbf{R}} R_1 / \|R_1\|^2.$$

Let $L_{\rho_1}, \dots, L_{\rho_b}$ be the score statistics computed with $0 = \rho_1 < \rho_2 < \dots < \rho_b = 1$. Denote q_ρ as the $(1 - T)$ -th percentile of the distribution of L_ρ , which can be computed based on moment matching (Liu *et al.*, 2009). Let $\tilde{q}_1 = F^{-1}(1 - T | \chi_1^2)$, where $F(\cdot | \chi_1^2)$ is the distribution function of 1-DF chi-square distribution. Note that $L_1 = \|R_1\|^2 \eta_0$. Hence $q_1 = \|R_1\|^2 \tilde{q}_1$. The significance p-value based on the test statistic T is

$$1 - \Pr(L_{\rho_1} < q_{\rho_1}, \dots, L_{\rho_b} < q_{\rho_b}) = 1 - E \left[\Pr \left(\eta_2 + 2\eta_1 < \min_{v < b} \frac{q_{\rho_v} - \tau_{\rho_v} \eta_0}{1 - \rho_v} \mid \eta_0 \right) I(\eta_0 < \tilde{q}_1) \right],$$

where η_0 follows the 1-DF chi-square distribution, and $I(\cdot)$ is an indicator function. Denote $\mu = \sum_{j=1}^m \lambda_j$, $\sigma^2 = 2 \sum_{j=1}^m \lambda_j^2$, $\sigma_0^2 = \sigma^2 + 4 \text{tr}[\tilde{\mathbf{R}} H_1 \tilde{\mathbf{R}} (\mathbf{I} - H_1)]$. Let

$$\delta(x) = \left(\min_{v < b} \frac{q_{\rho_v} - \tau_{\rho_v} x}{1 - \rho_v} - \mu \right) \frac{\sigma}{\sigma_0} + \mu.$$

The p-value is computed as

$$1 - \int_0^{\tilde{q}_1} M(\delta(x)) f(x | \chi_1^2) dx,$$

where $f(\cdot | \chi_1^2)$ is the density of 1-DF chi-square distribution, and $M(\cdot)$ is the distribution function of 1-DF chi-square mixture with coefficients $(\lambda_1, \dots, \lambda_m)$. Here we want to emphasize that special care is needed for $\rho = 1$. When $\rho_b = 1$ is included in the minimum p-

value search, we have an indicator $I(\eta_0 < q_1^{\tilde{}})$ in the expectation, and the integration is in interval $[0, q_1^{\tilde{}}]$. Otherwise the integration is over $[0, q_1^{\tilde{}} = \infty)$.

Table 1

Type I error divided by the nominal significance level for rare variant set analysis: $n_e = 2500$ cases and $n_c = 2500$ controls. The SKAT/SKAT-O and burden tests used (1,25) weight, and the proposed SKAT_L/SKAT-O_L used (1.5,25.5) weight.

α	10^{-5}	10^{-4}	10^{-3}
SKAT	0.82	0.85	0.92
SKAT-O	0.91	1.02	1.07
SKAT _L	0.92	1.10	1.08
SKAT-O _L	0.96	1.00	1.11
Burden	0.94	0.98	1.00

Table 2

Type I error divided by the nominal significance level for rare variant set analysis: $n_e = 1700$ cases and $n_c = 3300$ controls. The SKAT/SKAT-O and burden tests used (1,25) weight, and the proposed SKAT_L/SKAT-O_L used (1.5,25.5) weight.

α	10^{-5}	10^{-4}	10^{-3}
SKAT	0.86	0.91	0.93
SKAT-O	0.89	0.98	1.04
SKAT _L	0.96	1.12	1.11
SKAT-O _L	0.88	1.07	1.10
Burden	0.91	0.97	1.01

Power comparison of rare variant set analysis: $n_e = n_c = 2500$, equal proportions of protective and deleterious variants. The highest powered tests in each row are bold-faced.

Table 3

$\theta = 0.05, d = -0.6$				
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L Burden
10^{-6}	0.1654	0.1373	0.2000	0.1661 0.0028
10^{-5}	0.2279	0.1995	0.2627	0.2336 0.0067
10^{-4}	0.3080	0.2808	0.3521	0.3191 0.0181
10^{-3}	0.4286	0.3994	0.4740	0.4398 0.0451
$\theta = 0.1, d = -0.5$				
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L Burden
10^{-6}	0.2469	0.2041	0.2940	0.2442 0.0081
10^{-5}	0.3446	0.3030	0.3906	0.3506 0.0164
10^{-4}	0.4612	0.4260	0.5051	0.4691 0.0352
10^{-3}	0.6031	0.5657	0.6453	0.6092 0.0750
$\theta = 0.2, d = -0.4$				
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L Burden
10^{-6}	0.3481	0.2952	0.3965	0.3381 0.0161
10^{-5}	0.4742	0.4239	0.5189	0.4695 0.0302
10^{-4}	0.6122	0.5698	0.6513	0.6130 0.0634
10^{-3}	0.7577	0.7283	0.7885	0.7613 0.1174
$\theta = 0.5, d = -0.25$				
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L Burden
10^{-6}	0.2959	0.2409	0.3379	0.2783 0.0139
10^{-5}	0.4306	0.3796	0.4724	0.4196 0.0279

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

10^{-4}	0.5871	0.5432	0.6272	0.5827	0.0568
10^{-3}	0.7554	0.7234	0.7842	0.7542	0.1125

Power comparison of rare variant set analysis: $n_e = n_c = 2500$, unequal proportions of protective and deleterious variants. The highest powered tests are bold-faced.

Table 4

$\theta = 0.05, d = -0.6$					
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L	Burden
10^{-6}	0.0858	0.0699	0.1054	0.0842	0.0031
10^{-5}	0.1287	0.1110	0.1508	0.1290	0.0077
10^{-4}	0.1858	0.1659	0.2096	0.1878	0.0162
10^{-3}	0.2781	0.2492	0.3026	0.2785	0.0367
$\theta = 0.1, d = -0.5$					
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L	Burden
10^{-6}	0.1475	0.1257	0.1725	0.1452	0.0112
10^{-5}	0.2114	0.1887	0.2390	0.2128	0.0214
10^{-4}	0.3059	0.2768	0.3373	0.3061	0.0405
10^{-3}	0.4278	0.3988	0.4587	0.4319	0.0820
$\theta = 0.2, d = -0.4$					
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L	Burden
10^{-6}	0.2510	0.2321	0.2796	0.2598	0.0522
10^{-5}	0.3389	0.3200	0.3756	0.3528	0.0820
10^{-4}	0.4587	0.4446	0.4938	0.4768	0.1311
10^{-3}	0.6030	0.5961	0.6349	0.6248	0.2170
$\theta = 0.5, d = -0.25$					
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L	Burden
10^{-6}	0.3220	0.3804	0.3577	0.4076	0.2244
10^{-5}	0.4310	0.5050	0.4675	0.5348	0.3132

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

10^{-4}	0.5657	0.6483	0.5964	0.6722	0.4256
10^{-3}	0.7177	0.7879	0.7443	0.8063	0.5649

Power comparison of rare variant set analysis: $n_e = 1700$, $n_c = 3300$, equal proportions of protective and deleterious variants. The highest powered tests are bold-faced.

Table 5

$\theta = 0.05, d = -0.6$				
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L Burden
10^{-6}	0.1624	0.1363	0.1655	0.1356 0.0052
10^{-5}	0.2171	0.1913	0.2244	0.1942 0.0095
10^{-4}	0.2933	0.2640	0.3050	0.2728 0.0229
10^{-3}	0.4027	0.3734	0.4191	0.3850 0.0531
$\theta = 0.1, d = -0.5$				
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L Burden
10^{-6}	0.2364	0.2003	0.2460	0.2036 0.0123
10^{-5}	0.3238	0.2890	0.3337	0.2937 0.0230
10^{-4}	0.4325	0.3896	0.4471	0.4086 0.0488
10^{-3}	0.5700	0.5365	0.5843	0.5488 0.0914
$\theta = 0.2, d = -0.4$				
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L Burden
10^{-6}	0.3103	0.2631	0.3253	0.2666 0.0211
10^{-5}	0.4249	0.3823	0.4414	0.3896 0.0383
10^{-4}	0.5643	0.5240	0.5840	0.5372 0.0702
10^{-3}	0.7175	0.6831	0.7338	0.6946 0.1281
$\theta = 0.5, d = -0.25$				
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L Burden
10^{-6}	0.2638	0.2194	0.2786	0.2251 0.0166
10^{-5}	0.3844	0.3380	0.4022	0.3501 0.0335

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

10^{-4}	0.5355	0.4904	0.5559	0.5093	0.0652
10^{-3}	0.7131	0.6758	0.7303	0.6914	0.1238

Power comparison of rare variant set analysis: $n_e = 1700$, $n_c = 3300$, unequal proportions of protective and deleterious variants. The highest powered tests are bold-faced.

Table 6

$\theta = 0.05, d = -0.6$				
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L Burden
10^{-6}	0.0402	0.0314	0.0633	0.0482 0.0007
10^{-5}	0.0710	0.0564	0.1000	0.0842 0.0022
10^{-4}	0.1176	0.1006	0.1548	0.1355 0.0050
10^{-3}	0.1964	0.1743	0.2370	0.2153 0.0177
$\theta = 0.1, d = -0.5$				
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L Burden
10^{-6}	0.0746	0.0598	0.1093	0.0878 0.0036
10^{-5}	0.1242	0.1027	0.1686	0.1474 0.0078
10^{-4}	0.2009	0.1785	0.2535	0.2306 0.0177
10^{-3}	0.3177	0.2910	0.3784	0.3512 0.0470
$\theta = 0.2, d = -0.4$				
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L Burden
10^{-6}	0.1288	0.1120	0.1856	0.1675 0.0166
10^{-5}	0.2074	0.1922	0.2762	0.2632 0.0355
10^{-4}	0.3185	0.3103	0.3909	0.3816 0.0663
10^{-3}	0.4728	0.4636	0.5423	0.5367 0.1360
$\theta = 0.5, d = -0.25$				
α	SKAT	SKAT-O	SKAT _L	SKAT-O _L Burden
10^{-6}	0.1548	0.2098	0.2339	0.2961 0.1181
10^{-5}	0.2510	0.3333	0.3366	0.4270 0.1913

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

10^{-4}	0.3819	0.4948	0.4700	0.5760	0.3032
10^{-3}	0.5563	0.6708	0.6412	0.7365	0.4561