

# Optimal tests for rare variant effects in sequencing association studies

SEUNGGEUN LEE

*Department of Biostatistics, Harvard School of Public Health, Boston, MA 02115, USA*

MICHAEL C. WU

*Department of Biostatistics, University of North Carolina, Chapel Hill, NC 27599, USA*

XIHONG LIN\*

*Department of Biostatistics, Harvard School of Public Health, Boston, MA 02115, USA*

xlin@hsph.harvard.edu

## SUMMARY

With development of massively parallel sequencing technologies, there is a substantial need for developing powerful rare variant association tests. Common approaches include burden and non-burden tests. Burden tests assume all rare variants in the target region have effects on the phenotype in the same direction and of similar magnitude. The recently proposed sequence kernel association test (SKAT) (Wu, M. C., and others, 2011. Rare-variant association testing for sequencing data with the SKAT. *The American Journal of Human Genetics* 89, 82–93], an extension of the C-alpha test (Neale, B. M., and others, 2011. Testing for an unusual distribution of rare variants. *PLoS Genetics* 7, 161–165], provides a robust test that is particularly powerful in the presence of protective and deleterious variants and null variants, but is less powerful than burden tests when a large number of variants in a region are causal and in the same direction. As the underlying biological mechanisms are unknown in practice and vary from one gene to another across the genome, it is of substantial practical interest to develop a test that is optimal for both scenarios. In this paper, we propose a class of tests that include burden tests and SKAT as special cases, and derive an optimal test within this class that maximizes power. We show that this optimal test outperforms burden tests and SKAT in a wide range of scenarios. The results are illustrated using simulation studies and triglyceride data from the Dallas Heart Study. In addition, we have derived sample size/power calculation formula for SKAT with a new family of kernels to facilitate designing new sequence association studies.

**Keywords:** Burden tests; Correlated effects; Kernel association test; Rare variants; Score test.

## 1. INTRODUCTION

Advances in high-throughput sequencing technology are reshaping the landscape of medical and human genetics research. In contrast to genome-wide association studies (GWAS), which involves genotyping

\*To whom correspondence should be addressed.

preselected single-nucleotide polymorphisms (SNPs) that are relatively common, these new sequencing technologies enable us to also genotype rare genetic variants. Rare genetic variants, here defined as variants with minor allele frequency (MAF)  $< 1\text{--}5\%$ , have been shown to play a crucial role in complex trait etiology (Cohen and others, 2004). Despite their importance, testing for associations between rare variants and traits has proven challenging. Since standard individual variant tests, typically used for analysis of SNPs, are underpowered to detect rare variant effects due to the low allele frequencies and the large numbers of rare variants in the genome, region-based analysis has become the standard approach for analyzing rare variants in sequencing studies (Li and Leal, 2008).

The earliest and most commonly used class of region-based rare variant tests are the burden tests, which collapse or summarize the rare variants within the region as a single genetic variable which can then be tested for association with any trait of interest (Li and Leal, 2008; Madsen and Browning, 2009; Morris and Zeggini, 2010; Price and others, 2010). For example, the combined multivariate and collapsing method (Li and Leal, 2008) collapses information on all rare variants within a region by counting the number of rare alleles. Many extensions and variations on these methods exist. A key limitation of burden tests is that they suffer a substantial loss of power in the presence of large number of non-causal variants or in the presence of both protective and deleterious variants (Neale and others, 2011).

Recognizing the inherent limitations of burden-based methods, Wu and others (2011) recently proposed the sequence kernel association test (SKAT), which builds upon the kernel machine regression framework, to test rare variants associations. As a kernel machine-based test, SKAT aggregates genetic information across the region using a kernel function and uses a computationally efficient variance component test to test for association. Wu and others (2011) also showed that SKAT is a generalization of the classical C-alpha test (Neyman and Scott, 1966; Neale and others, 2011).

Although SKAT offers improved power over burden-based tests in many cases, if a large proportion of the rare variants in a region are truly causal and influence the phenotype in the same direction, then burden tests can have higher power than SKAT (Basu and Pan, 2011). This scenario can arise in several current sequencing studies, e.g. whole exome sequencing studies. This is because standard evolutionary principles and population genetics models indicate that majority of rare missense mutations are moderately deleterious (Kryukov and others, 2007). In addition, bioinformatic tools are often used to restrict testing to variants that are likely harmful variants. However, such prior knowledge is often lacking in practice. The underlying biological mechanisms vary from one gene to another across the genome and often unknown. It is hence of substantial practical interest to develop a data-adaptive test that is optimal for both scenarios when scanning the genome in genome-wide sequencing association studies.

In this paper, by exploiting the relationship between burden-based tests and SKAT, we propose a data-adaptive optimal test within a class of tests that include both burden tests and SKAT as special cases. Specifically, we consider a class of tests that is an arbitrary linear combination of burden test and SKAT statistics, and identify the optimal test within this class to maximize power. We show that this new class of tests can be formulated as a generalized family of SKAT tests by incorporating a correlation structure of variant effects through a family of kernels. It reduces to the burden tests when the effects of variants are perfectly correlated. We derive the optimal test (SKAT-O) by estimating the correlation parameter in the kernel matrix to maximize the power, which corresponds to the estimated weight in the linear combination of the burden test and SKAT test statistics that maximizes power. We derive the theoretical distribution of the SKAT-O test statistic, which allows us to calculate the  $p$ -value analytically with high accuracy in the tail. This is advantageous for analyzing genome-wide sequencing data by avoiding computationally intensive methods, such as resampling or permutation, to calculate  $p$ -values, especially in the tail required to reach genome-wide significance.

In addition, we also consider the problem of designing future sequencing association studies. In particular, to design a new sequence association study, it is important to be able to estimate the required sample size to achieve proper statistical power. Although power and sample size calculation can be done via

simulation, this computer-intensive approach is not desirable. Therefore, we derive the analytical formula for the statistical power of SKAT under the newly proposed family of kernels, and by inverting the power function, we can compute the necessary sample size to adequately power future studies.

## 2. OPTIMAL TEST

For simplicity, we assume that we are interested in testing whether the rare variants in a single region are associated with a complex trait. In a large-scale study of multiple regions, the same methods can be applied with the appropriate adjustment of multiple testing.

### 2.1 Rare variants testing methods

Assume that  $n$  subjects are sequenced in a region with  $p$  genotyped rare variants. For the  $i$ th subject, let  $y_i$  denote a phenotype variable,  $\mathbf{G}_i = (g_{i1}, \dots, g_{ip})$  the genotypes for the  $p$  variants ( $g_{ij} = 0, 1, 2$  for 0, 1, or 2 copies of the minor allele),  $\mathbf{X}_i = (x_{i1}, \dots, x_{iq})$  the covariates for which we would like to adjust (e.g. demographic or environmental variables). To relate genotypes to continuous/categorical phenotypes, we use the generalized linear model (GLM), such that  $y_i$  independently follows an exponential family distribution with first two moments  $E(y_i) = \mu_i$  and  $\text{Var}(y_i) = \phi v(\mu_i)$ , and a link function

$$g(\mu_i) = \mathbf{X}_i \boldsymbol{\alpha} + \mathbf{G}_i \boldsymbol{\beta}, \quad (2.1)$$

where  $v(\cdot)$  is a variance function.  $\boldsymbol{\alpha}$  and  $\boldsymbol{\beta} = (\beta_1, \dots, \beta_p)'$  are the vectors of regression coefficients for the covariates and rare variants, respectively. Under the GLM, the association between the  $p$  rare variants and the phenotype  $y$  can be tested by evaluating the null hypothesis that  $H_0: \boldsymbol{\beta} = (\beta_1, \dots, \beta_p)' = \mathbf{0}$ . However, the standard  $p$  degree of freedom (d.f.) test may lose power when  $p$  is large. To reduce the d.f., additional assumptions need to be made.

Popular burden-based tests reduce the d.f. by making the assumption that each  $\beta_j$  is a function of the MAFs such that  $\beta_j = w(m_j) = w_j \beta_0$ , where  $m_j$  is the MAF of the  $j$ th variant. Then (2.1) becomes

$$g(\mu_i) = \mathbf{X}_i \boldsymbol{\alpha} + \beta_0 \sum_{j=1}^p w_j g_{ij}, \quad (2.2)$$

and the association between the genetic variants and the phenotype can be tested by conducting a 1 d.f. test with  $H_0: \beta_0 = 0$ . We refer to this test as the weighted counting burden test (WBT). The WBT assumes that all variants are causal with the same direction of association and common  $\beta_0$ . Violation of these assumptions can result in a significant loss of power.

SKAT takes a different approach to reducing the d.f. It assumes that each  $\beta_j$  independently follows an arbitrary distribution with mean zero and variance  $w_j^2 \psi$ , where  $w_j$  is a fixed number that may depend on MAF. Under this assumption, the null hypothesis  $H_0: \boldsymbol{\beta} = \mathbf{0}$  is equivalent to  $H_0: \psi = 0$ , i.e. variance component test in generalized linear mixed models by treating  $\boldsymbol{\beta}$  as random effects. Suppose  $\mathbf{X}$  is the  $n \times q$  covariates matrix,  $\mathbf{G} = [\mathbf{G}_1, \dots, \mathbf{G}_n]'$  is the  $n \times p$  genotype matrix,  $\mathbf{W} = \text{diag}[w_1, \dots, w_p]$  is a  $p \times p$  diagonal matrix of weights and  $\mathbf{K} = \mathbf{G} \mathbf{W} \mathbf{W}' \mathbf{G}'$  is an  $n \times n$  weighted linear kernel matrix. Wu and others (2011) proposed to use a class of flexible weight functions of the MAF using the beta density function as  $w_j = \text{Beta}(\text{MAF}_j, a_1, a_2)$ , where the parameters  $a_1$  and  $a_2$  are pre-specified, and  $\text{MAF}_j$  are estimated using the sample MAF of the  $j$ th variant.

We define the working vector by  $\mathbf{y}^* = \mathbf{X} \boldsymbol{\alpha} + \boldsymbol{\Delta}(\mathbf{y} - \boldsymbol{\mu})$ , where  $\boldsymbol{\Delta} = \text{diag}\{g'(\mu_i)\}$ , and the variance matrix by  $\mathbf{V} = \text{diag}\{\phi v(\mu_i)[g'(\mu_i)]^2\}$ . Their estimates under the null hypothesis are  $\tilde{\mathbf{y}} = \mathbf{X} \hat{\boldsymbol{\alpha}} + \hat{\boldsymbol{\Delta}}(\mathbf{y} - \hat{\boldsymbol{\mu}})$ ,  $\hat{\boldsymbol{\Delta}} = \text{diag}\{g'(\hat{\mu}_i)\}$ , and  $\hat{\mathbf{V}} = \text{diag}\{\hat{\phi} v(\hat{\mu}_i)[g'(\hat{\mu}_i)]^2\}$ , respectively, where  $\hat{\boldsymbol{\alpha}}$  is an  $q \times 1$  vector of estimates

of  $\alpha$ ,  $\hat{\mu}$  is an  $n \times 1$  vector of estimates of  $\mu$  and  $\hat{\phi}$  is an estimate of  $\phi$ . All estimates are obtained under the null hypothesis. Following Zhang and Lin (2003), the score test statistic of the variance component  $\psi$  is

$$Q = (\tilde{y} - \mathbf{X}\hat{\alpha})' \hat{\mathbf{V}}^{-1} \mathbf{K} \hat{\mathbf{V}}^{-1} (\tilde{y} - \mathbf{X}\hat{\alpha}) = (\mathbf{y} - \hat{\mu})' \hat{\Delta} \hat{\mathbf{V}}^{-1} \mathbf{K} \hat{\mathbf{V}}^{-1} \hat{\Delta} (\mathbf{y} - \hat{\mu}). \quad (2.3)$$

When  $g(\cdot)$  is a canonical link function, (2.3) can be simplified to  $Q = (\mathbf{y} - \hat{\mu})' \mathbf{K} (\mathbf{y} - \hat{\mu}) / \hat{\phi}^2$ . For binary and Poisson data,  $\phi = 1$ .

## 2.2 New family of kernels

As shown in the previous section, the weighted linear kernel is constructed under the assumption that  $\beta_j$ s are independent. If a large percentage of variants in the target region are associated with the phenotype with the same direction of effect, burden tests can outperform SKAT because the current kernels used by SKAT do not account for correlation in  $\beta$ . Therefore, we propose a new family of kernels that explicitly incorporates correlation among the variant effects.

We propose to allow  $\beta$  to follow a multivariate distribution with exchangeable correlation structure. Then the correlation matrix of  $\beta$  is  $\mathbf{R}_\rho = (1 - \rho)\mathbf{I} + \rho\mathbf{1}\mathbf{1}'$ . With this correlation structure, the SKAT test statistic is a function of  $\rho$ :

$$Q_\rho = (\mathbf{y} - \hat{\mu})' \hat{\Delta} \hat{\mathbf{V}}^{-1} \mathbf{K}_\rho \hat{\mathbf{V}}^{-1} \hat{\Delta} (\mathbf{y} - \hat{\mu}), \quad (2.4)$$

where  $\mathbf{K}_\rho = \mathbf{GWR}_\rho \mathbf{WG}'$ . When  $\rho = 0$ ,  $\mathbf{K}_\rho$  results in the weighted linear kernel SKAT. When  $\rho = 1$ , SKAT test statistic is

$$Q_\rho = (\mathbf{y} - \hat{\mu})' \hat{\Delta} \hat{\mathbf{V}}^{-1} \mathbf{GW} \mathbf{1}\mathbf{1}' \mathbf{WG}' \hat{\mathbf{V}}^{-1} \hat{\Delta} (\mathbf{y} - \hat{\mu}) = \left[ \sum_{i=1}^n \frac{y_i - \hat{\mu}_i}{\hat{\phi} v(\hat{\mu}_i) g'(\hat{\mu}_i)} \sum_{j=1}^p w_j g_{ij} \right]^2,$$

which is equivalent to the square of the score test statistic of WBT. Thus, both tests can be framed within this new family of kernels. In fact, one can easily show that  $Q_\rho$  is a linear combination of SKAT and burden test, i.e.  $Q_\rho = (1 - \rho)Q_{\text{SKAT}} + \rho Q_{\text{burden}}$ .

For a fixed  $\rho$ ,  $Q_\rho$  follows a mixture of  $\chi^2$  distributions (Liu and others, 2007, 2008; Tzeng and Zhang, 2007). Specifically, if  $(\lambda_1, \dots, \lambda_m)$  are the eigenvalues of  $\hat{\mathbf{V}}^{-1/2} \mathbf{K}_\rho \hat{\mathbf{V}}^{-1/2}$ , then the null distribution of  $Q_\rho$  can be closely approximated by  $\sum \lambda_j \chi_{1,j}^2$ , where  $\{\chi_{1,j}^2\}$  are independent  $\chi_1^2$  random variables. To reduce small sample bias, the restricted maximum likelihood (REML) estimator of the variance component can be used (Zhang and Lin, 2003). Define  $\mathbf{P} = \hat{\mathbf{V}}^{-1} - \hat{\mathbf{V}}^{-1} \tilde{\mathbf{X}} (\tilde{\mathbf{X}}' \hat{\mathbf{V}}^{-1} \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}' \hat{\mathbf{V}}^{-1}$ , where  $\tilde{\mathbf{X}} = [\mathbf{1} \ \mathbf{X}]$ , we use the eigenvalues of  $\mathbf{P}^{1/2} \mathbf{K}_\rho \mathbf{P}^{1/2}$  to obtain the null distribution of  $Q_\rho$ . A  $p$ -value can be obtained by matching moments (Liu and others, 2009) or by inverting the characteristic function (Davies, 1980).

To understand the role of the parameter  $\rho$ , we derive in Supplementary Appendix C (supplementary material available at *Biostatistics* online) the analytic relationship between  $\rho$  and the given regression coefficient  $\beta$  as a function of the proportion of causal variants in a region (i.e. the proportion of the  $\beta$  coefficients that are non-zero) and the proportion of causal variants that are protective (i.e. the proportion of the non-zero  $\beta$  coefficients that are negative), assuming the magnitude of the causal variant effects as a function of the MAFs. We illustrate this relationship and its use in the power calculation section.

## 2.3 Optimal test (SKAT-O)

In practice, we rarely have any information about  $\rho$ . Thus, we need a procedure to select  $\rho$  to maximize power. The resulting optimal test corresponds to a best linear combination of SKAT and burden tests that maximizes power. This is in general a challenging problem because  $\rho$  disappears under the null hypothesis.

Davies (1987) studied this problem and proposed to use a maximum of score test statistic as a test statistic. This approach, however, is not directly applicable here due to the different kurtoses of the  $Q_\rho$ s. As  $\rho$  increase, the null distribution of  $Q_\rho$  has a heavier tail, and adjustment for this is difficult.

We employ an different approach and use the minimum of  $p$ -values as a test statistic rather than the score statistics. Specifically, the test statistic is

$$T = \inf_{0 \leq \rho \leq 1} p_\rho, \quad (2.5)$$

where  $p_\rho$  is the  $p$ -value computed based on  $Q_\rho$ .  $T$  can be obtained by simple grid search across a range of  $\rho$ : set a grid  $0 = \rho_1 < \rho_2 < \dots < \rho_b = 1$ , then the test statistic  $T = \min\{p_{\rho_1}, \dots, p_{\rho_b}\}$ .

**2.3.1 Null distribution of the test statistic.** In order to obtain the null distribution of  $T$  and conduct a hypothesis test, let  $\mathbf{Z} = \hat{\mathbf{V}}^{-1/2} \mathbf{G}\mathbf{W}$  and  $\bar{\mathbf{z}} = (\bar{z}_1, \dots, \bar{z}_n)'$ , where  $\bar{z}_i = \sum_{j=1}^p z_{ij}/p$ .  $\mathbf{M} = \bar{\mathbf{z}}(\bar{\mathbf{z}}'\bar{\mathbf{z}})^{-1}\bar{\mathbf{z}}'$  is a projection matrix onto a space of spanned by  $\bar{\mathbf{z}}$ . We further let

$$\tau(\rho) = p^2 \rho \bar{\mathbf{z}}'\bar{\mathbf{z}} + \frac{1-\rho}{\bar{\mathbf{z}}'\bar{\mathbf{z}}} \sum_{j=1}^p (\bar{\mathbf{z}}'\mathbf{z}_{.j})^2,$$

where  $\mathbf{z}_{.j}$  is the  $j$ th column of  $\mathbf{Z}$ , and  $Q_{\rho_1}, \dots, Q_{\rho_b}$  are the score test statistics computed with different  $\rho_v$ , ( $v = 1, \dots, b$ ). Then we show in Supplementary B (supplementary material available at *Biostatistics* online) that under the null hypothesis  $Q_{\rho_v}$  is asymptotically the same as

$$(1 - \rho_v) \left( \sum_{k=1}^m \lambda_k \eta_k + \zeta \right) + \tau(\rho_v) \eta_0 = (1 - \rho) \kappa + \tau(\rho_v) \eta_0, \quad (2.6)$$

where  $\{\lambda_1, \dots, \lambda_m\}$  are non-zero eigenvalues of  $\mathbf{Z}'(\mathbf{I} - \mathbf{M})\mathbf{Z}$ ,  $\eta_k$  ( $k = 0, \dots, m$ ) are i.i.d  $\chi_1^2$  random variables,  $\kappa = \sum_{k=1}^m \lambda_k \eta_k + \zeta$ , and  $\zeta$  satisfies the following conditions:

$$E(\zeta) = 0, \quad \text{Var}(\zeta) = 4 \text{trace}(\mathbf{Z}'\mathbf{M}\mathbf{Z}\mathbf{Z}'(\mathbf{I} - \mathbf{M})\mathbf{Z}),$$

$$\text{Corr} \left( \sum_{k=1}^m \lambda_k \eta_k, \zeta \right) = 0, \quad \text{and} \quad \text{Corr}(\eta_0, \zeta) = 0.$$

Since the Pearson correlation between  $\kappa$  and  $\eta_0$  is zero, we can approximate  $Q_\rho$  as the mixture of two independent random variables. We can approximate the distribution of  $\kappa$  by using the characteristic function inversion method (Davies, 1980) after adjusting for the extra variance term of  $\zeta$ . Letting  $q_{\min}(\rho_v)$  denote the  $(1 - T)$ th percentile of the distribution of  $Q_{\rho_v}$  for each  $\rho_v$ , the  $p$ -value based on the test statistic  $T$  is

$$\begin{aligned} & 1 - P(Q_{\rho_1} < q_{\min}(\rho_1), \dots, Q_{\rho_b} < q_{\min}(\rho_b)) \\ & = 1 - E[P(\kappa < \min\{(q_{\min}(\rho_v) - \tau(\rho_v)\eta_0)/(1 - \rho_v)\} | \eta_0)], \end{aligned}$$

which can be obtained by one-dimensional numerical integration, which can be easily calculated. When we compute  $p_{\rho_v}$  and  $q_{\min}(\rho_v)$ , we approximate each marginal distribution of  $Q_{\rho_v}$  by modifying the moment-matching method of Liu and others (2009). In particular, instead of matching the first 3 moments, we match the mean, variance and kurtosis to improve the approximation in the tail area. To adjust for small

sample bias, one can use REML estimates of the variance components, such that  $\mathbf{Z} = \mathbf{P}^{1/2}\mathbf{G}\mathbf{W}$ . The following algorithm provides the detailed description of the proposed method.

Step 1: Set a grid  $0 = \rho_1 < \rho_2 < \dots < \rho_b = 1$ .

Step 2: Compute  $Q_{\rho_1}, \dots, Q_{\rho_b}$ ,  $\mathbf{Z}$ , and  $\mathbf{M}$ . Here  $\mathbf{Z} = \mathbf{P}^{1/2}\mathbf{G}\mathbf{W}$ .

Step 3: Compute  $\lambda_k$ s,  $\tau(\rho_v)$ , and

$$\mu_Q = \sum_{k=1}^m \lambda_k, \quad \sigma_\zeta = 2\sqrt{\text{trace}(\mathbf{Z}'\mathbf{M}\mathbf{Z}\mathbf{Z}'(\mathbf{I} - \mathbf{M})\mathbf{Z})}, \quad \text{and} \quad \sigma_Q = \sqrt{2 \sum_{k=1}^m \lambda_k^2 + \sigma_\zeta^2}.$$

Step 4: Calculate  $p_{\rho_v}$ ,  $T$  and  $q_{\min}(\rho_v)$ s using the modified moment matching approximation.

Step 5: Numerically integrate  $F(\delta(x)|\lambda)f(x|\chi_1^2)$ , where

$$\delta(x) = (\min\{(q_{\min}(\rho_v) - \tau(\rho_v)x)/(1 - \rho_v)\} - \mu_Q) \frac{\sqrt{\sigma_Q^2 - \sigma_\zeta^2}}{\sigma_Q} + \mu_Q,$$

$f(x|\chi_1^2)$  is a density function of  $\chi_1^2$ , and  $F(\delta(x)|\lambda)$  is a distribution function of a mixture of chi-square distribution,  $\sum \lambda_k \chi_k^2$ . The  $p$ -value is found as

$$p\text{-value} = 1 - \int F(\delta(x)|\lambda)f(x|\chi_1^2) dx.$$

### 3. SAMPLE SIZE AND POWER CALCULATIONS FOR DESIGNING SEQUENCING ASSOCIATION STUDIES

Estimating the necessary sample size to adequately power a study is an important part of designing new sequencing association studies. In this section, we derive the analytical formula for the statistical power of SKAT. We restrict our interest to continuous traits study designs and dichotomous traits study designs, such as case-control study designs, since both are commonly used designs in association studies. We further note that we consider only SKAT with fixed  $\rho$ , since at the design stage, researchers generally specify anticipated alternative hypotheses, i.e. a specific  $\rho$  based on the scenario they have in mind. The detailed power calculation method can be found in Supplementary Appendix A (supplementary material available at *Biostatistics* online). The required sample size to achieve a fixed power level can be easily computed by inverting the power function.

The proposed formula can be used to calculate statistical power efficiently for specified sample sizes and  $\alpha$  level, given prior information on the genetic architecture of the genomic region of interest. Also required are the proportion of causal variants within the region (i.e. the fraction of non-zero  $\beta$  coefficients), and the proportion of causal variants that are protective (i.e. the fraction of non-zero  $\beta$  coefficients that are negative). Both of these can be easily posited by investigators based on prior belief, and detailed discussions are provided in [Wu and others \(2011\)](#). What is more challenging to specify is the particular  $\rho$  value to be used in power calculations. As discussed in Section 2.2 and Supplementary Appendix C (supplementary material available at *Biostatistics* online), if the proportion of nonzero  $\beta$ s is  $p_1$ , the proportion of positive  $\beta$ s among the non-zero  $\beta$ s is  $p_2$ , and the magnitude of the causal variant effects is a function of the MAFs, the optimal  $\rho$  can be estimated as

$$\rho = p_1^2(2p_2 - 1)^2. \quad (3.1)$$

Note that the power computed with this theoretical estimate of  $\rho$  will differ slightly from the power of the optimal test (2.5) which is based on the data driven optimal  $\rho$ , since we are fixing  $\rho$  to obtain the power. However, our simulation studies (Section 4.3) suggest that this difference tends to be small.

4. SIMULATIONS AND REAL DATA ANALYSIS

4.1 Type I error rate

For all simulations, sequence genotypes were generated from 10 000 chromosomes over 1 Mb regions by the calibrated coalescent model with mimicking the linkage disequilibrium structure of European ancestry samples (Schaffner and others, 2005). Continuous phenotypes were generated from the null linear model

$$y_i = 0.5X_1 + 0.5X_2 + \epsilon_i, \quad \epsilon_i \sim N(0, 1),$$

and binary phenotypes were generated from the null logistic model

$$\text{logit } P(y_i = 1) = \alpha_0 + 0.5X_1 + 0.5X_2,$$

where  $X_1$  was a continuous covariate generated from  $N(0, 1)$ ,  $X_2$  was a binary covariate generated from Bernoulli(0.5) and  $\alpha_0$  was chosen to make penetrance 0.01. For binary trait simulations, we generated retrospective case–control data with half being cases and half being controls. Since the average length of a gene is around 3 kb, we randomly selected 3 kb regions across the 1-Mb chromosome. For each model, we generated a total of 10 000 datasets.

We applied 6 different methods to each of the simulated datasets: proposed test with default  $w_j = \text{beta}(\text{MAF}_j; 1, 25)$  weights (SKAT-O); SKAT-O with flat weights (rSKAT-O); original SKAT ( $\rho = 0$ ) with  $\text{beta}(1, 25)$  weights (SKAT); SKAT ( $\rho = 0$ ) with flat weights (rSKAT); counting-based burden test ( $N$ ); and weighted counting burden test ( $W$ ). For rSKAT-O, rSKAT, and counting-based burden test ( $N$ ), only variants with observed MAF  $< 0.03$  were used. The equal size grid of 11 points (from 0 to 1) were used to obtain test statistics of the SKAT-O and rSKAT-O to search for optimal  $\rho$ . Table 1 shows that all 6 methods well-controlled type I error rates with  $\alpha = 0.05$  and  $\alpha = 0.01$ .

To investigate type I error rates at very stringent genome-wide  $\alpha$  levels, we conducted extensive simulations under a slightly different setting (Supplementary Appendix D, available at *Biostatistics* online). Table 2 shows that SKAT-O can accurately control type I error with moderate  $\alpha$  levels, but produces slightly inflated type I error rates at very small  $\alpha$  levels.

Table 1. Type I error estimates of six different methods to test an association between randomly selected 3 kb regions with continuous and binary traits. Each entry represents type I error rate estimates as the proportion of  $p$ -values smaller than  $\alpha$  under the null hypothesis based on 10 000 simulated datasets

Sample size	Level $\alpha$	SKAT-O	SKAT	rSKAT-O	rSKAT	$N$	$W$
Continuous trait							
2000	0.05	0.051	0.051	0.050	0.048	0.048	0.047
	0.01	0.009	0.009	0.009	0.009	0.009	0.009
5000	0.05	0.051	0.049	0.054	0.050	0.049	0.051
	0.01	0.011	0.010	0.011	0.009	0.009	0.009
Dichotomous trait							
2000	0.05	0.051	0.049	0.049	0.048	0.047	0.047
	0.01	0.011	0.009	0.012	0.010	0.010	0.011
5000	0.05	0.047	0.043	0.046	0.044	0.047	0.048
	0.01	0.010	0.009	0.010	0.008	0.009	0.010



Table 2. Type I error estimates of SKAT-O to test an association between randomly selected 3 kb regions with continuous and binary traits at stringent  $\alpha$  level. The sample size was 2000. Each entry represents type I error rate estimates as the proportion of  $p$ -values smaller than  $\alpha$  under the null hypothesis based on  $10^7$  simulated phenotypes

Level $\alpha$	Continuous traits	Binary traits
$10^{-3}$	$1.1 \times 10^{-3}$	$1.1 \times 10^{-3}$
$10^{-5}$	$1.7 \times 10^{-5}$	$1.2 \times 10^{-5}$
$2.5 \times 10^{-6}$	$4.1 \times 10^{-6}$	$3.3 \times 10^{-6}$

#### 4.2 Power

As with the type I error simulations, we randomly selected 3 kb regions from the broader 1 Mb region, but we then randomly chose causal variants from among the variants with true MAF  $< 0.03$ . The continuous phenotype were simulated from

$$y_i = 0.5X_1 + 0.5X_2 + \beta_1g_1 + \cdots + \beta_sg_s + \epsilon_i, \quad \epsilon_i \sim N(0, 1)$$

and dichotomous phenotypes were simulated from

$$\text{logit } P(y_i = 1) = \alpha_0 + 0.5X_1 + 0.5X_2 + \beta_1g_1 + \cdots + \beta_sg_s,$$

where  $(g_1, \dots, g_s)$  were selected causal variants. Covariates  $X_1$  and  $X_2$  follow the same distribution in the type I error simulation, and  $\alpha_0$  was chosen to make prevalence 0.01.

We considered simulations in which 10%, 20%, or 50% of rare variants were causal. Since it is assumed that rarer variants are more likely to have large effect sizes, we set  $\beta_j = c|\log_{10}(m_j)|$ , where  $m_j$  is the MAF of the  $j$ th variant. For continuous trait simulations, we set  $c = 0.6$ , when 10% of the rare variants were causal, which gives maximum  $\beta = 2.4$  for variants with MAF  $= 10^{-4}$ . We used  $c = 0.3$  and  $c = 0.2$ , when 20% and 50% of the rare variants were causal to compensate for the increased number of causal variants. For dichotomous trait simulations, we set  $c = \ln \frac{13}{4} = 0.64$ , when 10% of the rare variants were causal, which gives maximum OR = 13 for variants with MAF  $= 10^{-4}$ . We scaled down  $c$  with larger percentages of causal variants. We allowed the sample size to vary as  $n = 1000, 2000$ , and 5000. Datasets were generated 1000 times for each configuration. We applied the same 6 methods used in the type I error simulations to each dataset, and power was estimated as the proportion of  $p$ -values less than  $\alpha = 2.5 \times 10^{-6}$ .

Figure 1 shows the empirical power under all considered configurations when non-zero  $\beta$  coefficients are all positive, i.e. causal variants are all in the same directions. When the percentage of causal SNPs was low, both the original SKATs (SKAT and rSKAT) and the proposed SKAT-O (SKAT-O and rSKAT-O) had higher power than the burden tests. When the proportion of causal SNPs increased, the burden tests performed better, and the original SKAT has lower power than burden tests when 50% of variants are causal. The SKAT-O and rSKAT-O perform better (when 20% of rare variants are causal) or similar to the burden tests (when 50% of rare variants are causal). This suggests that the performance of SKAT-O is closer to the burden tests in the presence of a larger proportion of causal variants. The higher power of the weighted test over the unweighted test also suggests that appropriate weighting can increase the power.

We also conducted simulations in which 20% of causal variants have negative  $\beta$ s (and 80% have positive  $\beta$ s). Results for these simulations are presented in Figure 2 and show that as expected, burden tests lose a significant amount of power since the effects of the causal variants cancel out due to the presence of negative  $\beta$  coefficients. Both the original and the new optimal SKAT and rSKAT outperform the burden tests no matter where the percentage of causal variants is small or large. In this case, the performance



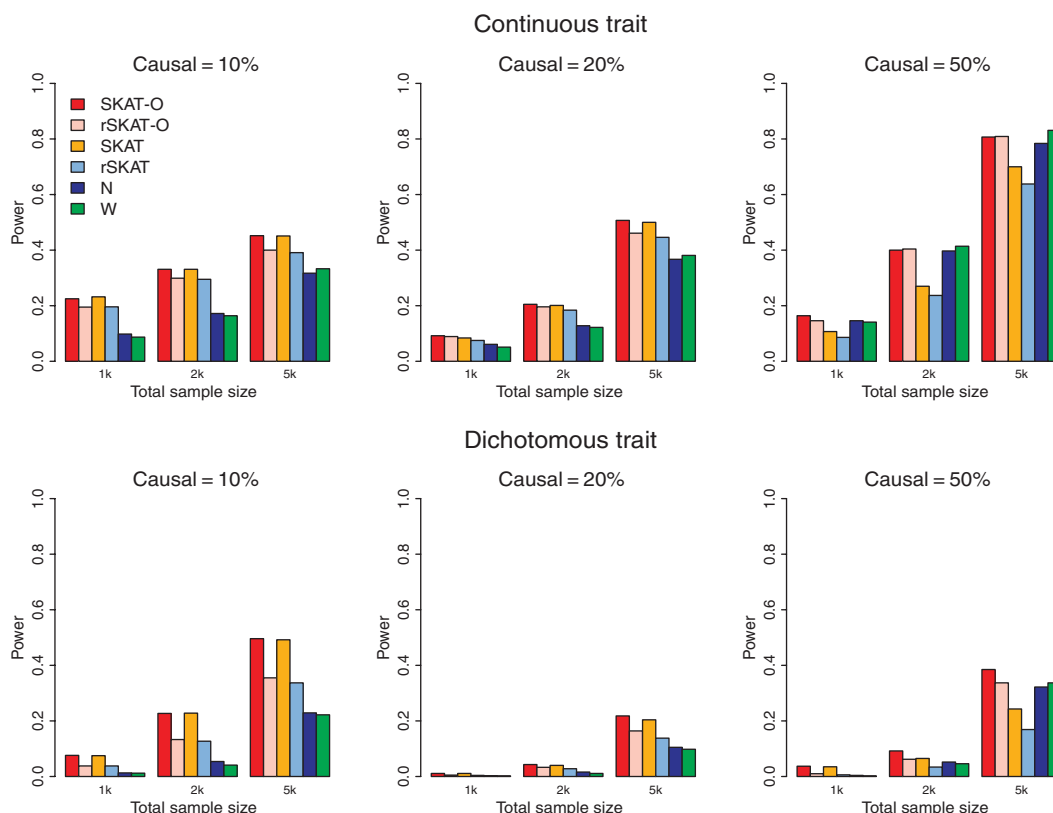


Fig. 1. Empirical power of SKAT-O and competing methods at  $\alpha = 2.5 \times 10^{-6}$  using simulation studies when region size = 3 kb and  $\beta_{\pm} = 100/0$ . Top panel considers continuous phenotypes and bottom panel considers dichotomous phenotypes. From left to right, the plots consider the setting in which 10% of rare variants were causal, 20% of rare variants were causal, and 50% of rare variants were causal. The detailed simulation setups are described in “Simulations and Real Data Analysis”.

of SKAT-O is closer to SKAT. When 10% or 20% of variants were causal variants, SKAT-O and rSKAT-O had slightly lower power than SKAT and rSKAT. With 50% of variants being causal, SKAT-O and rSKAT-O had slightly higher power than SKAT and rSKAT. Compared to Figure 1, SKAT and rSKAT did not suffer any power loss in continuous trait simulation and had slightly lower power in dichotomous trait simulation which resulted from lower enrichment for causal variants due to fixed prevalence and the presence of protective variants. However, burden tests lose considerable power due to the fact that causal variants  $\beta$  coefficients are in mixed directions. In addition, we present the  $\rho$  values selected by SKAT-O in Supplementary Figure S4 (supplementary material available at *Biostatistics* online), which shows that SKAT-O generally selects large  $\rho$  values when the percentage of causal variants is high and  $\beta_{\pm} = 100/0$  and selects small  $\rho$ s when either the percentage of causal variants is low or  $\beta_{\pm} = 80/20$ . Overall, our simulation results confirm that the proposed SKAT-O performed very well under broad circumstances.

#### 4.3 Sample size and power calculation

We investigated the effect of different  $\rho$  values on the power under various models (Supplementary Appendix E, available at *Biostatistics* online). Supplementary Figures S2 and S3 (supplementary material

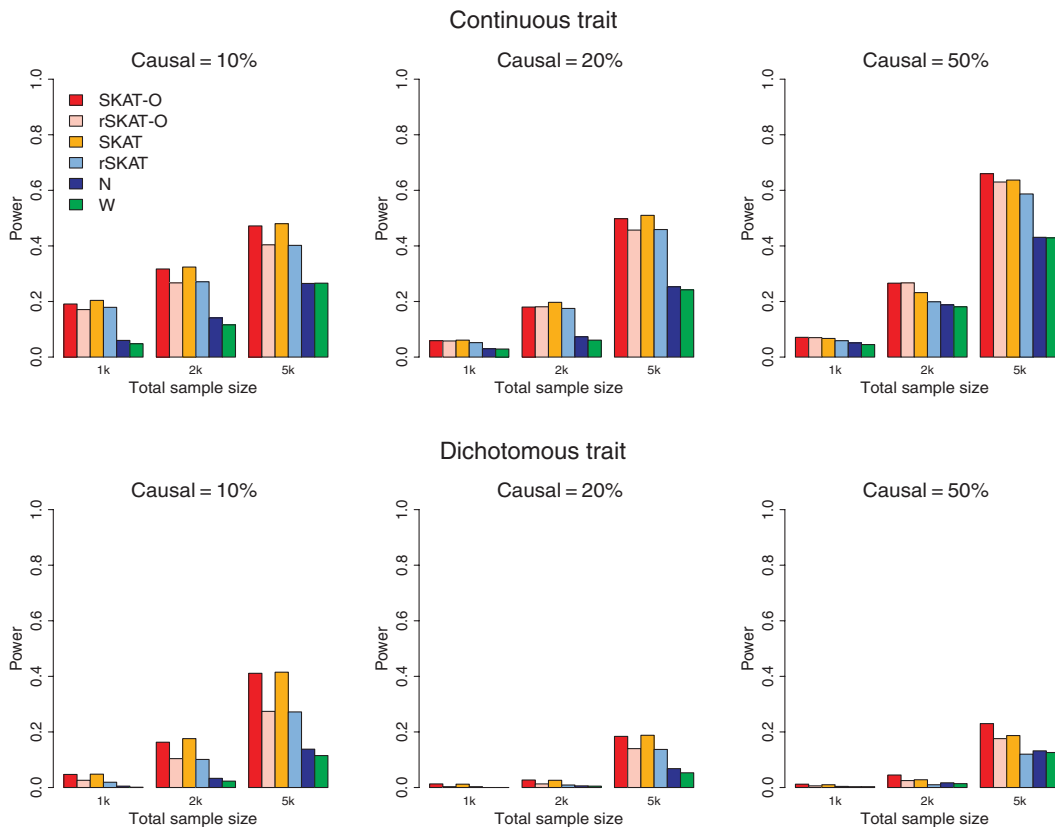


Fig. 2. Empirical power of SKAT-O and competing methods at  $\alpha = 2.5 \times 10^{-6}$  using simulation studies when region size = 3 kb and  $\beta \pm = 80/20$ . Top panel considers continuous phenotypes and bottom panel considers dichotomous phenotypes. From left to right, the plots consider the setting in which 10% of rare variants were causal, 20% of rare variants were causal, and 50% of rare variants were causal. The detailed simulation setups are described in “Simulations and Real Data Analysis”.

available at *Biostatistics* online) show that the test with  $\rho = 0$  (original SKAT) or  $\rho = 1$  (burden tests) can have reduced power, depending on the model assumptions. For example, if only 10% of the variants are causal, then the test with  $\rho = 1$  is significantly less powerful than the test with  $\rho = 0$ . In contrast, when 50% of the variants are causal and all of the non-zero  $\beta$  coefficients are positive, the test with  $\rho = 0$  had lower power than the test with  $\rho = 1$ . In all scenarios,  $\rho$  computed from the proposed formula ( $\rho = \text{estimated}$ ) was most powerful. The figures also show that the theoretical power under the theoretical optimal  $\rho$  closely approximates the empirical power of SKAT-O verifying the adequacy of using equation (3.1) to select  $\rho$  for power and sample size calculations.

We also conducted simulations to evaluate the accuracy of the power calculation formula given  $\rho$  values, and details can be found in Supplementary Appendix F (supplementary material available at *Biostatistics* online).

#### 4.4 Real data application

We applied the proposed SKAT-O and other competing methods to the resequencing data from the Dallas Heart Study (Victor and others, 2004) to test for association between serum triglyceride (TG) levels and

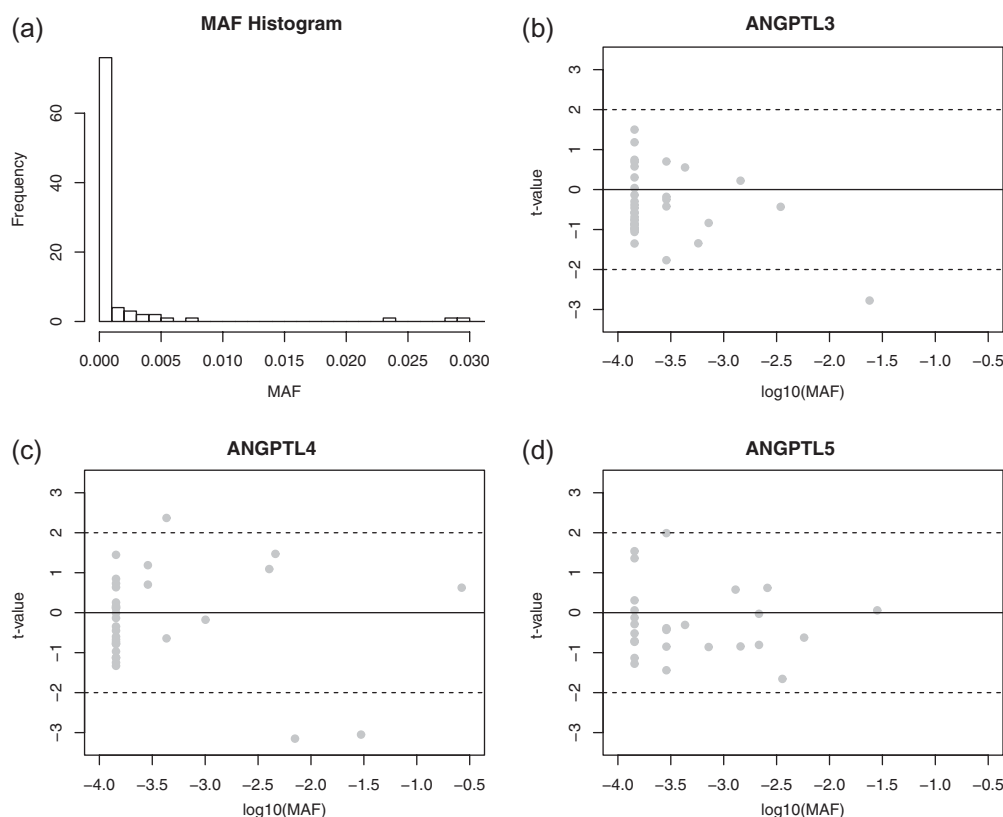


Fig. 3. Single variant analysis results of Dallas Heart Study data. (a) Histogram of minor allele frequencies of 92 variants with  $\text{MAF} < 0.03$ . (b)–(d) Plots of  $\log_{10}(\text{MAF})$  versus  $t$ -statistic values of each variant of ANGPTL3, 4, and 5 genes. The dashed line represents the 95% confidence interval of no-association.

rare variants in 3 genes (ANGPTL3, ANGPTL4, and ANGPTL5). The resequencing dataset has sequence information on 93 observed variants in the 3 genes from each of 3476 individuals in three ethnic groups (white = 1043, black = 1832, and hispanic = 601) (Romeo and others, 2007). A total of 35, 32, and 26 variants reside in the ANGPTL3, ANGPTL4, and ANGPTL5 genes, respectively. All variants except one have  $\text{MAF} < 0.03$ . The histogram of the estimated allele frequencies of the 92 variants with  $\text{MAF} < 0.03$  is presented in Figure 3(a) and clearly indicates that the majority of variants are very rare.

We first performed single-variant association analysis between each variant and log-transformed TG level to explore whether there were variants with different directions of effect. Specifically, we regressed the trait value on the variant while adjusting for gender and ethnicity and computed the  $t$ -statistic based on the regression coefficient (Figures 3(b)–(d)). There is no clear evidence that ANGPTL3 and ANGPTL5 have variants with opposite effects; however, since some of the  $t$ -statistics for the variants in ANGPTL4 have opposing signs and relatively large magnitude, suggesting that there is the potential for variants to have different directions of effect.

We conducted two different analyses. First, we pooled all 93 variants from across the 3 genes and tested for their cumulative effect on log TG level. Second, we considered each gene separately and tested the association between rare variants in each gene and log TG level. We again applied the 6 methods used in the simulation studies with adjusting gender and ethnicity. The results (Table 3) show that when the

Table 3. Analysis of Dallas Heart Study sequence data. Each entry represents a  $p$ -value from each method after adjusting gender and ethnicity. The selected  $\rho$  values by SKAT-O and rSKAT-O are presented in the parentheses. "ALL" indicates analysis results of joint test of all 93 variants in 3 genes

Gene	SKAT-O	rSKAT-O	SKAT	rSKAT	$N$	$W$
Continuous TG level						
ALL	$1.8 \times 10^{-5}$ ( $\rho = 0.1$ )	$4.6 \times 10^{-5}$ ( $\rho = 0.2$ )	$9.5 \times 10^{-5}$	$2.9 \times 10^{-4}$	$7.2 \times 10^{-5}$	$2.3 \times 10^{-4}$
ANGPTL3	$2.6 \times 10^{-3}$ ( $\rho = 0.5$ )	$1.3 \times 10^{-3}$ ( $\rho = 1$ )	$8.9 \times 10^{-3}$	$6.2 \times 10^{-3}$	$1.1 \times 10^{-3}$	$1.9 \times 10^{-3}$
ANGPTL4	$1.7 \times 10^{-4}$ ( $\rho = 0$ )	$8.5 \times 10^{-4}$ ( $\rho = 0.1$ )	$9.5 \times 10^{-5}$	$6.9 \times 10^{-4}$	$4.9 \times 10^{-3}$	$3.1 \times 10^{-2}$
ANGPTL5	$3.4 \times 10^{-1}$ ( $\rho = 1$ )	$5.0 \times 10^{-1}$ ( $\rho = 1$ )	$7.6 \times 10^{-1}$	$8.5 \times 10^{-1}$	$3.3 \times 10^{-1}$	$2.0 \times 10^{-1}$
Dichotomous TG level						
ALL	$1.1 \times 10^{-4}$ ( $\rho = 0.1$ )	$2.8 \times 10^{-4}$ ( $\rho = 0.1$ )	$1.4 \times 10^{-4}$	$2.4 \times 10^{-4}$	$2.2 \times 10^{-3}$	$2.7 \times 10^{-3}$
ANGPTL3	$4.7 \times 10^{-2}$ ( $\rho = 0.4$ )	$3.2 \times 10^{-2}$ ( $\rho = 0.8$ )	$4.7 \times 10^{-2}$	$4.1 \times 10^{-2}$	$2.3 \times 10^{-2}$	$3.6 \times 10^{-2}$
ANGPTL4	$1.6 \times 10^{-4}$ ( $\rho = 0$ )	$3.6 \times 10^{-4}$ ( $\rho = 0.1$ )	$9.4 \times 10^{-5}$	$3.3 \times 10^{-4}$	$3.1 \times 10^{-3}$	$2.2 \times 10^{-2}$
ANGPTL5	$4.6 \times 10^{-1}$ ( $\rho = 0$ )	$3.1 \times 10^{-1}$ ( $\rho = 0$ )	$3.0 \times 10^{-1}$	$1.9 \times 10^{-1}$	$8.7 \times 10^{-1}$	$4.0 \times 10^{-1}$

variants in all three genes were pooled to form a single region for analysis, SKAT-O was by far the most powerful. Our individual variant analysis results suggested that the variants in ANGPTL3 and ANGPTL5 may affect log TG unidirectionally so it is unsurprising that the burden tests ( $W, N$ ) had comparable or better performance than SKAT and rSKAT for testing the variants in ANGPTL3 and ANGPTL5. Also as expected due to the apparent presence of variants with opposing effects, SKAT and rSKAT performed better than the burden tests for testing the association between ANGPTL4 and log TG level. However, SKAT-O and rSKAT-O performed very well across all settings and had only slightly larger  $p$ -values than the best test for the particular setting.

Although log TG is a continuous variable, purely for illustration we also dichotomized the log TG level by taking the the highest and lowest quartiles of each of the six sex-ethnicity groups and using high/low log TG as a dichotomous outcome. The results were qualitatively similar to the results keeping log TG continuous (Table 3).

## 5. DISCUSSION AND CONCLUSION

In this paper, we propose a new family of kernels that incorporates correlation among the effects of causal variants. Based on this new family of kernels, we develop the optimal testing procedure that uses the minimum  $p$ -values from different kernels as a test statistic. In simulation and real data analysis, we show that the proposed optimal test often outperformed the existing burden test and weighted linear kernel SKAT. In addition, we derived sample size/power formula for SKAT for designing new sequence association studies.

In whole exome or whole-genome sequencing studies, we scan the genome by testing a large number of genes. We cannot expect all the genes/regions to follow the same genetic model of association: some are likely to have many causal variants with the effects in the same direction while others may have few causal variants or the causal variants may have the effects in different directions. A good example is the Dallas Heart Study data in which ANGPTL3 and 4 seem to follow different genetic association models. Thus, the proposed SKAT-O can be an attractive choice for many situations, because it adapts to the underlying biological model by selecting  $\rho$  based on the data.

We note that the proposed kernels consider only the compound symmetry covariance structure for the effect of rare variants. The simple structure of compound symmetry allows us to efficiently compute  $p$ -values of the optimal test. For example, it would take only 3–4 h to analyze a whole exome sequencing

study with 20 000 genes and 2000 samples on a laptop. Extension of the method to accommodate other correlation structures is a problem of future interest.

The power/sample size calculation formula presented in this paper is derived with fixed  $\rho$ , and thus we acknowledge that this formula may not be applicable for the proposed SKAT-O. To derive power/sample size formula for the optimal test is challenging due to the  $\zeta$  term, and we suggest a practical approach to pre-select  $\rho$  based on assumptions on underlying genetic structure of association. In particular, researchers can use the simple formula we developed in this paper to select a proper  $\rho$  parameter based on the expected effects of the variants in power calculations. We note that for simplicity, this formula was derived under the assumption that the magnitudes of the non-zero  $\beta$  coefficients are equal to the weights used in the test. However, our simulation studies show that the power is robust to this assumption and is not very sensitive to modest changes in  $\rho$ , and the proposed approach can closely approximate the empirical power of SKAT-O.

As shown in simulation and real data analysis, proper weighting can improve power to detect rare variants association. In this paper, we use the beta(1, 25) density function as a weight function that has been proposed by Wu and others (2011). Since the weight is fixed prior to conducting the association test, the type I error rate is not inflated. It is possible, however, the power can be lost if the weight is misspecified. Although an attractive solution is to choose the weight adaptively using data and obtain  $p$ -values through permutation or resampling, the computational expense is not desirable, particularly for genome-wide studies, and it remains of interest to select the weight that maximizes power.

## 6. SOFTWARE

The software can be found at (<http://www.hsph.harvard.edu/research/skat/>).

## SUPPLEMENTARY MATERIAL

Supplementary material is available at <http://biostatistics.oxfordjournals.org>.

## ACKNOWLEDGMENTS

We thank the Associate Editor and the referees for constructive comments which have improved the paper. *Conflict of Interest:* None declared.

## FUNDING

This work was supported by the National Institutes of Health [R37 CA076404 and P01 CA134294 to S.L. and X.L., and R01 HG006292 to M.C.W.].

## REFERENCES

- BASU, S. AND PAN, W. (2011). Comparison of statistical tests for disease association with rare variants. *Genetic Epidemiology* **35**, 606–619.
- COHEN, J. C., KISS, R. S., PERTSEMLIDIS, A., MARCEL, Y. L., MCPHERSON, R. AND HOBBS, H. H. (2004). Multiple rare alleles contribute to low plasma levels of HDL cholesterol. *Science* **305**, 869.
- DAVIES, R. B. (1980). Algorithm AS 155: the distribution of a linear combination of  $\chi^2$  random variables. *Applied Statistics* **29**, 323–333.
- DAVIES, R. B. (1987). Hypothesis testing when a nuisance parameter is present only under the alternative. *Biometrika* **74**, 33.

- KRYUKOV, G. V., PENNACCHIO, L. A. AND SUNYAEV, S. R. (2007). Most rare missense alleles are deleterious in humans: implications for complex disease and association studies. *The American Journal of Human Genetics* **80**, 727–739.
- LI, B. AND LEAL, S. M. (2008). Methods for detecting associations with rare variants for common diseases: application to analysis of sequence data. *The American Journal of Human Genetics* **83**, 311–321.
- LIU, D., GHOSH, D. AND LIN, X. (2008). Estimation and testing for the effect of a genetic pathway on a disease outcome using logistic kernel machine regression via logistic mixed models. *BMC Bioinformatics* **9**.
- LIU, D., LIN, X. AND GHOSH, D. (2007). Semiparametric regression of multidimensional genetic pathway data: least-squares kernel machines and linear mixed models. *Biometrics* **63**, 1079–1088.
- LIU, H., TANG, Y. AND ZHANG, H.H. (2009). A new chi-square approximation to the distribution of non-negative definite quadratic forms in non-central normal variables. *Computational Statistics & Data Analysis* **53**, 853–856.
- MADSEN, B. E. AND BROWNING, S. R. (2009). A groupwise association test for rare mutations using a weighted sum statistic. *PLoS Genet* **5**, e1000384.
- MORRIS, A. P. AND ZEGGINI, E. (2010). An evaluation of statistical approaches to rare variant analysis in genetic association studies. *Genetic Epidemiology* **34**, 188.
- NEALE, B. M., RIVAS, M. A., VOIGHT, B. F., ALTSHULER, D., DEVLIN, B., ORHO-MELANDER, M., KATHIRESAN, S., PURCELL, S. M., ROEDER, K., DALY, M. J. (2011). Testing for an unusual distribution of rare variants. *PLoS Genetics* **7**, 161–165.
- NEYMAN, J. AND SCOTT, E. (1966). On the use of  $c(\alpha)$  optimal tests of composite hypotheses. *Bulletin of the International Statistical Institute* **41**.
- PRICE, A. L., KRYUKOV, G. V., DE BAKKER, P. I. W., PURCELL, S. M., STAPLES, J., WEI, L. J. AND SUNYAEV, S. R. (2010). Pooled association tests for rare variants in exon-resequencing studies. *The American Journal of Human Genetics* **86**, 832–838.
- ROMEO, S., PENNACCHIO, L. A., FU, Y., BOERWINKLE, E., TYBJAERG-HANSEN, A., HOBBS, H. H. AND COHEN, J. C. (2007). Population-based resequencing of angptl4 uncovers variations that reduce triglycerides and increase hdl. *Nature Genetics* **39**, 513–516.
- SCHAFFNER, S. F., FOO, C., GABRIEL, S., REICH, D., DALY, M. J. AND ALTSHULER, D. (2005). Calibrating a coalescent simulation of human genome sequence variation. *Genome Research* **15**, 1576.
- TZENG, J. Y. AND ZHANG, D. (2007). Haplotype-based association analysis via variance-components score test. *The American Journal of Human Genetics* **81**, 927–938.
- VICTOR, R. G., HALEY, R. W., WILLETT, D. L., PESHOCK, R. M., VAETH, P. C., LEONARD, D., BASIT, M., COOPER, R. S., IANNACCHIONE, V. G., VISSCHER, W. A., STAAB, J. M., HOBBS, H. H., DALLAS HEART STUDY INVESTIGATORS. (2004). The Dallas Heart Study: a population-based probability sample for the multidisciplinary study of ethnic differences in cardiovascular health\* 1. *The American Journal of Cardiology* **93**, 1473–1480.
- WU, M. C., LEE, S., CAI, T., LI, Y., BOEHNKE, M. AND LIN, X. (2011). Rare-variant association testing for sequencing data with the sequence kernel association test. *The American Journal of Human Genetics* **89**, 82–93.
- ZHANG, D. AND LIN, X. (2003). Hypothesis testing in semiparametric additive mixed models. *Biostatistics* **4**, 57.

[Received September 13, 2011; revised January 17, 2012; accepted for publication April 23, 2012]