# Efficient Approximation of P-value of the Maximum of Correlated Tests, with Applications to Genome-Wide Association Studies

Qizhai Li[1,2], Gang Zheng[3], Zhaohai Li[4,5] and Kai Yu[1,*]

[1]*Biostatistics Branch, Division of Cancer Epidemiology and Genetics, National Cancer Institute, Bethesda, MD 20892, USA*

[2]*Academy of Mathematics and Systems Science, Chinese Academy of Sciences, Beijing 100080, P. R. China*

[3]*Office of Biostatistics Research, Office of the Director, National Heart, Lung and Blood Institute, Bethesda, MD 20892, USA*

[4]*Department of Statistics, George Washington University, Washington, D.C., 20052, USA*

[5]*Biometry and Mathematical Statistics Branch, National Institute of Child Health and Human Development, Bethesda, MD 20892, USA*

## Summary

Genome-wide association study (GWAS), typically involving 100,000 to 500,000 single-nucleotide polymorphisms (SNPs), is a powerful approach to identify disease susceptibility loci. In a GWAS, single-marker analysis, which tests one SNP at a time, is usually used as the first stage to screen SNPs across the genome in order to identify a small fraction of promising SNPs with relatively low p-values for further and more focused studies. For single-marker analysis, the trend test derived for an additive genetic model is often used. This may not be robust when the additive assumption is not appropriate for the true underlying disease model. A robust test, MAX, based on the maximum of three trend test statistics derived for recessive, additive, and dominant models, has been proposed recently for GWAS. But its p-value has to be evaluated through a resampling-based procedure, which is computationally challenging for the analysis of GWAS. Obtaining the p-value for MAX with adjustment for the covariates can be even more time-consuming. In this article, we provide a simple approximation for the p-value of the MAX test with or without adjusting for the covariates. The new method avoids resampling steps and thus makes the MAX test readily applicable to GWAS. We use simulation studies as well as real datasets on 17 confirmed disease-associated SNPs to assess the accuracy of the proposed method. We also apply the method to the GWAS of coronary artery disease.

Keywords: MAX, robust genome-wide scan, case-control

## Introduction

Genome-wide association study (GWAS), typically using hundreds of thousands of single nucleotide polymorphisms (SNPs) across the genome, has become a powerful tool for identifying genes or genetic markers underlying disease susceptibility (Klein et al. 2005; Hunter et al. 2007; Sladek et al. 2007; Yeager et al. 2007; The Wellcome Trust Case Control Consortium (WTCCC) 2007). In a typical current GWAS, a panel of 100K–500K SNPs is often genotyped on thousands of individuals. Single-marker analysis, testing the association between the outcome and an individual SNP, is usually used for selecting a subset of SNPs for further investigation (Hoh & Ott 2003; Marchini et al. 2005; Schaid et al. 2005; Wang et al. 2007; Skol et al. 2006; Yu et al. 2007). For example, in a two-stage GWAS (Skol et al. 2006), SNPs whose p-values (obtained from the single-marker analysis in the first stage) are less than a given threshold are evaluated further in an independent sample in the second stage.

A typical test statistic used in single-marker analysis for case-control studies is the Cochran-Armitage trend test (CATT), derived under the assumption of an additive mode of inheritance (Sasieni, 1997; Slager & Schaid, 2001; Zheng et al. 2006a). Since the CATT has an asymptotic normal distribution under the null hypothesis, ranking SNPs based on their test statistics is equivalent to ranking them on their p-values. The CATT for the additive model, however, is

*Corresponding author: Kai Yu, Ph.D., 6120 Executive Boulevard, Room 8050, Rockville, MD 20852. E-mail: yuka@mail.nih.gov

not very robust under other modes of inheritance, e.g., recessive or dominant models. A search for disease-related SNPs with their risk effects governed by a particular disease model might miss SNPs following other risk patterns. Furthermore, for complex diseases with low penetrance, usually none of the above simplified models is appropriate. Under these circumstances, efficiency robust tests, which retain high power across all scientifically plausible genetic models, are preferable (Sladek et al. 2007; Zheng et al. 2003; Zheng et al. 2006a). The theory of efficiency robust tests was summarized in Gastwirth (1985) and Freidlin et al. (1999). One commonly used robust test is based on the MAX statistic, the maximum of three CATTs derived under the recessive, additive, and dominant models, respectively. Empirical results show the advantages of using the MAX statistic over the CATT, derived for the additive model, to prioritize SNPs or to detect disease-associated SNPs (Zheng et al. 2006a).

Under the null hypothesis of no association, the MAX statistic does not follow the standard normal distribution asymptotically. Thus a computationally intensive resampling-based procedure is required to estimate its p-value. For example, in a GWAS of type 2 diabetes, Sladek et al. (2007) conducted 10,000 permutations per SNP to estimate p-values of MAX tests. They identified 59 SNPs, based on a p-value threshold around the level of $10^{-4}$, for further replication in an independent sample. They then used 10,000,000 permutation steps to estimate the p-values associated with the MAX test on each of the 59 chosen SNPs, based on the replication sample. The reason for this extremely large number of permutation steps was to ensure a reliable estimation for any p-value falling below the level of $10^{-6}$. Given situations where the p-value of MAX is not available and a fixed number of SNPs need to be selected for the next-stage study, Zheng et al. (2007) proposed using the MAX statistic rather than its p-value as the basis for the ranking. This approach is easy to carry out without any Monte Carlo simulation. However, the asymptotic null distribution for MAX depends on the genotypic distribution of the study SNP and is SNP-dependent. Therefore, the ranks of SNPs based on their MAX statistics are not weighted on the same scale. It would be more appropriate to rank SNPs based on their p-values.

In many GWAS, in order to account for the other covariates' effects, the logistic regression model is commonly used for the evaluation of individual markers' marginal effect. A similar MAX statistic can be defined based on three Wald (score, or likelihood ratio) test statistics, derived under the dominant, recessive, and additive genotypic effect models, respectively. Clearly, a more computationally intensive resampling procedure is required to estimate the p-value for this type of MAX statistic.

Although using the MAX statistic has various advantages over the CATT derived under an additive model, it is computationally challenging to apply it to a large-scale GWAS. In this article, we propose a simple approach to approximate the p-value of the MAX statistic without Monte Carlo simulation. The approximation formula, called the *Rhombus formula*, is designed to estimate the two-sided test p-value for the MAX statistic. This *Rhombus formula* is an extension of the *W-formula* of Efron (1997), which was originally derived to approximate the one-sided test p-value of the MAX statistic and had been applied to family-based association tests (Yan et al. 2008). To apply this rhombus formula, we need to estimate the covariance matrix for the three CATT (or Wald) tests corresponding to the additive, recessive, and dominant models. Zheng et al. (2006a) provided an analytic formula to estimate the covariate matrix for CATT-based tests. For Wald tests with adjustment for other covariate effects, we propose to use the approach of Pepe et al. (1999), which was based on the generalized estimating equation (GEE) method (Liang & Zeger, 1986), to estimate their covariance matrix numerically. We conducted extensive simulation studies to evaluate the accuracy of the proposed rhombus formula in the setting of the GWAS. To illustrate the application of our methods, we applied the results to 17 confirmed disease-associated SNPs from three GWAS and to a real dataset from a GWAS for coronary artery disease (CAD) with about 350K SNPs (WTCCC, 2007).

## Methods

### A MAX Test Statistic Based on Trend Tests

Suppose that in a GWAS with $r$ cases and $s$ controls, we have genotypes measured on a large panel of SNPs for each individual. Let $AA$, $Aa$, and $aa$ be three possible genotypes for a given SNP. The CATT derived for the additive model is often applied to detect the disease-associated SNPs or to prioritize SNPs for further analysis, if there are no other covariates to be adjusted for. For a given SNP, we denote its genotype frequencies in the case and control groups as shown in Table 1.

Under the notations listed in Table 1, a general form of the CATT can be written as

$$Z_\varphi = \frac{n^{1/2} \sum_{i=0}^{2} \varphi_i (s r_i - r s_i)}{\left\{ r s \left[ n \sum_{i=0}^{2} \varphi_i^2 n_i - \left( \sum_{i=0}^{2} x_i n_i \right)^2 \right] \right\}^{1/2}}, \qquad (1)$$

**Table 1** Notation for Genotype Frequencies

|         | $AA$  | $Aa$  | $aa$  | Total |
|---------|-------|-------|-------|-------|
| Case    | $r_0$ | $r_1$ | $r_2$ | $r$   |
| Control | $s_0$ | $s_1$ | $s_2$ | $s$   |
| Total   | $n_0$ | $n_1$ | $n_2$ | $n$   |

where $\varphi = (\varphi_0, \varphi_1, \varphi_2)$ is a pre-determined genotype score. Since the trend test is invariant under a linear transformation of $\varphi$ with $\varphi_0 \leq \varphi_1 \leq \varphi_2$, we can always set $\varphi_0 = 0$ and $\varphi_2 = 1$, varying the value for $\varphi_1$ between 0 and 1. Under the null hypothesis of no association, $H_0$, $Z_\varphi$ has an asymptotic normal distribution $N(0,1)$. Usually, we do not know which allele has high risk, so a two-sided test is recommended. Results from Sasieni (1997) and Zheng et al. (2006a) showed that the optimal choices of $\varphi_1$ for the recessive, additive, and dominant models are $\varphi_1 = 0$, 1/2, and 1, respectively. Denote the corresponding CATTs by $Z_{\text{REC}}$, $Z_{\text{ADD}}$, and $Z_{\text{DOM}}$.

Based on the three CATTs, $(Z_{\text{REC}}, Z_{\text{ADD}}, Z_{\text{DOM}})$, a robust test (Sladek et al. 2007; Zheng et al. 2006a) is given by

$$Z_{\text{MAX}} = \max\left(|Z_{\text{REC}}|, |Z_{\text{ADD}}|, |Z_{\text{DOM}}|\right).$$

To estimate the p-value of $Z_{\text{MAX}}$ using the rhombus formula described later, we need to evaluate pairwise asymptotic correlation coefficients among $(Z_{\text{REC}}, Z_{\text{ADD}}, Z_{\text{DOM}})$. From Zheng et al. (2006a), under the null hypothesis, we have

$$\text{cor}_{H_0}(Z_{\text{ADD}}, Z_{\text{DOM}}) = \frac{p_0(p_1 + 2p_2)}{\sqrt{p_0(1-p_0)}\sqrt{p_0(p_1 + 2p_2) + p_2(p_1 + 2p_0)}},$$

$$\text{cor}_{H_0}(Z_{\text{REC}}, Z_{\text{DOM}}) = \frac{p_0 p_2}{\sqrt{p_0(1-p_0)}\sqrt{p_2(1-p_2)}}, \text{ and} \qquad (2)$$

$$\text{cor}_{H_0}(Z_{\text{REC}}, Z_{\text{ADD}}) = \frac{p_2(2p_0 + p_1)}{\sqrt{p_2(1-p_2)}\sqrt{p_0(p_1 + 2p_2) + p_2(p_1 + 2p_0)}}$$

where $p_0$, $p_1$ and $p_2$ are probabilities of genotypes $AA$, $Aa$, and $aa$ in the population, respectively, which can be estimated by the observed frequencies using the combined case and control samples, $\hat{p}_i = \frac{n_i}{n}$ for $i = 0, 1, 2$ (see Table 1).

## A MAX Test Statistic Based on Wald Tests

When there are other covariates to be adjusted for, the following logistic regression model can be used,

$$\log \frac{\Pr(\gamma = 1 | \mathbf{z}, g)}{1 - \Pr(\gamma = 1 | \mathbf{z}, g)} = \alpha + \mathbf{z}^T \gamma + g^T \beta, \qquad (3)$$

where $\gamma$, $\mathbf{z}$, and $g$ are, respectively, the outcome variable (case or control), the column vector of non-genetic covariates, and the genotype variable at the study SNP. Similar to the CATT, the genotype in model (3) can be coded according to the following three schemes: (i) $g^{(R)} = 1$ for $AA$, and 0 for either $Aa$ or $aa$, which is based on the recessive model in terms of the odd ratio; (ii) $g^{(A)}$ equals the number of copies (0, 1, or 2) of allele $A$, which is based on the additive (in logit scale) model; and (iii) $g^{(D)} = 1$ for either $AA$ or $Aa$, and 0 for $aa$, which is based on the dominant model. For each type of predictor, say $g^{(R)}$, denote the corresponding coefficients in model (3) by $(\alpha_{(R)}, \gamma_{(R)}, \beta_{(R)})$. The null hypothesis for $g^{(R)}$, for example, can be written as $H_0 : \beta_{(R)} = 0$. The standard likelihood ratio test, score test, and Wald test can be used to test this null hypothesis while adjusting for the effect of $\mathbf{z}$.

Here we focus on the Wald test. Depending on which model is assumed, $g^{(R)}$, $g^{(A)}$, or $g^{(D)}$, we could have three different Wald

tests, denoted by $W_{\text{REC}}$, $W_{\text{ADD}}$, and $W_{\text{DOM}}$, each of which is asymptotically optimal under the assumed model. To have a more robust test when the underlying genetic model is unknown, we define the following MAX statistic based on the Wald test,

$$W_{\text{MAX}} = \max(|W_{\text{REC}}|, |W_{\text{ADD}}|, |W_{\text{DOM}}|).$$

In order to evaluate the p-value of $W_{\text{MAX}}$ using the rhombus formula, which will be described later, we need to estimate pairwise asymptotic correlation coefficients among $(W_{\text{REC}}, W_{\text{ADD}}, W_{\text{DOM}})$. However, unlike the case for the CATT, we do not have explicit formulas for the correlation coefficients. Instead, we propose to use the approach of Pepe et al. (1999) to estimate them numerically.

## Covariance Matrix Estimation

Pepe et al. (1999) originally proposed to use the GEE method to compare several predictors in terms of their strength of association with a common outcome. In our application, we have three predictors, $g^{(R)}$, $g^{(A)}$, and $g^{(D)}$, and one common outcome $\gamma$. The association between $\gamma$ and each of $g^{(R)}$, $g^{(A)}$, and $g^{(D)}$ is measured by $\beta_{(R)}$, $\beta_{(A)}$, and $\beta_{(D)}$, respectively. Pepe et al.'s (1999) procedure provides a way to estimate the covariance matrix for the estimates of $(\beta_{(R)}, \beta_{(A)}, \beta_{(D)})$, and thus to estimate the correlation coefficients among the three Wald test statistics.

Here is an outline of how to apply the procedure of Pepe et al. (1999) Let $\{(\gamma_i, \mathbf{z}_i, g_i) : i = 1, \ldots, n\}$ be observed values for a sample of $n$ subjects. Then form the following coefficient vector by combining coefficients from the three models (recessive, additive, and dominant),

$$\theta = \begin{pmatrix} \alpha_{(R)} & \alpha_{(A)} & \alpha_{(D)} & \gamma_{(R)}^T & \gamma_{(A)}^T & \gamma_{(D)}^T & \beta_{(R)} & \beta_{(A)} & \beta_{(D)} \end{pmatrix}^T.$$

For the $i$th subject, based on its non-genetic covariates $\mathbf{z}_i$ and three predictors $g_i^{(R)}$, $g_i^{(A)}$, and $g_i^{(D)}$, create the following three rows of expanded covariates, corresponding to the coefficient vector $\theta$,

$$\begin{pmatrix} \mathbf{x}_{i,1}^T \\ \mathbf{x}_{i,2}^T \\ \mathbf{x}_{i,3}^T \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & \mathbf{z}_i^T & \mathbf{0}^T & \mathbf{0}^T & g_i^{(R)} & 0 & 0 \\ 0 & 1 & 0 & \mathbf{0}^T & \mathbf{z}_i^T & \mathbf{0}^T & 0 & g_i^{(A)} & 0 \\ 0 & 0 & 1 & \mathbf{0}^T & \mathbf{0}^T & \mathbf{z}_i^T & 0 & 0 & g_i^{(D)} \end{pmatrix}.$$

We can estimate $\theta$ by solving the following estimating equations,

$$\mathbf{R}(\mathbf{x}, \theta) = \sum_{i=1}^n \sum_{k=1}^3 \mathbf{x}_{i,k}(\gamma_i - \mu_{i,k}(\theta)) = \mathbf{0}, \qquad (4)$$

where $\mu_{i,k}(\theta) = \frac{\exp(\mathbf{x}_{i,k}^T \theta)}{1 + \exp(\mathbf{x}_{i,k}^T \theta)}$. Let $\hat{\theta} = (\hat{\alpha}_{(R)} \quad \hat{\alpha}_{(A)} \quad \hat{\alpha}_{(D)} \quad \hat{\gamma}_{(R)}^T \quad \hat{\gamma}_{(A)}^T \quad \hat{\gamma}_{(D)}^T \quad \hat{\beta}_{(R)} \quad \hat{\beta}_{(A)} \quad \hat{\beta}_{(D)})^T$ be the estimate based on the above estimation equation. Its covariance matrix $V(\hat{\theta})$ can be estimated by the following sandwich estimate,

$$V(\hat{\theta}) = (M^{-1}(\hat{\theta}))^T \sum(\hat{\theta}) M^{-1}(\hat{\theta}),$$

where $M(\hat{\theta})$ is the information matrix and is defined as

$$M(\hat{\theta}) = \sum_{i=1}^{n} \sum_{k=1}^{3} \mathbf{x}_{i,k} \mathbf{x}_{i,k}^T \mu_{i,k}(\hat{\theta})(1 - \mu_{i,k}(\hat{\theta})),$$

and

$$\sum(\hat{\theta}) = \sum_{i=1}^{n} \left[ \sum_{k=1}^{3} \mathbf{x}_{i,k}(\gamma_i - \mu_{i,k}(\hat{\theta})) \right] \left[ \sum_{k=1}^{3} \mathbf{x}_{i,k}(\gamma_i - \mu_{i,k}(\hat{\theta})) \right]^T.$$

We can then define the three Wald tests as $W_{\text{REC}} = \hat{\beta}_{(R)}/\sqrt{v_{(R)}}$, $W_{\text{ADD}} = \hat{\beta}_{(A)}/\sqrt{v_{(A)}}$, and $W_{\text{DOM}} = \hat{\beta}_{(D)}/\sqrt{v_{(D)}}$, with $v_{(R)}$, $v_{(A)}$, and $v_{(D)}$ being the 7th to 9th diagonal elements in the covariance matrix $V(\hat{\theta})$. Each Wald test has an asymptotic $N(0,1)$ under the null hypothesis. Their correlation matrix can be obtained (after rescaling) from the corresponding principal submatrix of $V(\hat{\theta})$.

It can be verified that $(\hat{\alpha}_{(R)}, \hat{\gamma}_{(R)}^T, \hat{\beta}_{(R)})$ is the same as the maximum likelihood estimate (MLE) based on the logistic regression model given by (3) with genotype coded by $g^{(R)}$. The same is true for estimates under the other two genetic models. Thus by solving the estimation equation (4), we simultaneously obtain the MLE for $(\alpha, \gamma^T, \beta)$ under three different models based on (3). However, the estimated variances $(v_{(R)}, v_{(A)}, v_{(D)})$ for $(\hat{\beta}_{(R)}, \hat{\beta}_{(A)}, \hat{\beta}_{(D)})$ are based on the robust sandwich estimate and are different from the ones based on the information matrix derived from (3). Thus our definition for the Wald statistic is slightly different from the standard Wald statistic, although the two definitions are asymptotically equivalent.

## A Rhombus Formula to Approximate the P-value of MAX

Once we obtain the estimates of pairwise correlation coefficients among $(Z_{\text{ADD}}, Z_{\text{DOM}}, Z_{\text{REC}})$ or $(W_{\text{REC}}, W_{\text{ADD}}, W_{\text{DOM}})$, we can use the approximation method developed in this section to calculate the p-value of $Z_{\text{MAX}}$ or $W_{\text{MAX}}$. We describe the method in its general form. Assume that there are $k$ ($k = 3$ for our application) test statistics, $T_1, T_2, \ldots, T_k$, each of which is used to test the null hypothesis $H_0$, under which these $k$ test statistics approximately follow the standard normal distribution whose density and probability functions are denoted by $\phi(x)$ and $\Phi(t)$, respectively. We further assume that the correlation coefficient cor $(T_i, T_j)$, for $i, j \in \{1, 2, \ldots, k\}$, is known or can be estimated consistently. Let $T_{\max} = \max\{|T_1|, |T_2|, \ldots, |T_k|\}$ be the MAX statistic. Given $T_{\max} = t$, we are interested in calculating the p-value $\Pr(T_{\max} > t)$.

Letting $T^*_{\max} = \max\{T_1, T_2, \ldots, T_k\}$, Efron (1997) derived a formula (called the $W$-formula) to approximate $\Pr(T^*_{\max} > t^*)$, with $t^*$ being an observed value for $T^*_{\max}$. Thus Efron (1997) dealt with a one-sided rejection region, whereas we try to calculate the probability for a rejection region that is symmetric about the origin point. Following the techniques of Efron (1997), we derived a tight upper bound for $\Pr(T_{\max} > t)$. The derivation is given in the Appendix. Corresponding to Efron's $W$-formula, we call ours the rhombus formula, which is given by

$$\Pr(T_{\max} > t) \leq (k - 2)[\Phi(t) - \Phi(-t) - 1] - \frac{4\phi(t)(k - 1)}{t}$$

$$+ \frac{4\phi(t)}{t} \sum_{i=1}^{k-1} \left\{ 2\Phi\left(\frac{t L_{i(i+1)}}{2}\right) + e^{-\frac{t^2 L_{i(i+1)}^2}{8}} \left[\Phi\left(\frac{t(\pi - L_{i(i+1)})}{2}\right) - \Phi\left(\frac{t L_{i(i+1)}}{2}\right)\right] \right\} I\left\{0 \leq L_{i(i+1)} \leq \frac{\pi}{2}\right\}$$

$$+ \frac{4\phi(t)}{t} \sum_{i=1}^{k-1} \left\{ 2\Phi\left(\frac{t(\pi - L_{i(i+1)})}{2}\right) + e^{-\frac{t^2[\pi - L_{i(i+1)}]^2}{8}} \left[\Phi\left(\frac{t L_{i(i+1)}}{2}\right) - \Phi\left(\frac{t(\pi - L_{i(i+1)})}{2}\right)\right] \right\} I\left\{\frac{\pi}{2} \leq L_{i(i+1)} \leq \pi\right\}. \quad (5)$$

where $L_{ij} = \arccos(\text{cor}(T_i, T_j))$ and $I\{\cdot\}$ is an indicator function. We can see from (5) that the estimated upper bound for the p-value depends on how these $k$ test statistics are indexed, but $\Pr(T_{\max} > t)$ is independent of the index. To have a tighter upper bound, in practice, we can compare the upper bound evaluated under all possible orderings of the $k$ test statistics, choosing the smallest value as the estimation for the p-value. This strategy is feasible for $Z_{\text{MAX}}$ and $W_{\text{MAX}}$ with $k = 3$.

It should be pointed out that the rhombus formula provides a theoretical upper bound for the true p-value if $(T_1, T_2, \ldots, T_k)$ follows a joint normal distribution with a known correlation matrix. In real applications, $(T_1, T_2, \ldots, T_k)$ is asymptotically normal. The correlation matrix is also estimated. Therefore the true p-value of the MAX test is not necessarily less than the bound calculated by the rhombus formula. However, from the numerical examples in both simulations and real data applications, we observed that the values given by the rhombus formula tended to overestimate the true p-values.

## Simulations Design

To evaluate the accuracy of the rhombus formula for approximating the p-value associated with $Z_{\text{MAX}}$ and $W_{\text{MAX}}$, we conducted simulation studies to estimate the empirical type I error rate under various significance levels. We simulated genotypes for 1,000,000 null SNPs for cases and controls, with various sample

sizes. We considered two scenarios: S1, all cases and controls were sampled from a homogeneous population; S2, the study population consisted of two subpopulations. For each scenario, Hardy-Weinberg equilibrium (HWE) was assumed within each subpopulation. Under S1, the MAX test $Z_{MAX}$ was applied, with its p-value estimated by the rhombus formula (5). Under S2, we used $W_{MAX}$ with an adjustment for the (known) subpopulation structure, i.e., we entered the covariate $\mathbf{z}$ in model (3), with $\mathbf{z} = 0$ for subjects from subpopulation 1, and $\mathbf{z} = 1$ for subjects from subpopulation 2. Note that the purpose of S2 is not to evaluate the effect of the population substructure, but to demonstrate the use of the $W_{MAX}$ test. We assumed the number of cases ($r$) was the same as that of controls ($s$), with $r = s = 500, 1,000, 1,500,$ and 2,000. Under S1, we assumed that minor allele frequencies (MAFs) of all 1,000,000 SNPs were independently generated from the uniform distribution U[0.1, 0.5], and we randomly assigned genotypes to cases and controls according to the genotype frequencies under HWE. Under S2, for any given SNP, its MAFs in these two subpopulations were generated by two independent random draws from a Beta distribution with two parameters, $p(1 - F_{ST})/F_{ST}$ and $(1 - p)(1 - F_{ST})/F_{ST}$, where $F_{ST} = 0.01$ (a typical value for divergent European populations), and $p$ was the ancestral population MAF drawn from U[0.1, 0.5] (Price et al. 2006).We further assumed that 60% and 40% of cases were chosen from subpopulations 1 and 2, respectively, while 40% and 60% of controls were sampled from subpopulations 1 and 2, respectively. Under each scenario, based on the p-value estimated by the rhombus formula, we estimated the empirical Type I error for MAX by averaging results over 1,000,000 null SNPs. The nominal level $\alpha$ was set to 0.0001, 0.001, 0.01, 0.05, and 0.1.

### Application to 17 Disease–Associated SNPs

We applied the MAX test to 17 SNPs whose association with various complex diseases had been confirmed, including 8 SNPs associated with type 2 diabetes (Sladek et al. 2007), 6 SNPs associated with breast cancer (Hunter et al. 2007), and 3 SNPs associated with prostate cancer (Yeager et al. 2007). For each of the above SNPs, we obtained its genotype counts and applied the MAX test $Z_{MAX}$. We used the rhombus formula as well as the two resampling-based approaches (the parametric bootstrap and

the permutation procedure) to estimate the p-values. The bootstrap method generates genotype counts for the cases and controls under the null hypothesis of no association, based on genotype frequencies in the pooled samples. The permutation procedure just randomly shuffles the case/control status among all individuals. We used 10,000,000 bootstrap or permutation steps to ensure a reliable estimation of the p-value for each SNP.

### Application to the GWAS of Coronary Artery Disease (CAD)

CAD is one of the most common heart diseases. It is the main cause of death among the elderly. WTCCC (2007) reported results of a GWAS with 459,446 SNPs for CAD. To demonstrate the application of the MAX test in GWAS, we applied the test based on $Z_{MAX}$ to the CAD study and estimated p-values using the rhombus formula. We focused on 343,413 SNPs, after removing SNPs without SNP IDs, SNPs with a genotype frequency below 5 in any cell listed in Table 1 and SNPs with bad clustering properties.

## Results

### Simulation Results

Table 2 reports the empirical Type I error results when the p-value of the MAX statistic is approximated by the rhombus formula. It shows that the rhombus formula can estimate the p-value reasonably well. Similar conclusions can be made when we draw MAF from more restrictive intervals rather than uniformly from [0.1, 0.5] (results not shown). It is not surprising that the rhombus formula tends to overestimate the p-value. As a result, the empirical Type I error based on the estimated p-value is lower than the nominal value most of the time. From Table 2, it appears that the rhombus formula is especially appropriate when true p-values are relatively small (less than 0.1). For example, when the nominal level is 0.05, the largest absolute difference between our estimated p-values and 0.05 is less than 0.004 (with N = 500), while the largest absolute difference

**Table 2** Empirical Type I Error Based on 1,000,000 Replicates

| Using the MAX test $Z_{MAX}$ with p-values estimated by the rhombus formula | | | | | | |
|---|---|---|---|---|---|---|
| N | 0.0001 | 0.001 | 0.01 | 0.05 | 0.1 | 0.2 |
| 500 | $1.13 \times 10^{-4}$ | $9.17 \times 10^{-4}$ | 0.0093 | 0.0463 | 0.0915 | 0.1768 |
| 1,000 | $9.90 \times 10^{-5}$ | $9.46 \times 10^{-4}$ | 0.0096 | 0.0476 | 0.0934 | 0.1791 |
| 1,500 | $1.07 \times 10^{-4}$ | $1.03 \times 10^{-3}$ | 0.0099 | 0.0482 | 0.0936 | 0.1789 |
| 2,000 | $9.70 \times 10^{-5}$ | $1.02 \times 10^{-3}$ | 0.0100 | 0.0484 | 0.0939 | 0.1797 |
| Using the MAX test $W_{MAX}$ with p-values estimated by the rhombus formula | | | | | | |
| 500 | $9.00 \times 10^{-5}$ | $9.26 \times 10^{-4}$ | 0.0093 | 0.0475 | 0.0928 | 0.1788 |
| 1,000 | $9.50 \times 10^{-5}$ | $9.25 \times 10^{-4}$ | 0.0098 | 0.0478 | 0.0935 | 0.1795 |
| 1,500 | $9.40 \times 10^{-5}$ | $9.59 \times 10^{-4}$ | 0.0097 | 0.0476 | 0.0931 | 0.1793 |
| 2,000 | $9.30 \times 10^{-5}$ | $9.57 \times 10^{-4}$ | 0.0098 | 0.0477 | 0.0934 | 0.1797 |

**Table 3** P-values of Indentiied SNPs in GWASs of Diabetes, Breast, and Prostate Cancers

| | $r_0$ | $r_1$ | $r_2$ | $s_0$ | $s_1$ | $s_2$ | Rhombus Formula | 10,000,000 Bootstraps | 10,000,000 Permutations |
|---|---|---|---|---|---|---|---|---|---|
| **8 confirmed SNPs associated with Type 2 Diabetes[3]** | | | | | | | | | |
| rs7903146 | 197 | 348 | 149 | 335 | 254 | 65 | $1.58 \times 10^{-18}$ | $< 1 \times 10^{-7}$ | $< 1 \times 10^{-7}$ |
| rs13266634 | 54 | 229 | 411 | 53 | 293 | 307 | $1.84 \times 10^{-5}$ | $1.52 \times 10^{-5}$ | $1.42 \times 10^{-5}$ |
| rs1111875 | 77 | 302 | 315 | 119 | 308 | 227 | $6.78 \times 10^{-6}$ | $5.40 \times 10^{-6}$ | $7.10 \times 10^{-6}$ |
| rs7923837 | 66 | 300 | 328 | 116 | 296 | 242 | $2.28 \times 10^{-6}$ | $2.20 \times 10^{-6}$ | $2.33 \times 10^{-6}$ |
| rs7480010 | 301 | 327 | 66 | 363 | 246 | 45 | $2.18 \times 10^{-5}$ | $1.76 \times 10^{-5}$ | $2.19 \times 10^{-5}$ |
| rs3740878 | 25 | 273 | 386 | 65 | 249 | 353 | $1.84 \times 10^{-5}$ | $1.70 \times 10^{-5}$ | $1.52 \times 10^{-5}$ |
| rs11037909 | 25 | 274 | 387 | 65 | 251 | 353 | $1.85 \times 10^{-5}$ | $1.81 \times 10^{-5}$ | $1.71 \times 10^{-5}$ |
| rs1113132 | 25 | 271 | 390 | 63 | 251 | 355 | $4.12 \times 10^{-5}$ | $3.68 \times 10^{-5}$ | $3.66 \times 10^{-5}$ |
| **6 reported SNPs associated with breast cancer[2]** | | | | | | | | | |
| rs10510126 | 955 | 180 | 10 | 854 | 272 | 14 | $1.42 \times 10^{-6}$ | $0.90 \times 10^{-6}$ | $0.50 \times 10^{-6}$ |
| rs12505080 | 608 | 477 | 50 | 628 | 408 | 99 | $8.27 \times 10^{-5}$ | $7.92 \times 10^{-5}$ | $7.33 \times 10^{-5}$ |
| rs17157903 | 777 | 316 | 18 | 862 | 220 | 26 | $6.20 \times 10^{-5}$ | $4.95 \times 10^{-5}$ | $5.69 \times 10^{-5}$ |
| rs1219648 | 352 | 543 | 250 | 433 | 538 | 170 | $4.80 \times 10^{-6}$ | $4.30 \times 10^{-6}$ | $4.10 \times 10^{-6}$ |
| rs7696175 | 353 | 605 | 187 | 396 | 496 | 249 | $1.98 \times 10^{-3}$ | $2.07 \times 10^{-3}$ | $2.14 \times 10^{-3}$ |
| Rs2420946 | 357 | 546 | 242 | 440 | 537 | 165 | $5.14 \times 10^{-6}$ | $5.60 \times 10^{-6}$ | $3.80 \times 10^{-6}$ |
| **3 reported SNPs associated with prostate cancer[4]** | | | | | | | | | |
| rs1447295 | 25 | 283 | 864 | 10 | 218 | 929 | $1.10 \times 10^{-4}$ | $0.88 \times 10^{-4}$ | $0.80 \times 10^{-4}$ |
| rs6983267 | 351 | 598 | 223 | 277 | 579 | 301 | $2.06 \times 10^{-5}$ | $2.12 \times 10^{-5}$ | $2.36 \times 10^{-5}$ |
| rs7837688 | 861 | 283 | 27 | 939 | 206 | 11 | $6.67 \times 10^{-6}$ | $3.70 \times 10^{-6}$ | $3.00 \times 10^{-6}$ |

becomes 0.023 (with N = 500) when the nominal level is 0.2. Thus, the rhombus formula becomes modestly conservative for approximating less extreme p-values. This demonstrates that the rhombus formula is not only particularly useful for GWAS, where the main focus is on the SNPs with small p-values; it can also be applied to candidate studies with a nominal level below 0.05.

## Estimated P-Values for 17 Disease–Associated SNPs

Table 3 reports the p-values obtained by the rhombus formula and two resampling-based methods for 17 confirmed SNPs from three genetic studies of type 2 diabetes (Sladek et al. 2007), breast cancer (Hunter et al. 2007), and prostate cancer (Yeager et al. 2007). Table 3 shows that p-values from the three methods generally agree well, especially when the minimum genotype count observed in the cases and controls is larger than 20. The p-value estimated by the rhombus formula tends to be slightly larger than that obtained by the resampling-based procedures. This is consistent with the fact that the rhombus formula provides a theoretical upper bound under the normality assumption. One advantage of using the rhombus formula to estimate the p-value is that it can provide a reasonably accurate approximation when the true p-value is less than $10^{-6}$, which requires more than 10,000,000 permutation steps.
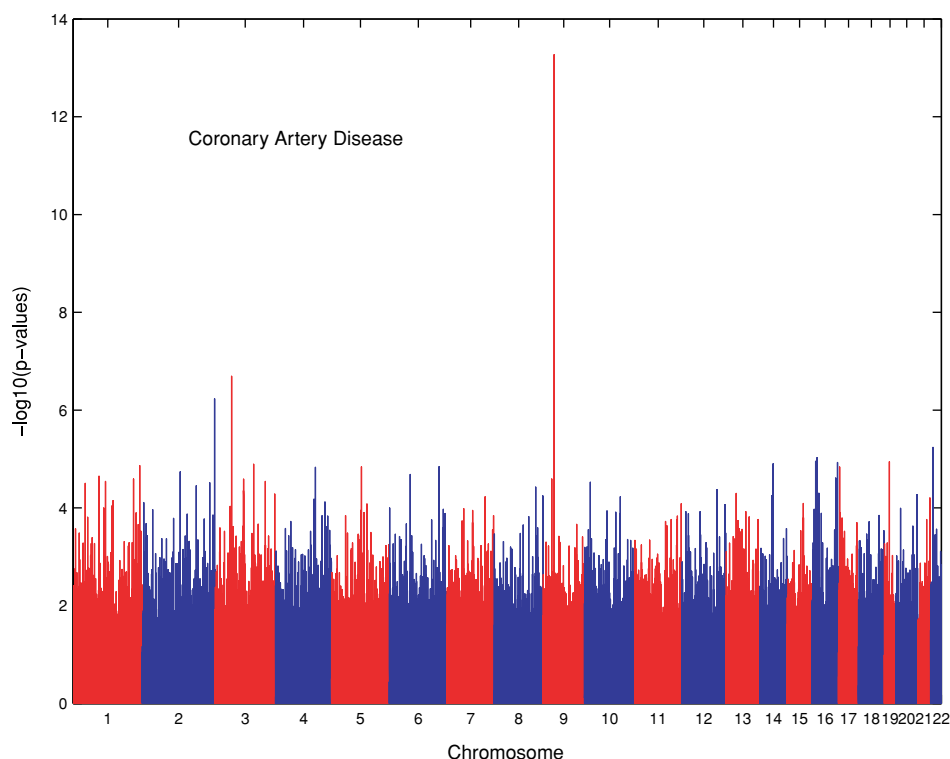
For example, for SNP rs7903146 in the study of type 2 diabetes (Sladek et al. 2007), the resampling-based estimates show that the p-value is less than $10^{-7}$ because no simulated MAX statistics were more extreme than the observed MAX. Using the rhombus formula, we estimated the p-value to be $1.58 \times 10^{-18}$. This information is useful in replication studies and meta-analyses when p-values from several studies are to be combined.

## Application to GWAS of CAD

Figure 1 plots the estimated p-values based on the rhombus formula for 343,413 SNPs according to their positions along each chromosome. Table 4 lists all 22 SNPs with an estimated p-value (by the rhombus formula) below $10^{-5}$. Also presented in Table 4 are p-values estimated by the two resampling-based methods. From Table 4, it can be seen that results from the rhombus formula agree well with estimates from the two resampling-based methods. For those SNPs with extremely small pvalues (say less than $10^{-7}$), it is not computationally feasible to use resampling-based methods. The rhombus formula has no such limitations.
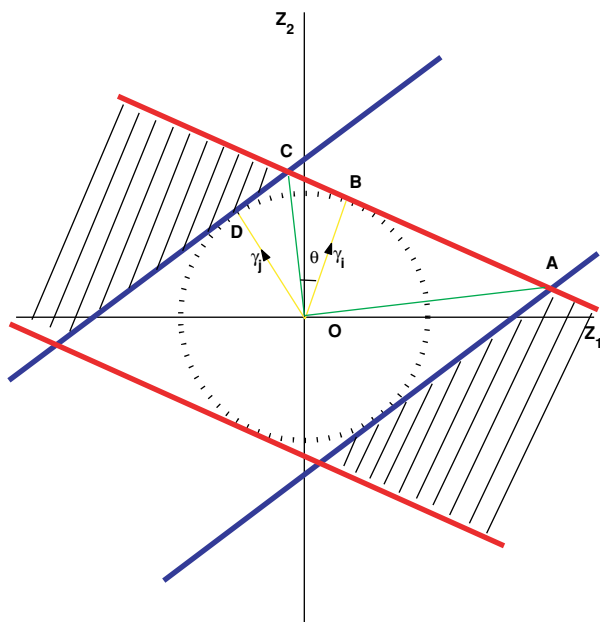
## Discussion

A computationally feasible single-marker analysis is usually required at the initial stage of a GWAS to scan through the

**Figure 1** P-values of MAX statistics along the chromosome for 343,413 SNPs in the GWAS of CAD.

**Table 4** Estimated P-values for 22 Chosen SNPs from The GWAS of CAD

| SNP ID | Chromosome | Position | $r_0$ | $r_1$ | $r_2$ | $s_0$ | $s_1$ | $s_2$ | Rhombus Formula | 10,000,000 Bootstrap | 10,000,000 Permutation |
|---|---|---|---|---|---|---|---|---|---|---|---|
| rs4854090 | 2 | 240888623 | 43 | 788 | 1025 | 155 | 1127 | 1584 | $5.88 \times 10^{-7}$ | $7.00 \times 10^{-7}$ | $8.00 \times 10^{-7}$ |
| rs5007171 | 3 | 45124712 | 29 | 684 | 1173 | 126 | 966 | 1803 | $2.03 \times 10^{-7}$ | $4.00 \times 10^{-7}$ | $<10^{-7}$ |
| rs7044859 | 9 | 22008781 | 449 | 972 | 502 | 539 | 1447 | 950 | $1.34 \times 10^{-7}$ | $<10^{-7}$ | $2.00 \times 10^{-7}$ |
| rs523096 | 9 | 22009129 | 649 | 960 | 314 | 857 | 1442 | 637 | $3.81 \times 10^{-6}$ | $3.70 \times 10^{-6}$ | $3.90 \times 10^{-6}$ |
| rs518394 | 9 | 22009673 | 648 | 964 | 310 | 856 | 1445 | 631 | $3.77 \times 10^{-6}$ | $2.30 \times 10^{-6}$ | $3.90 \times 10^{-6}$ |
| rs10757264 | 9 | 22009732 | 509 | 976 | 437 | 623 | 1485 | 828 | $5.21 \times 10^{-7}$ | $4.00 \times 10^{-7}$ | $2.00 \times 10^{-7}$ |
| rs10965212 | 9 | 22013795 | 513 | 950 | 459 | 603 | 1428 | 895 | $2.74 \times 10^{-9}$ | $<10^{-7}$ | $<10^{-7}$ |
| rs1292136 | 9 | 22014351 | 382 | 951 | 590 | 746 | 1482 | 708 | $1.40 \times 10^{-8}$ | $<10^{-7}$ | $<10^{-7}$ |
| rs7049105 | 9 | 22018801 | 506 | 950 | 461 | 593 | 1439 | 897 | $3.02 \times 10^{-9}$ | $<10^{-7}$ | $<10^{-7}$ |
| rs10965215 | 9 | 22019445 | 506 | 954 | 463 | 588 | 1447 | 902 | $1.48 \times 10^{-9}$ | $<10^{-7}$ | $<10^{-7}$ |
| rs564398 | 9 | 22019547 | 721 | 925 | 277 | 921 | 1428 | 583 | $4.99 \times 10^{-8}$ | $<10^{-7}$ | $<10^{-7}$ |
| rs7865618 | 9 | 22021005 | 706 | 929 | 289 | 888 | 1429 | 619 | $3.80 \times 10^{-9}$ | $<10^{-7}$ | $<10^{-7}$ |
| rs10965219 | 9 | 22043687 | 530 | 956 | 439 | 613 | 1444 | 876 | $1.45 \times 10^{-10}$ | $<10^{-7}$ | $<10^{-7}$ |
| rs9632884 | 9 | 22062301 | 592 | 949 | 381 | 693 | 1422 | 818 | $1.13 \times 10^{-12}$ | $<10^{-7}$ | $<10^{-7}$ |
| rs6475606 | 9 | 22071850 | 588 | 957 | 380 | 683 | 1425 | 830 | $1.31 \times 10^{-13}$ | $<10^{-7}$ | $<10^{-7}$ |
| rs4977574 | 9 | 22088574 | 382 | 937 | 605 | 804 | 1435 | 698 | $1.27 \times 10^{-12}$ | $<10^{-7}$ | $<10^{-7}$ |
| rs2891168 | 9 | 22088619 | 383 | 938 | 605 | 803 | 1435 | 698 | $1.75 \times 10^{-12}$ | $<10^{-7}$ | $<10^{-7}$ |
| rs1333042 | 9 | 22093813 | 365 | 939 | 617 | 770 | 1426 | 724 | $7.23 \times 10^{-12}$ | $<10^{-7}$ | $<10^{-7}$ |
| rs1333048 | 9 | 22115347 | 619 | 951 | 354 | 730 | 1423 | 781 | $3.80 \times 10^{-13}$ | $<10^{-7}$ | $<10^{-7}$ |
| rs1333049 | 9 | 22115503 | 586 | 960 | 378 | 676 | 1431 | 829 | $5.36 \times 10^{-14}$ | $<10^{-7}$ | $<10^{-7}$ |
| rs17697005 | 16 | 12268766 | 154 | 987 | 715 | 361 | 1398 | 1105 | $9.36 \times 10^{-6}$ | $8.70 \times 10^{-6}$ | $9.80 \times 10^{-6}$ |
| rs688034 | 22 | 25014189 | 823 | 832 | 265 | 1386 | 1266 | 276 | $5.81 \times 10^{-6}$ | $4.80 \times 10^{-6}$ | $6.30 \times 10^{-6}$ |

**Figure 2** Graphical representation of
$\bar{E}_i E_j = \{|T_i| < t, |T_j| > t\}$ (The shaped region), where
$T_i = \gamma'_i(Z_1, Z_2)'$, $T_j = \gamma'_j(Z_1, Z_2)'$, and $(Z_1, Z_2) \sim N(\mathbf{0}, \mathbf{I}_2)$,
$\mathbf{0} = (0, 0)'$, $\mathbf{I}_2 = \left(\begin{smallmatrix} 1 & 0 \\ 0 & 1 \end{smallmatrix}\right)$.

genome in order to prioritize SNPs for subsequent studies or to identify SNPs that reach the global significance level. The MAX statistic, which is defined as the maximum of several statistics targeting alternative hypotheses for the disease model, is a good candidate for the single-marker analysis because it can retain high power across a wide range of disease models. This robustness is particularly attractive for a GWAS, since it is unlikely that all disease-associated markers follow the same disease model. But its application in GWAS is limited by the difficulty in assessing the significance level of the MAX test. In this paper, we derive a simple approximation formula, called the rhombus formula, for estimating the p-value of the MAX test. Nevertheless, multiple-integration (Conneely & Boehnke, 2007) could be an alternative method for calculating adjusted p-values. Compared with multiple-integration, our method has an analytic expression and is more convenient to use. Our method can be applied to the MAX test with or without adjustment for the effect of covariates, based on three CATTs derived for three alternative disease penetrance models. It doesn't require resampling steps (permutation or bootstrap) and thus is readily applicable to GWAS.

The rhombus formula provides a theoretical upper bound for p-values under the normal assumptions. In real applications, using this upper bound tends to overestimate the true p-value. This formula is particularly suitable for approximating low p-values, but it is less accurate for estimating p-values above 0.2, as is evident from the simulation studies. However this defect should not limit its application in GWAS, where we are interested primarily in identifying SNPs with relatively low p-values.

With the rhombus formula, the MAX test can be used routinely in GWAS. In this paper, we focus mainly on the operational aspects of the MAX test, such as how to estimate the p-value and how to do the MAX test with adjustment for covariate effects. It is important to evaluate the impact of using the MAX test on the design and analysis of the GWAS. Although it is not straightforward to derive an analytic power calculation formula for the MAX test, it is computationally feasible to evaluate its power and other properties through simulation studies, using the rhombus formula.

In practice, case-control studies are susceptible to various confounding effects. One issue in case-control design is population stratification, which leads to spurious associations when the allele (genotype) frequencies and disease prevalence change across subpopulations. Various approaches (e.g., Devlin & Roeder, 1999; Price et al. 2006; Zheng et al. 2006b, and Li & Yu, 2008) have been proposed for the correction of population stratification. We are currently investigating the effect of population stratification on the MAX test and how to apply those correction methods with MAX to GWAS.

## Acknowledgements

## References

Conneely, K. N. & Boehnke, M. (2007) So many correlated tests, so little time Rapid adjustment of p-values for multiple correlated tests. *Am J Hum Genet* **81**, 1158–1168.

Devlin, B. & Roeder, K. (1999) Genomic control for association studies. *Biometrics* **55**, 997–1004.

Efron, B. (1997) The length heuristic for simultaneous hypothesis tests. *Biometrika* **84**, 143–157.

Freidlin, B., Podgor, M. J. & Gastwirth, J. L. (1999) Efficiency robust tests for survival or ordered categorical data. *Biometrics* **55**, 883–886

Gastwirth, J. L. (1985) The use of maximin efficiency robust tests in combining contingency tables and survival analysis. *J Am Stat Assoc* **80**, 380–384.

Hoh, J. & Ott, J. (2003) Mathematical multi-locus approaches to localizing complex human trait genes. *Nat Rev Genet* **4**, 701–709.

Hunter, D. J., Kraft, P., Jacobs, K. B., Cox, D. G., Yeager, M., Hankinson, S. E., Wacholder, S., Wang, Z., Welch, R., Hutchinson, A. et al. (2007) A genome-wide association study identifies alleles in FGFR2 associated with risk of sporadic postmenopausal breast cancer. *Nat Genet* **39**, 870–874.

Klein, R. J., Zeiss, C., Chew, E. Y., Tsai, J. Y., Sackler, R. S., Haynes, C., Henning, A. K., SanGiovanni, J. P., Mane, S. M., Mayne, S. T. et al. (2005) Complement factor H polymorphism in aged-related macular degeneration. *Science* **308**, 385–389.

Li, Q. Z. & Yu, K. (2008) Improved correction for population stratification in genome-wide association studies by identifying hidden population structures. *Genet Epidemiol* (In press).

Liang, K. Y. & Zeger, S. L. (1986) longitudinal data analysis using generalized linear models. *Biometrika* **73**, 13–22.

Marchini, J., Donnelly, P. & Cardon, L. R. (2005) Genome-wide strategies for detecting multiple loci that influence complex diseases. *Nat Genet* **37**, 413–417.

Pepe, M. S., Whitaker, B. C. & Seidel, K. (1999) Estimating and comparing univariate associations with application to the prediction of adult obesity. *Stat Med* **18**, 163–173.

Price, A. L., Patterson, N. J., Plenge, R. M., Weinblatt, M. E., Shadick, N. A., Reich, D. (2006) Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet* **38**, 904–909.

Sasieni, P. D. (1997) From genotypes to genes: Doubling the sample size. *Biometrics* **53**, 1253–1261.

Schaid, D. J., McDonnell, S. K., Hebbring, S. J., Cunningham, J. M. & Thibodeau, S. N. (2005) Nonparametric tests of association of multiple genes with human diseases. *Am J Hum Genet* **76**, 780–793.

Skol, A. D., Scott, L. J., Abecasis, G. R. & Boehnke, M. (2006) Joint analysis is more efficient than replication-based analysis for two-stage genome-wide association studies. *Nat Genet* **38**, 209–213.

Sladek, R., Rocheleau, G., Rung, J., Dina, C., Shen, L., Serre, D., Boutin, P., Vincent, D., Belisle, A., Hadjadj, S. et al. (2007) A genome-wide association study identifies novel risk loci for type 2 diabetes. *Nature* **445**, 881–885.

Slager, S. L. & Schaid, D. J. (2001) Case-control studies of genetic markers: power and sample size approximations for Armitage's test for trend. *Hum Hered* **52**, 149–153.

The Wellcome Trust Case Control Consortium (2007) Genome-wide association study of 14,000 cases of seven common diseases and 3000 shared controls. *Nature* **447**, 661–678.

Wang, T., Zhu, X. & Elston, R. C. (2007) Improving power in contrasting linkage-disequilibrium patterns between cases and controls. *Am J Hum Genet* **80**, 911–920.

Yan, L. H., Zheng, G. & Li, Z. H. (2008) Two-stage group sequential robust tests in family-based association studies: controlling type I error. *Ann Hum Genet* (Tentatively accepted).

Yeager, M., Orr, N., Hayes, R. B., Jacobs, K. B., Kraft, P. et al. (2007) Genome-wide association study of prostate cancer identifies a second risk locus at 8q24. *Nat Genet* **39**, 645–649.

Yu, K., Chatterjee, N., Wheeler, W., Li, Q. Z., Wang, S. et al. (2007) Flexible Design for Following Up Positive Findings. *Am J Hum Genet* **81**, 540–551.

Zhang, S., Zhu, X. & Zhao, H. (2003) On a Semiparametric Test to Detect Associations Between Quantitative Traits and Candidate Genes Using Unrelated Individuals. *Genet Epidemiol* **24**, 44–56.

Zheng, G., Freidlin, B., Li, Z. H. & Gastwirth, J. L. (2003) Choice of scores in trend tests for case-control studies of candidate gene associations. *Biometrical J* **45**, 335–348.

Zheng, G., Friedlin, B. & Gastwirth, J. L. (2006a) Comparison of robust tests for genetic association using case-control studies. *IMS Lecture Notes-Monograph Series. 2nd Eric L. Lehmann Symposium − Optimality* **49**, 253–265.

Zheng, G., Freidlin, B. & Gastwirth, J. L. (2006b) Robust genomic control for association studies. *Am J Hum Genet* **78**, 350–356.

Zheng, G., Joo, J., Lin, J. P., Stylianou, M. & Waclawiw, M. A. et al. (2007) Robust ranks of true associations in genome-wide case-control association studies. *BMC Proceedings* **1**(Suppl 1), S165.

## Appendix

**Derivation of the Rhombus Formula:** For any given $t$, define $2k$ events:

$$E_1 = \{|T_1| > t\}, \bar{E}_1 = \{|T_1| \leq t\}, \cdots,$$
$$E_k = \{|T_k| > t\}, \bar{E}_k = \{|T_k| \leq t\}$$

Then $\Pr(T_{\max} > t) = \Pr\left(\bigcup_{i=1}^{k} E_i\right)$

$$\leq \Pr(E_1) + \Pr(\bar{E}_1 E_2)$$
$$+ \Pr(\bar{E}_2 E_3) + \cdots + \Pr(\bar{E}_{k-1} E_k).$$

Let $\mathbf{Z} \sim N(\mathbf{0}, \mathbf{I}_2)$, where $\mathbf{0} = (0, 0)^T$ and $\mathbf{I}_2 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$. Then $T_i$ and $T_j$ can be expressed as $T_i = \gamma_i^T \mathbf{Z}$ and $T_j = \gamma_j^T \mathbf{Z}$, where $\gamma_i$ and $\gamma_j$ are vectors such that $\text{cor}(T_i, T_j) = \gamma_i^T \gamma_j$. As illustrated in Figure 2, $E_i$ and $E_j$ are two half-spaces. Their boundaries are two tangents that are perpendicular to $\gamma_i$ and $\gamma_j$, respectively. The distance between the tangent points and the original point is $t$. $\bar{E}_i E_j$ is the shaded region lying between the two parallel lines. Let $\Omega$ be the middle rhombus-shaped region. Then, we have

$$\Pr(\bar{E}_i E_j) = \Phi(t) - \Phi(-t) - \Pr(\Omega).$$

$\Pr(\Omega)$ is four times the probability of the event defined by triangle AOC. Expressing $\mathbf{Z}$ in polar coordinates, we have

$$\Pr(\bar{E}_i E_j) = \Phi(t) - \Phi(-t) - 4$$
$$\times \left[\int_0^{\frac{L_{ij}}{2}} \int_0^{t \sec(\theta)} \frac{1}{2\pi} r\, e^{-\frac{r^2}{2}}\, dr\, d\theta \right.$$
$$\left. + \int_0^{\frac{\pi - L_{ij}}{2}} \int_0^{t \sec(\theta)} \frac{1}{2\pi} r\, e^{-\frac{r^2}{2}}\, dr\, d\theta \right]$$

$$= \Phi(t) - \Phi(-t) - 1 + \frac{2}{\pi}\left[\int_0^{\frac{L_{ij}}{2}} e^{-\frac{t^2 \sec^2(\theta)}{2}} d\theta \right.$$

$$\left. + \int_0^{\frac{\pi - L_{ij}}{2}} e^{-\frac{t^2 \sec^2(\theta)}{2}} d\theta\right].$$

Since $\sec^2(\theta) \geq 1 + \theta^2$ for $0 < \theta < \frac{\pi}{2}$, we have

$$\Pr\left(\bar{E}_i E_j\right) \leq \Phi(t) - \Phi(-t) - 1$$

$$+ \frac{2}{\pi}\left[2 \int_0^{\frac{L_{ij}}{2}} e^{-\frac{t^2(1+\theta^2)}{2}} d\theta\right.$$

$$\left. + \int_{\frac{L_{ij}}{2}}^{\frac{\pi - L_{ij}}{2}} e^{-\frac{t^2(1+\theta^2)}{2}} d\theta\right] I\left\{0 \leq L_{ij} \leq \frac{\pi}{2}\right\}$$

$$+ \frac{2}{\pi}\left[\int_{\frac{\pi - L_{ij}}{2}}^{\frac{L_{ij}}{2}} e^{-\frac{t^2(1+\theta^2)}{2}} d\theta\right.$$

$$\left. + 2 \int_0^{\frac{\pi - L_{ij}}{2}} e^{-\frac{t^2(1+\theta^2)}{2}} d\theta\right] I\left\{\frac{\pi}{2} < L_{ij} \leq \pi\right\}$$

$$\leq \Phi(t) - \Phi(-t) - 1 + 4\phi(t)\left\{\frac{2\left[\Phi\left(\frac{t L_{ij}}{2}\right) - 0.5\right]}{t}\right.$$

$$\left. + \frac{e^{-\frac{t^2 L_{ij}^2}{8}}\left[\Phi\left(\frac{t(\pi - L_{ij})}{2}\right) - \Phi\left(\frac{t L_{ij}}{2}\right)\right]}{t}\right\} I\left\{0 \leq L_{ij} \leq \frac{\pi}{2}\right\}$$

$$+ 4\phi(t)\left\{\frac{e^{-\frac{t^2(\pi - L_{ij})^2}{8}}\left[\Phi\left(\frac{t L_{ij}}{2}\right) - \Phi\left(\frac{t(\pi - L_{ij})}{2}\right)\right]}{t}\right.$$

$$\left. + \frac{2\left[\Phi\left(\frac{t(\pi - L_{ij})}{2}\right) - 0.5\right]}{t}\right\} I\left\{\frac{\pi}{2} \leq L_{ij} \leq \pi\right\}.$$

Hence, we obtain

$$\Pr\left(T_{\max} > t\right) \leq (k - 2)\left[\Phi(t) - \Phi(-t) - 1\right]$$

$$- \frac{4\phi(t)(k - 1)}{t}$$

$$+ \frac{4\phi(t)}{t}\sum_{i=1}^{k-1}\left\{2\Phi\left(\frac{t L_{i(i+1)}}{2}\right)\right.$$

$$\left. + e^{-\frac{t^2 L_{i(i+1)}^2}{8}}\left[\Phi\left(\frac{t(\pi - L_{i(i+1)})}{2}\right) - \Phi\left(\frac{t L_{i(i+1)}}{2}\right)\right]\right\}$$

$$\times I\left\{0 \leq L_{i(i+1)} \leq \frac{\pi}{2}\right\}$$

$$+ \frac{4\phi(t)}{t}\sum_{i=1}^{k-1}\left\{2\Phi\left(\frac{t(\pi - L_{i(i+1)})}{2}\right)\right.$$

$$\left. + e^{-\frac{t^2[\pi - L_{i(i+1)}]^2}{8}}\left[\Phi\left(\frac{t L_{i(i+1)}}{2}\right) - \Phi\left(\frac{t(\pi - L_{i(i+1)})}{2}\right)\right]\right\}$$

$$\times I\left\{\frac{\pi}{2} \leq L_{i(i+1)} \leq \pi\right\}.$$