# Accepted Manuscript

Numerical evaluation of methods approximating the distribution of a large quadratic form in normal variables

Tong Chen, Thomas Lumley

Please cite this article as: T. Chen and T. Lumley, Numerical evaluation of methods approximating the distribution of a large quadratic form in normal variables. *Computational Statistics and Data Analysis* (2019), https://doi.org/10.1016/j.csda.2019.05.002

# Numerical evaluation of methods approximating the distribution of a large quadratic form in normal variables

Tong Chen[a,*], Thomas Lumley[a]

[a]*University of Auckland*

## Abstract

Quadratic forms of Gaussian variables occur in a wide range of applications in statistics. They can be expressed as a linear combination of chi-squareds. The coefficients in the linear combination are the eigenvalues $\lambda_1, ..., \lambda_n$ of $\Sigma A$, where $A$ is the matrix representing the quadratic form and $\Sigma$ is the covariance matrix of the Gaussians. The previous literature mostly deals with approximations for small quadratic forms ($n < 10$) and moderate p-values ($p > 10^{-2}$). Motivated by genetic applications, moderate to large quadratic forms ($300 < n < 12,000$) and small to very small p-values ($p < 10^{-4}$) are studied. Existing methods are compared under these settings and a leading-eigenvalue approximation, which only takes the largest $k$ eigenvalues, is shown to have the computational advantage without any important loss in accuracy. For time complexity, a leading-eigenvalue approximation reduces the computational complexity from $O(n^3)$ to $O(n^2k)$ on extracting eigenvalues and avoids speed problems with computing the sum of $n$ terms. For accuracy, the existing methods have some limits on calculating small p-values under large quadratic forms. Moment methods are inaccurate for very small p-values, and Farebrother's method is not usable if the minimum eigenvalue is much smaller than others. Davies's method is usable for p-values down to machine epsilon. The saddlepoint approximation is proved to have bounded relative error for any $A$ and $\Sigma$ in the extreme right tail, so it is usable for arbitrarily small p-values.

*Keywords:* small p-values, leading-eigenvalue approximation, accuracy, computational complexity

## 1. Introduction

A quadratic form can be expressed as $Q(X) = X^\top A X$, where $X = (X_1, \ldots, X_n)^\top$ is a multivariate normal random vector with mean vector $\mu = (\mu_1, \ldots, \mu_n)$ and covariance matrix $\Sigma$, and $A$ is a $n \times n$ symmetric and non-negative definite matrix. The question of interest is to estimate the upper tail probability of $Q(X)$

$$Pr(Q(X) > q), \tag{1}$$

where $q$ is a scalar.

The distribution of $Q(x)$ is a linear combination of noncentral $\chi_1^2$ variables, where the coefficients are the non-zero eigenvalues $\lambda_1, \ldots \lambda_n$ of matrix $M = \Sigma A = XX^T$. When $\mu = \mathbf{0}_n$, it is a linear combination of central $\chi_1^2$ variables.

These quadratic forms often occur when a set of asymptotically Normal test statistics are combined using a weight matrix other than the inverse of their covariance matrix. A famous example is the Rao-Scott test (Rao and Scott, 1981) in survey statistics. The true variance matrix of the individual test statistics tends to be poorly estimated; the Rao-Scott test replaces it with the variance matrix under iid sampling. In genomics, the Sequence Kernel Association Test (SKAT) (Wu et al., 2011) evaluates the association between rare variants and phenotype. It replaces the true variance matrix with a set of weights that ignore correlation and upweight less-common variants, corresponding to a diagonal matrix $A$.

The null distribution of these tests is a weighted sum of central $\chi_1^2$ variables, where the coefficients are the eigenvalues of $M$. Many methods are proposed to evaluate the upper tail probability of the distribution of $Q(X)$. We classified these existing methods into three categories: 'exact' methods (Davies, 1980; Farebrother, 1984; Bausch, 2013), moment methods (see, eg., the Satterthwaite approximation and Liu et al. (2009)) and a saddlepoint approximation (Kuonen, 1999).

The 'exact' methods are exact in the sense that an approximation with arbitrary accuracy could be obtained if arbitrary precision arithmetic were available. Davies (1980) exploited the fact that the characteristic function of a sum is the product of characteristic functions, so the characteristic function for a weighted sum of $\chi_1^2$ variables is straightforward to obtain. Farebrother (1984) showed that the tail probability can be written as an infinite series of central chi-squared distributions, by writing the linear combination as a mixture (Robbins and Pitman, 1949). Bausch (2013) showed that a linear combination of gamma densities form an algebra under convolutions and derived the density for weighted sums of $\chi_k^2$ variables.

The Satterthwaite approximation approximates the distribution of $Q(X)$ by $a\chi_d^2$ with $a$ and $d$ chosen to give the correct mean and variance. Liu et al. (2009) proposed a four-moment approximation using a noncentral chi-squared distribution of the form $a + b\chi_d^2(\nu)$, where $a$ is an offset, $b$ is a scaling parameter and $\nu$ is the non-centrality parameter. Kuonen (1999) derived a form of saddlepoint approximation to the sum. The accuracy of these approximations has been previous studied (Kuonen, 1999; Duchesne and De Micheaux, 2010; Bausch, 2013), but only for small quadratic forms ($n < 10$) and moderate p-values.

However, genetics studies often involve a large number of terms ($n > 1000$) and small p-values ($p < 10^{-4}$) raising concerns about both time complexity and accuracy. For time complexity, extracting all set of eigenvalues scales as cube of sample size $n$ and it would take more time to compute a tail probability when the number of terms $n$ is large. For accuracy, moment methods are anti-conservative in the right tail of the distribution.

Recently, a companion paper (Lumley et al., 2018) developed a leading-eigenvalue approximation to solve above problems. This method is mainly developed for large quadratic forms and ends up with less computational time without any important loss in accuracy. This is done by extracting the largest $k$ eigenvalues using a low-rank stochastic singular

2

value decomposition (SSVD) (Halko et al., 2011) and utilizing the cheap Satterthwaite approximation to approximate the rest $n - k$ terms.

This work is motivated by genetic problems which often involve large quadratic forms with thousands or tens of thousands of terms, under which the existing methods would have a computational deficiency and may be less accurate. The main objective is to find an optimal way to perform convolutions for large quadratic forms. We provide empirical evidence for the existing methods and a leading-eigenvalue approximation under moderate and large quadratic forms. Evaluations and discussions of the existing methods under large quadratic forms are made in Section 2. In Section 3, accuracy and computational complexity of a leading-eigenvalue approximation are discussed. Impact of sparsity, rank and definiteness of matrix $M$ is discussed in Section 4. Discussions are made in Section 5.

R codes for producing numerical examples can be found in Supplementary information and are available from `https://github.com/T0ngChen/LargeQuadraticForm`.

## 2. Existing methods under genetic settings

This section evaluates the performance of the existing methods in the right tail of the distribution. Davies's (1980) and Farebrother's (1984) methods are usable even for thousands of terms and achieve close to their nominal accuracy as long as the right tail probability is much larger than machine epsilon. As they compute Equation (1) from $1 - Pr(Q(X) < q)$, they break down completely if the extreme right tail probabilities are near or beyond machine epsilon. The value of machine epsilon mentioned in this work is $2^{-52} \approx 2 \times 10^{-16}$.

We observed that Farebrother's (1984) method ended up with fault indicator 1 when it was evaluated using the quadratic forms $Q_1$–$Q_6$ generated in this section. The fault indicator 1 represents the calculation has non-fatal underflow of a variable called $a_0$ (Farebrother, 1984). If $Q(X)$ is a weighted sum of central $\chi_1^2$ variables, the quantity $a_0$ in Farebrother's (1984) algorithm can be simplified to

$$ a_0 = \exp\left( \frac{1}{2} \left( n \log \lambda_n - \sum_i^n \log \lambda_i \right) \right), $$

where $\lambda_1, \ldots, \lambda_n$ are sorted eigenvalues in descending order. For large quadratic forms, if $\lambda_n$ is much smaller than other eigenvalues, a large $n$ can cause the variable $a_0$ to underflow to 0. We cannot use Farebrother's (1984) method as a reference because the leading eigenvalues are much larger than the minimum eigenvalue in our simulated genome sequence data. Bausch's (2013) method has rounding errors especially in the left tail with double precision for moderate and large quadratic forms and is slow with multiple precision (see Supplementary information). We hereafter choose Davies's (1980) method as a reference to conduct numerical studies.

To evaluate the performance of these approximation methods, we simulated human genome sequence data using the Markov Coalescent Simulator (Chen et al., 2009). This was done by fixing the number of rows $s$ (people) then choosing the length to make the number of columns $m$ (variants) approximately equal to the number of rows. We discarded

variants with minor allele frequency greater than 5% to filter rare variants, giving a large sparse matrix. We generated six data sets $Q_1$–$Q_6$ and their dimensions are shown in Table 1. When we extracted eigenvalues, they are set to be zero if they are smaller than $10^{-10}$. As the exact methods have internal estimates of accuracy, we compare Davies's (1980) method with simulation p-value to verify Davies's (1980) method is exact and can be used for moderate and large quadratic forms (see Supplementary information).

| Name | $Q_1$ | $Q_2$ | $Q_3$ | $Q_4$ | $Q_5$ | $Q_6$ |
|------|-------|-------|-------|-------|-------|-------|
| $s$  | 500   | 1000  | 2000  | 7000  | 9000  | 20000 |
| $m$  | 470   | 987   | 1643  | 7352  | 8887  | 22456 |
| $n$  | 305   | 637   | 1063  | 3985  | 4834  | 11259 |

Table 1: The dimensions of simulated human genome sequence data, where $s$ is the number of people, $m$ is the number of variants and $n$ is the number of non-zero eigenvalues.

Next, we compare the accuracy of the Satterthwaite approximation, Liu–Tang–Zhang's (2009) four-moment approximation and Kuonen's (1999) saddlepoint approximation when p-value is greater than machine epsilon. R (R Core Team, 2017) packages `survey` (Lumley, 2011) and `CompQuadForm` (Duchesne and De Micheaux, 2010) are used to perform analysis in this section. The eigenvalues of $Q_1$ to $Q_6$ are extracted using a full eigendecomposition.
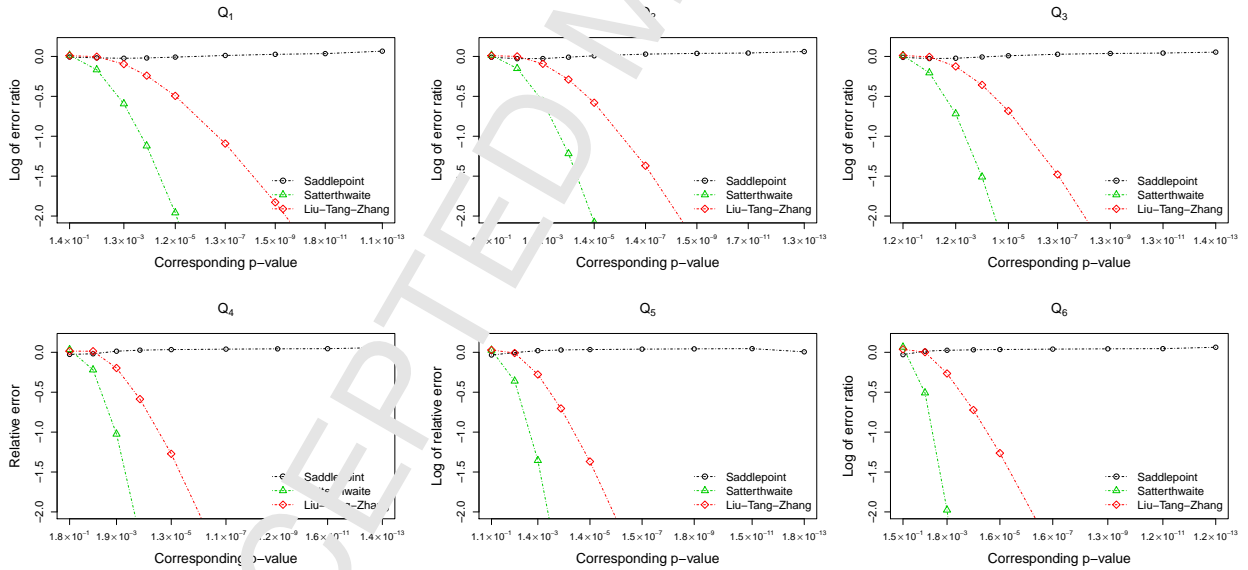


Figure 1 Comparisons between methods for the quadratic forms $Q_1$–$Q_6$ when p-value is greater than machine epsilon. Exact values are computed using Davies's (1980) method with accuracy $10^{-16}$.

Results are presented in Figure 1 and Table 1 in Supplementary information. The x-axis represents corresponding underlying true p-value from $10^{-1}$ down to $10^{-13}$. The y-axis represents the logarithm of error ratio to the base 10. It is computed by generating the underlying true p-values using Davies's (1980) method and then calculating the logarithm

4

ratio of each method to Davies's (1980) method. Our numerical studies show that Kuonen's (1999) saddlepoint approximation is highly accurate. The maximum logarithm of error ratio is less than 0.07 in all cases. Moment methods have better performance than Kuonen's (1999) saddlepoint approximation in the left tail but are anti-conservative in the right tail. The Satterthwaite approximation is accurate if the p-value is greater than $10^{-1}$. Liu–Tang–Zhang's (2009) four-moment approximation performs better and is accurate until $10^{-2}$. After that, the logarithm of error ratio for both two moment methods increases very fast. Figure 1 also shows that moment methods tend to have better performance for moderate quadratic forms than large quadratic forms under our simulated human genome sequence data.

If the p-value is smaller than machine epsilon, neither the moment methods nor the exact methods work. All that's left is the saddlepoint approximation. To analyse the extreme right tail performance of the saddlepoint approximation, we consider exponential tail rates. A linear combination of chi-squared variables have an exponential tail in the sense of Berman (1992), with tail rate $1/2\lambda_1$. We show in Appendix A that the saddlepoint approximation has the same exponential tail rate in the extreme right tail, so that the relative error in $Pr(Q(X) > q)$ is bounded as $q \to \infty$, for any $A$ and $\Sigma$. Kuonen (1999) showed that the relative error is of order $o(n^{-3/2})$, so the approximation improves with increasing $n$, and the saddlepoint approximation can be used as a reference in the extreme right tail.

## 3. A leading-eigenvalue approximation under genetic settings

This section explores accuracy and time complexity for a leading-eigenvalue approximation. It approximates the distribution of $Q(X)$ by formula (4) in Lumley et al. (2018) which is

$$T \sim \left( \sum_{i=1}^{k} \lambda_i \chi_1^2 \right) + a\chi_d^2, \tag{2}$$

where $\lambda_1, ..., \lambda_k$ are the largest $k$ eigenvalues of matrix $M$, $a = (\sum_{k+1}^{n} \lambda_i^2)/(\sum_{k+1}^{n} \lambda_i)$ and $d = (\sum_{k+1}^{n} \lambda_i)^2/(\sum_{k+1}^{n} \lambda_i^2)$.

The leading eigenvalues are extracted using a low-rank SSVD (Halko et al., 2011). In Equation (2), the leading terms can be combined using either the exact methods or the saddlepoint approximation, and the remainder term is obtained by the Satterthwaite approximation. Followed by the performance of approximation methods discussed in Section 2, the leading terms are combined using Davies's (1980) method if the p-value is greater than machine epsilon and using Kuonen's (1999) saddlepoint approximation if the p-value is near or beyond machine epsilon. If the number of leading eigenvalues $k$ is much smaller than $n$, a leading-eigenvalue approximation also works well with Farebrother's (1984) method, because the variable $a_0$ would not underflow to zero for small quadratic forms.

For accuracy, we compare the leading-eigenvalue approximation with Davies's (1980) method when p-value is greater than machine epsilon and with Kuonen's (1999) saddlepoint approximation in the extreme right tail using data generated in Section 2. R package bigQF (Lumley, 2019) is used to do the leading-eigenvalue approximation. SSVD uses 50, 50, 100, 100, 200 and 200 eigenvalues for quadratic forms $Q_1$–$Q_6$ respectively.

5

| | $q$ | $D$ | $L_D$ | $R_{L_D}$ | | $q$ | $D$ | $L_D$ | $R_{L_D}$ |
|---|---|---|---|---|---|---|---|---|---|
| $Q_1$ | $1.2\times10^{04}$ | $1.647\times10^{-04}$ | $1.654\times10^{-04}$ | 0.005 | $Q_2$ | $4.0\times10^{04}$ | $1.214\times10^{-04}$ | $1.213\times10^{-04}$ | 0.000 |
| | $1.6\times10^{04}$ | $1.511\times10^{-06}$ | $1.518\times10^{-06}$ | 0.005 | | $5.4\times10^{04}$ | $3.277\times10^{-07}$ | $3.276\times10^{-07}$ | -0.001 |
| | $1.9\times10^{04}$ | $1.473\times10^{-08}$ | $1.480\times10^{-08}$ | 0.005 | | $6.8\times10^{04}$ | $1.022\times10^{-09}$ | $1.021\times10^{-09}$ | -0.001 |
| | $2.3\times10^{04}$ | $1.513\times10^{-10}$ | $1.520\times10^{-10}$ | 0.004 | | $8.2\times10^{04}$ | $3.395\times10^{-12}$ | $3.390\times10^{-12}$ | -0.001 |
| $Q_3$ | $1.1\times10^{05}$ | $1.515\times10^{-04}$ | $1.512\times10^{-04}$ | -0.002 | $Q_4$ | $1.2\times10^{06}$ | $4.396\times10^{-04}$ | $4.361\times10^{-04}$ | -0.008 |
| | $1.5\times10^{05}$ | $2.770\times10^{-07}$ | $2.766\times10^{-07}$ | -0.002 | | $1.7\times10^{06}$ | $4.158\times10^{-07}$ | $4.127\times10^{-07}$ | -0.007 |
| | $1.9\times10^{05}$ | $5.894\times10^{-10}$ | $5.885\times10^{-10}$ | -0.001 | | $2.2\times10^{06}$ | $4.625\times10^{-10}$ | $4.592\times10^{-10}$ | -0.007 |
| | $2.3\times10^{05}$ | $1.348\times10^{-12}$ | $1.341\times10^{-12}$ | -0.005 | | $2.7\times10^{06}$ | $5.421\times10^{-13}$ | $5.438\times10^{-13}$ | 0.003 |
| $Q_5$ | $2.0\times10^{06}$ | $2.242\times10^{-05}$ | $2.242\times10^{-05}$ | 0.000 | $Q_6$ | $9.0\times10^{06}$ | $3.025\times10^{-04}$ | $3.024\times10^{-04}$ | -0.001 |
| | $2.5\times10^{06}$ | $1.515\times10^{-07}$ | $1.515\times10^{-07}$ | 0.000 | | $1.2\times10^{07}$ | $8.826\times10^{-07}$ | $8.820\times10^{-07}$ | -0.001 |
| | $3.0\times10^{06}$ | $1.091\times10^{-09}$ | $1.091\times10^{-09}$ | 0.000 | | $1.5\times10^{07}$ | $2.872\times10^{-09}$ | $2.870\times10^{-09}$ | -0.001 |
| | $3.5\times10^{06}$ | $8.134\times10^{-12}$ | $8.126\times10^{-12}$ | -0.001 | | $1.8\times10^{07}$ | $9.832\times10^{-12}$ | $9.806\times10^{-12}$ | -0.003 |

Table 2: Probability that the quadratic forms $Q_1$–$Q_6$ exceed $q$, $D$: exact value using Davies's (1980) method with accuracy $10^{-16}$; $L_D$: the leading-eigenvalue approximation where the leading eigenvalues are combined using Davies's (1980) method; $R_{L_D}$: $(L_D - D)/D$.

| | $q$ | $S$ | $L_S$ | $R_{L_S}$ | | $q$ | $S$ | $L_S$ | $R_{L_S}$ |
|---|---|---|---|---|---|---|---|---|---|
| $Q_1$ | $2.8\times10^{04}$ | $2.364\times10^{-13}$ | $2.374\times10^{-13}$ | 0.004 | $Q_2$ | $8.0\times10^{04}$ | $8.455\times10^{-12}$ | $8.449\times10^{-12}$ | -0.001 |
| | $3.8\times10^{04}$ | $8.871\times10^{-19}$ | $8.910\times10^{-19}$ | 0.004 | | $1.0\times10^{05}$ | $2.578\times10^{-15}$ | $2.576\times10^{-15}$ | -0.001 |
| | $4.8\times10^{04}$ | $3.510\times10^{-24}$ | $3.525\times10^{-24}$ | 0.004 | | $1.2\times10^{05}$ | $8.124\times10^{-19}$ | $8.117\times10^{-19}$ | -0.001 |
| | $5.8\times10^{04}$ | $1.430\times10^{-29}$ | $1.436\times10^{-29}$ | 0.004 | | $1.4\times10^{05}$ | $2.614\times10^{-22}$ | $2.612\times10^{-22}$ | -0.001 |
| $Q_3$ | $2.0\times10^{05}$ | $1.406\times10^{-10}$ | $1.404\times10^{-10}$ | -0.001 | $Q_4$ | $3.0\times10^{06}$ | $1.090\times10^{-14}$ | $1.082\times10^{-14}$ | -0.007 |
| | $2.5\times10^{05}$ | $7.213\times10^{-14}$ | $7.203\times10^{-14}$ | -0.001 | | $3.5\times10^{06}$ | $1.356\times10^{-17}$ | $1.347\times10^{-17}$ | -0.007 |
| | $3.0\times10^{05}$ | $3.838\times10^{-17}$ | $3.832\times10^{-17}$ | -0.001 | | $4.0\times10^{06}$ | $1.714\times10^{-20}$ | $1.701\times10^{-20}$ | -0.007 |
| | $3.5\times10^{05}$ | $2.088\times10^{-20}$ | $2.085\times10^{-20}$ | -0.001 | | $4.5\times10^{06}$ | $2.189\times10^{-23}$ | $2.174\times10^{-23}$ | -0.007 |
| $Q_5$ | $4.0\times10^{06}$ | $6.899\times10^{-14}$ | $6.907\times10^{-14}$ | 0.001 | $Q_6$ | $2.0\times10^{07}$ | $2.524\times10^{-13}$ | $2.538\times10^{-13}$ | 0.006 |
| | $4.5\times10^{06}$ | $5.342\times10^{-16}$ | $5.349\times10^{-16}$ | 0.001 | | $2.5\times10^{07}$ | $2.125\times10^{-17}$ | $2.137\times10^{-17}$ | 0.006 |
| | $5.0\times10^{06}$ | $4.177\times10^{-18}$ | $4.185\times10^{-18}$ | 0.001 | | $3.0\times10^{07}$ | $1.845\times10^{-21}$ | $1.856\times10^{-21}$ | 0.006 |
| | $5.5\times10^{06}$ | $3.291\times10^{-20}$ | $3.297\times10^{-20}$ | 0.001 | | $3.5\times10^{07}$ | $1.634\times10^{-25}$ | $1.643\times10^{-25}$ | 0.006 |

Table 3: Probability that the quadratic forms $Q_1$–$Q_6$ exceed $q$, $S$: approximation obtained by a full eigendecomposition of Kuonen's (1999) saddlepoint approximation; $L_S$: the leading-eigenvalue approximation where the leading eigenvalues are combined using Kuonen's (1999) saddlepoint approximation; $R_{L_S}$: $(L_S - S)/S$.

Results are shown in Table 2 and 3. The convolutions of leading terms are approximated by Davies's (1980) method in Table 2 and Kuonen's (1999) saddlepoint approximation in Table 3. The relative error is less than 1% for all examples in the whole probability range. So that the leading-eigenvalue approximation is consistent with Davies's (1980) method when the p-value is much larger than machine epsilon and with the saddlepoint approximation in the extreme right tail. There is no important loss in accuracy for the leading-eigenvalue approximation. In Table 3, comparisons are made at very small p-values where the order is smaller than $10^{-20}$. As discussed in Section 2, the relative error of Kuonen's (1999) saddlepoint approximation is uniformly bounded as $q \to \infty$ and the approximation improves with increasing $n$. Numerical examples in Section 2 also show that Kuonen's (1999) saddlepoint approximation is highly accurate. It is reasonable to assume it will have the same accuracy

6

as $q \to \infty$. Therefore, a saddlepoint approximation and a leading-eigenvalue approximation combine to be usable for all p-values and all large enough numbers of variables.

For time complexity, except the moment methods, implementation of other existing methods needs to extract all the eigenvalues. It would cost $O(n^3)$ time to do a full eigen-decomposition for $X$ or $X^2$ (Golub and Van Loan, 2012). The computational complexity for SSVD (Halko et al., 2011) is of order $O(n^2 k)$ to get the largest $k$ eigenvalues. As $\sum_1^n \lambda_i = \text{trace}(M)$ and $\sum_1^n \lambda_i^2 = \text{trace}((M)^2)$, the remainder term of Equation (2), which is approximated by the Satterthwaite approximation, also takes $O(n^2)$ time. So that a leading-eigenvalue approximation would reduce the computational complexity from $O(n^3)$ to $O(n^2 k)$.

Moment methods are implemented by matching moments. The Satterthwaite approximation can be calculated in $O(n^2)$ time, but Liu–Tang–Zhang's (2009) four-moment approximation is no faster than singular value decomposition (SVD) because computing the fourth moment would take as much work ($n^3$ operations) as getting all the eigenvalues.

Even after the eigenvalues are computed, there is also a speed problem in adding up thousands or tens of thousands of terms. In order to achieve the same accuracy, Davies's (1980) and Farebrother's (1984) methods would spend more computational time for large $n$, because Davies's (1980) method needs more integration terms and Farebrother's (1984) method needs more terms in truncated series. These two methods would also take more time to compute a small p-value as the number of terms they need is dependent on accuracy. The computational time of moment methods and the saddlepoint approximation does not increase when the p-value is getting smaller.

For Davies's (1980) method, it is slow to get high accuracy if the sum is dominated by a small number of eigenvalues and the number of terms $n$ in the sum is large, because the number of integration terms is highly dependent on accuracy in this context. Table 4 shows that, for large quadratic forms, in order to get high accuracy ($10^{-13}$ in our example), the computational time of Davies's (1980) method increases with the largest eigenvalue.

| Case | A | B | C | D | E |
|------|------|------|------|------|------|
| Time(s) | 0.01 | 0.05 | 0.41 | 3.96 | 34.77 |

Table 4: Computational time of computing a single p-value around $10^{-6}$ with accuracy at $10^{-13}$ using Davies's (1980) method. Case A uses eigenvalues of $Q_5$, case B, C, D and E are obtained by multiplying the largest eigenvalue of case A by $10, 10^2, 10^3$ and $10^4$ respectively.

A leading-eigenvalue approximation would do well in such situation because only the largest $k$ eigenvalues are combined using either the exact methods or the saddlepoint approximation. A leading-eigenvalue approximation has the computational advantage in both computing the eigenvalues and adding them up.

However, unless $n$ is greater than hundreds, there is no reason to use the leading-eigenvalue approximation as it does not save any time. We compare computational time of SSVD and SVD for $Q_1$–$Q_4$ and a small example $Q_0$ ($s = 2000; m = 67$) provided in the SKAT package (Lee et al., 2017). SSVD uses 50 eigenvalues for $Q_0$ and 100 eigenvalues for $Q_1$–$Q_4$.

7

As shown in Table 5, SSVD does not save time for small $n$. For moderate and large $n$, the choice of $k$ is not important as long as $k$ is large enough. The criterion for the choice of $k$ is provided in Section 3.3 of companion paper (Lumley et al., 2018). As Table 2 and 3 show, the relative error does not increase way out in the tails. So the criterion is also applicable here even the p-value is much smaller than the companion paper (Lumley et al., 2018).

| Qudratic form | $Q_0$ | $Q_1$ | $Q_2$ | $Q_3$ | $Q_4$ |
|---|---|---|---|---|---|
| SVD(s) | 0.03 | 0.24 | 2.22 | 13.40 | 84.56 |
| SSVD(s) | 0.12 | 0.24 | 0.84 | 2.59 | 36.06 |

Table 5: Comparisons of computational time between SSVD and SVD.

## 4. Impact of sparsity, rank and definiteness of matrix $M$

Sparsity would affect the speed of computing eigenvalues, but the leading-eigenvalue approximation still has the computational advantage over a full eigendecomposition for moderate and large quadratic forms. SSVD (Halko et al., 2011) takes $k$ matrix multiplications. Suppose $M = XX^T$, if matrix $X$ is sparse with $\alpha n^2$ non-zero entries, a matrix-vector multiplication takes $\alpha n^2$ time, so the leading eigenvalues can be computed in $O(\alpha n^2 k)$ time. The setting in the companion paper (Lumley et al., 2018) was for situations where $X$ or $M$ is not sparse, but matrix $X$ is the product of a sparse matrix and a projection on to residuals for an adjustment model. If the number of adjustment variables is $p$, the leading eigenvalues are available in $O(k(\alpha n^2 + np^2))$.

If $X$ is a general dense matrix, a matrix-vector multiplication takes $n^2$ operations, computing the $k$ leading eigenvalues would take $O(n^2 k)$ time. Therefore, for moderate or large quadratic forms, the leading-eigenvalue approximation is always faster than a full decomposition, and the advantage can be larger if the matrix $X$ has a special structure.

The rank of matrix $M$ does not affect computational complexity, because the leading-eigenvalue approximation is not simply a low-rank approximation. The matrices simulated in Section 2 are not full rank, but their ranks are still much larger than $k$ and computation would be the same if they were full rank. Figure 2 in the companion paper (Lumley et al., 2018) illustrated this by comparing a leading-eigenvalue approximation with a rank-$k$ approximation, showing that the low-rank approximation is much less accurate.

Davies's (1980) method, moment methods and the saddlepoint approximation are usable when matrix $M$ has negative eigenvalues but Farebrother's (1984) method is not usable in such a situation. A leading-eigenvalue approximation thereafter also works for negative definite, negative semi-definite and indefinite matrices as long as convolutions of the leading eigenvalues are calculated using either Davies's (1980) method or the saddlepoint approximation.

## 5. Discussion

Moment methods are inaccurate for very small p-values. They use a single $\chi_d^2$ distribution to approximate the distribution of $Q(X)$ giving a right tail that decreases faster than the

8

true distribution. Except for the Satterthwaite approximation, the other moment methods are no faster than getting all the eigenvalues: computing the third moment would take as much work as extracting all the eigenvalues.

Davies's (1980) and Farebrother's (1984) methods are exact when the p-value is much larger than machine epsilon. However, for large quadratic forms, Farebrother's (1984) method breaks down if the minimum eigenvalue is small and Davies's method is slow to obtain high accuracy if the sum is dominated by a small number of terms. A leading-eigenvalue approximation avoids above problems, so that it works well with both Davies's (1980) and Farebrother's (1984) methods.

The saddlepoint approximation ends up with highly accurate approximation results for very small p-values. We show it has the correct exponential rate in the extreme right tail, so the relative error is bounded as $q \to \infty$, for any $A$ and $\Sigma$. In our numerical examples, the maximum logarithm of error ratio is less than 0.07. Therefore, a saddlepoint approximation and a leading-eigenvalue approximation combine to be usable for all p-values and all large enough numbers of variables.

For large quadratic forms, a leading-eigenvalue approximation provides a computational advantage without any important loss in accuracy and convolutions of the leading eigenvalues can be approximated by either the exact methods or the saddlepoint approximation.

## ACKNOWLEDGMENTS

## Appendix A. Exponential tail rate of the saddlepoint approximation

**Theorem 1.** *The saddlepoint approximation has the correct exponential rate in the extreme right tail.*

PROOF. One form of saddlepoint approximation defined in Equation (3) of Kuonen (1999) can be expressed in terms of error function

$$S = Pr(Q(X) > q) = 1 - \Phi\left\{w + \frac{1}{w}\log\left(\frac{v}{w}\right)\right\} = \frac{1}{2} - \frac{1}{2}\mathrm{erf}\left(\frac{x}{\sqrt{2}}\right), \tag{A.1}$$

where $x = w + (1/w)\log\left(v/w\right)$, $w = \mathrm{sign}(\hat{\zeta})[2\{\hat{\zeta}q - K(\hat{\zeta})\}]^{\frac{1}{2}}$, $v = \hat{\zeta}\{K''(\hat{\zeta})\}^{\frac{1}{2}}$, $K(\zeta)$ is the cumulant generating function of $Q(X)$ and $\hat{\zeta}$ is the saddlepoint. When $x \gg 1$, the asymptotic form of error function can be expanded as (Decker, 1975)

$$\mathrm{erf}(x) = 1 - \frac{e^{-x^2}}{\sqrt{\pi}}\sum_{m=0}^{\infty}\frac{(-1)^m(2m-1)!!}{2^m}x^{-(2m+1)},$$

9

where $(2m-1)!!$ is the product of all odd numbers up to $2m-1$.

Retain the first term ($m = 0$) in above summation and then plug it into Equation (A.1). Using Theorem 3.1 of Berman (1992), the exponential tail rate then becomes

$$-\frac{\mathrm{d}\log S}{\mathrm{d}q} = \left(w + \frac{1}{w}\log(\frac{v}{w}) + \frac{1}{w + \frac{1}{w}\log(\frac{v}{w})}\right)\frac{\mathrm{d}x}{\mathrm{d}q}. \tag{A.2}$$

To get $w$ and $v$, $K(\zeta)$ and its derivatives should be deduced. $Q(X)$ is a linear combination of central $\chi_1^2$ variables, so that $K(\zeta) = -\frac{1}{2}\sum_{i=1}^{n}\log(1-2\zeta\lambda_i)$, where $\lambda_1 \ldots \lambda_n$ are the non-zero eigenvalues and $\zeta < \frac{1}{2}\min 1/\lambda_i$. As the saddlepoint is the value of $\zeta$ satisfying $K'(\hat{\zeta}) = q$, it can be simplified to

$$K'(\hat{\zeta}) = \sum_{i=1}^{r}\frac{\lambda_i}{1 - 2\hat{\zeta}\lambda_i} = q. \tag{A.3}$$

Equation (A.3) shows that as $q$ tends to infinity, $\zeta$ tends towards $1/2\lambda_1$, but it will be always less than $1/2\lambda_1$, where $\lambda_1$ is the largest eigenvalue. In above summation, the largest term is $\lambda_1/(1-2\hat{\zeta}\lambda_1)$, so that $\hat{\zeta}$ can be approximated by $(q-\lambda_1)/2\lambda_1 q$. Then the asymptotic expression of $w$ and $v$ can be written as $((q-\lambda_1)/\lambda_1)^{\frac{1}{2}}$ and $(q-\lambda_1)/\sqrt{2}\lambda_1$. As $q \to \infty$, $w$ and $v$ tend towards infinity as well. Plugging $\mathrm{d}x/\mathrm{d}q$, $w$ and $v$ into Equation (A.2), the tail rate can be expressed as

$$-\frac{\mathrm{d}\log S}{\mathrm{d}q} \approx \left(w + \frac{1}{w}\log(\frac{v}{w}) + \frac{1}{w + \frac{1}{w}\log(\frac{v}{w})}\right)\left(\frac{1}{2\lambda_1 w} - \frac{1}{2\lambda_1 w^3}\log(\frac{v}{w}) - \frac{1}{2\lambda_1 w^3} + \frac{1}{\sqrt{2}\lambda_1 vw}\right)$$

$$\approx w\frac{1}{2\lambda_1 w} = \frac{1}{2\lambda_1}.$$

## References

Bausch, J., 2013. On the efficient calculation of a linear combination of chi-square random variables with an application in counting string vacua. Journal of Physics A: Mathematical and Theoretical 46 (50), 505202.

Berman, S. M., 1992. The tail of the convolution of densities and its application to a model of HIV-latency time. The Annals of Applied Probability, 481–502.

Chen, G. K., Marjoram, P., Wall, J. D., 2009. Fast and flexible simulation of DNA sequence data. Genome Research 19 (1), 136–142.

Davies, R. B., 1980. Algorithm AS 155: The distribution of a linear combination of $\chi^2$ random variables. Journal of the Royal Statistical Society. Series C (Applied Statistics) 29 (3), 323–333.

Decker, D. L., 1975. Computer evaluation of the complementary error function. American Journal of Physics 43, 833–834.

Duchesne, P., De Micheaux, P. L., 2010. Computing the distribution of quadratic forms: Further comparisons between the Liu–Tang–Zhang approximation and exact methods. Computational Statistics & Data Analysis 54 (4), 858–862.

Farebrother, R., 1984. Algorithm AS 204: The distribution of a positive linear combination of $\chi^2$ random variables. Journal of the Royal Statistical Society. Series C (Applied Statistics) 33 (3), 332–339.

Golub, G. H., Van Loan, C. F., 2012. Matrix Computations. JHU Press.

Halko, N., Martinsson, P.-G., Tropp, J. A., 2011. Finding structure with randomness: Probabilistic algorithms for constructing approximate matrix decompositions. SIAM review 53 (2), 217–288.

Kuonen, D., 1999. Miscellanea. Saddlepoint approximations for distributions of quadratic forms in normal variables. Biometrika 86 (4), 929–935.

Lee, S., with contributions from Larisa Miropolsky, Wu, M., 2017. SKAT: SNP-Set (Sequence) Kernel Association Test. R package version 1.3.2.1.
URL https://CRAN.R-project.org/package=SKAT

Liu, H., Tang, Y., 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 (4), 853–856.

Lumley, T., 2011. Complex Surveys: A Guide to Analysis Using R. John Wiley & Sons.

Lumley, T., 2019. bigQF: Quadratic Forms in Large Matrices. R package version 1.3-3.
URL https://github.com/tslumley/bigQF

Lumley, T. S., Brody, J. A., Peloso, G. M., Morrison, A. C., Rice, K. M., 2018. FastSKAT: Sequence kernel association tests for very large sets of markers. Genetic Epidemiology.

R Core Team, 2017. R: A Language and Environment for Statistical Computing. R Foundation for Statistical Computing, Vienna, Austria.
URL https://www.R-project.org/

Rao, J. N., Scott, A. J., 1981. The analysis of categorical data from complex sample surveys: chi-squared tests for goodness of fit and independence in two-way tables. Journal of the American Statistical Association 76 (374), 221–230.

Robbins, H., Pitman, E., 1949. Application of the method of mixtures to quadratic forms in normal variates. The Annals of Mathematical Statistics, 552–560.

Wu, M. C., Lee, S., Cai, T., Li, Y., Boehnke, M., Lin, X., 2011. Rare-variant association testing for sequencing data with the sequence kernel association test. The American Journal of Human Genetics 89 (1), 82–93.