

Scalable multi-component linear mixed models with application to SNP heritability estimation

Ali Pazokitoroudi¹, Yue Wu¹, Kathryn S. Burch², Kangcheng Hou³, Bogdan Pasaniuc^{4,5,6}, and Sriram Sankararaman^{*1,5,6}

¹Department of Computer Science, UCLA, Los Angeles, California

²Bioinformatics Interdepartmental Program, UCLA, Los Angeles, California

³Department of Computer Science, Zhejiang University, Hangzhou, Zhejiang, China

⁴Department of Pathology and Laboratory Medicine, David Geffen School of Medicine, UCLA, Los Angeles, California

⁵Department of Human Genetics, David Geffen School of Medicine, UCLA, Los Angeles, California

⁶Department of Computational Medicine, David Geffen School of Medicine, UCLA, Los Angeles, California

Abstract

A central question in human genetics is to find the proportion of variation in a trait that can be explained by genetic variation. A number of methods have been developed to estimate this quantity, termed narrow-sense heritability, from genome-wide SNP data. Recently, it has become clear that estimates of narrow-sense heritability are sensitive to modeling assumptions that relate the effect sizes of a SNP to its minor allele frequency (MAF) and linkage disequilibrium (LD) patterns [3]. A principled approach to estimate heritability while accounting for variation in SNP effect sizes involves the application of linear Mixed Models (LMMs) with multiple variance components where each variance component represents the fraction of genetic variance explained by SNPs that belong to a given range of MAF and LD values. Beyond their importance in accurately estimating genome-wide SNP heritability, multiple variance component LMMs are useful in partitioning the contribution of genomic annotations to trait heritability which, in turn, can provide insights into biological processes that are associated with the trait.

Existing methods for fitting multi-component LMMs rely on maximizing the likelihood of the variance components. These methods pose major computational bottlenecks that makes it challenging to apply them to large-scale genomic datasets such as the UK Biobank which contains half a million individuals genotyped at tens of millions of SNPs.

We propose a scalable algorithm, RHE-reg-mc, to jointly estimate multiple variance components in LMMs. Our algorithm is a randomized method-of-moments estimator that has a runtime that is observed to scale as $\mathcal{O}\left(\frac{NMB}{\max(\log_3(N), \log_3(M))} + K^3\right)$ for N individuals, M SNPs, K variance components, and $B \approx 10$ being a parameter that controls the number of random matrix-vector multiplication. RHE-reg-mc also efficiently computes standard errors. We evaluate the accuracy and scalability of RHE-reg-mc for estimating the total heritability as well as in partitioning heritability. The ability to fit multiple variance components to SNPs partitioned according to their MAF and local LD allows RHE-reg-mc to obtain relatively unbiased estimates of SNP heritability under a wide range of models of genetic architecture. On the UK Biobank dataset consisting of $\approx 300,000$ individuals and $\approx 500,000$ SNPs, RHE-reg-mc can fit 250 variance components, corresponding to genetic variance explained by 1 MB blocks, in ≈ 40 minutes on standard hardware.

1 Introduction

Heritability is a central parameter in understanding the contribution of genetic variation to trait variation [22]. Narrow-sense heritability refers to the maximal proportion of variation in a trait that can be

explained by a linear function of genetic variation[22]. Over the last decade, there has been substantial attention focused on estimating narrow-sense heritability from genome-wide SNP genotype data [25]. These SNP heritability estimates are of great interest in understanding the genetic basis of complex traits and can inform strategies for designing future genetic studies. While a number of methods have been developed to estimate SNP heritability [25, 27, 11, 5, 19], recent studies have shown that the estimates of SNP heritability from these methods can be highly sensitive to modeling assumptions such as the joint distribution of the effect sizes at causal variants, their allele frequencies as well as the levels of correlation or linkage disequilibrium (LD) between causal variants and the genotyped SNPs[19, 3]. An observation from these studies is that methods that assume that the effect size of each SNP on a trait comes from the same distribution can yield biased estimates of SNP heritability [3]. One potential strategy to relax the assumption of identically distributed effect sizes consists of parametrizing the distribution of SNP effect sizes in terms of relevant covariates that are expected to influence their distributions. For example, methods that partition SNPs into bins based on minor allele frequency (MAF) and local LD patterns thereby allowing the distribution of effect sizes to vary as a function of MAF and LD tend to be robust to variation in the underlying genetic models [3]. Such partitioning strategies are also motivated by a number of recent studies suggest that the distribution of effect sizes of SNPs on a trait can have a complex relationship with their allele frequencies and LD [3, 4]. Beyond the dependency of SNP effect sizes on MAF and LD, the SNP effect sizes have been observed to vary as a function of a number of genomic annotations such as whether a SNP lies in the protein coding regions or in regions of open chromatin [7, 6, 4]. These observations, together, motivate statistical models that allow for the flexible modeling of genetic effects on phenotype where the genetic effects vary as a function of covariates such as MAF, LD or other genomic annotations.

Linear mixed models (LMMs) are an important class of models that permit the representation of flexible relationships between genetic variation and complex traits, *i.e.*, traits that are modulated by multiple genetic and environmental factors. LMMs have been applied successfully in linkage analysis in family studies, in genomic selection and risk prediction, as well as in association analysis to control for individual relatedness and population stratification.

LMMs [15] have been applied to estimate SNP heritability attributed to genome-wide SNPs (h_{SNP}^2) [25]. In the simplest setting, the LMM is endowed with two parameters, *i.e.*, variance components, corresponding to the phenotypic variance explained by genetic and residual factors respectively. These models, termed *single-component* LMMs, as they employ a single variance component to capture all the genetic effects, assume that the effect of each SNP on the phenotype is drawn independently from the same underlying distribution. However, when this assumption is violated, these single component LMMs are expected to yield biased estimates of SNP heritability. To overcome this limitation, multi-components LMMs have been proposed which assign SNPs into one of several variance components and assume that the effect sizes at SNPs assigned to a given variance component are drawn independently from the same distribution [24]. By binning SNPs according to their MAF and LD and assigning a variance component to each bin, these multi-components LMMs have been shown to yield heritability estimates that are accurate across a range of underlying genetic architectures.

While multi-component LMMs appear to be well-suited to estimate genome-wide SNP heritability, estimating their parameters, *i.e.*, the variance components, is computationally demanding. The most common approach to estimate the variance components is to search for parameter values that maximize the likelihood. Usually a particular form of maximum likelihood, the restricted maximum likelihood (REML) estimator [16], is preferred due to a reduced bias relative to the full maximum likelihood estimator. Computing maximum likelihood or REML estimators, however, can be challenging. Methods for computing maximum likelihood or REML rely on iterative optimization algorithms. These algorithms do not scale well to large data sets like the UK biobank which contains $\approx 100,000$ individuals and over a million SNPs. While a number of algorithmic approaches have been proposed for efficient inference in LMMs, many of these are designed to leverage the specific structure of single-component LMMs and cannot be applied to the multi-component setting [13, 28]. Thus, efficient parameter estimation in multi-component LMMs remains a challenging computational problem.

1.1 Our contribution

We propose a fast multi-component variance components estimation algorithm for linear mixed models with many variance components based on a randomized method-of-moments (MoM) estimator. Our estimator can be viewed as a generalization of the classic Haseman-Elston (HE) regression [8] estimator to the multi-component setting as well as of a recently proposed randomized version of HE regression for the single component setting[23].

Being a Method-of-Moments estimator, our proposed estimator is statistically less efficient than REML but it is computationally attractive. Recently, a scalable estimator of variance components for a single-component LMM based on a randomized version of HE regression (RHE-reg), was proposed. This method has a runtime complexity $\mathcal{O}\left(\frac{NMB}{\max(\log_3(N), \log_3(M))}\right)$ for N individuals and M SNPs and a parameter B that controls the number of random matrix-vector multiplications[23]. In this paper, we extend RHE-reg to the multiple component setting where we assume that SNPs have been classified into one of K different non-overlapping functional categories. Our method, RHE-reg-mc, estimates the variance component of each category in time $\mathcal{O}\left(\frac{NMB}{\max(\log_3(N), \log_3(M))} + K^3\right)$

The time complexity of RHE-reg-mc scales with the size of the genotype matrix as long as the number of components or partitions K is less than $(MN)^{1/3}$ where M and N are the number of SNPs and individuals respectively. For example, in the full set of UK Biobank genotypes ($M \approx 10^6$, $N \approx 10^5$), this allows us to define 10^4 variance components while maintaining scalability.

We apply RHE-reg-mc to the problem of estimating genome-wide SNP heritability when the underlying genetic architecture assumes that the effect sizes are a function of the MAF and LD patterns at the SNP. To account for variation in the SNP effect sizes, we bin SNPs according to their MAF and LD patterns (following previous approaches [26] and assign a distinct variance component for each bin. In small scale experiments containing around 14,000 individuals, we show that RHE-reg-mc yields approximately unbiased estimates of SNP heritability albeit with larger standard errors relative to multi-component REML methods. RHE-reg-mc remained relatively unbiased in large-scale experiments on about 337,000 individuals. Importantly, RHE-reg-mc was the only method that used individual-level genotypes that could run on these datasets. Our bench-marking experiments show that RHE-reg-mc is around 400 times faster than other state-of-the-art methods such as BOLT-REML on a dataset of 100,000 individuals and 500,000 SNPs even when the latter were required to estimate only a single component. Therefore, RHE-reg-mc is the only method that used individual-level genotypes that can run efficiently on large datasets such as full UK Biobank which contains half a million individuals and millions of SNPs. Beyond estimating SNP heritability, we also show that RHE-reg-mc accurately partitions heritability across 1 Mb regions which corresponds to jointly estimating 250 variance components.

2 Methods

2.1 Motivation

Various methods have been developed to estimate narrow sense heritability [1]. According to recent studies [18, 3], our assumptions about the levels of LD, MAF and effect sizes of underlying causal variants have significant effects on the accuracy of estimates of narrow sense heritability. However, methods which partition SNPs into bins based on minor allele frequency (MAF) and linkage disequilibrium (LD) [24] obtain reduced bias in their estimates of heritability.

To understand why the estimation of narrow-sense heritability is sensitive to assumptions about the underlying causal variants, and why single component LMMs are expected to yield biased estimates of SNP heritability, consider the following single-component linear mixed model:

$$\begin{aligned} \mathbf{y} | \boldsymbol{\epsilon}, \boldsymbol{\beta} &= \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \\ \boldsymbol{\epsilon} &\sim \mathcal{N}(\mathbf{0}, \sigma_{\epsilon}^2 \mathbf{I}_N) \\ \boldsymbol{\beta} &\sim \mathcal{N}\left(\mathbf{0}, \frac{\sigma_g^2}{M} \mathbf{I}_M\right) \end{aligned} \tag{1}$$

Here \mathbf{y} is a N -vector of phenotypes which is centered, and \mathbf{X} is a $N \times M$ matrix of standardized genotypes obtained by centering and scaling each column of the genotype matrix \mathbf{G} where $g_{i,j} \in \{0, 1, 2\}$ denotes the number of minor alleles carried by individual i at SNP j , N and M are the number of individuals and SNPs respectively. $\boldsymbol{\beta}$ is a M -vector of SNP effect sizes. σ_e^2 is the residual variance while σ_g^2 is the variance component corresponding to the M SNPs. In this model, the SNP heritability is defined as $h_{SNP}^2 = \frac{\sigma_g^2}{\sigma_g^2 + \sigma_e^2}$. This model assumes that the variance of the effect size is the same across all SNPs. When effect sizes are coupled with MAF or LD, the mismatch between the model assumptions and data can lead to biased heritability estimates.

To tackle this problem, the LMM can be extended to include multiple components as follows:

$$\mathbf{y} | \boldsymbol{\epsilon}, \boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_K = \sum_{k=1}^K \mathbf{X}_k \boldsymbol{\beta}_k + \boldsymbol{\epsilon} \quad (2)$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \sigma_e^2 \mathbf{I}_N)$$

$$\boldsymbol{\beta}_k \sim \mathcal{N}(\mathbf{0}, \frac{\sigma_k^2}{M_k} \mathbf{I}_{M_k}), k \in \{1, \dots, K\} \quad (3)$$

Here each of the M SNPs is assigned to one of K non-overlapping categories. Each category k contains M_k SNPs, $k \in \{1, \dots, K\}$, $\sum_k M_k = M$. Let \mathbf{G}_k denote the $N \times M_k$ genotype matrix for category k such that $g_{k,i,j} \in \{0, 1, 2\}$ denotes the number of minor alleles carried by individual i at SNP j in category k . Let \mathbf{X}_k be a $N \times M_k$ matrix of standardized genotypes obtained by centering and scaling each column of \mathbf{G}_k so that $\sum_i g_{k,i,m} = 0$ and $\sum_k g_{k,i,m}^2 = N$ for $m \in \{1, 2, \dots, M\}$. Let $\boldsymbol{\beta}_k$ be a M_k -vector of SNP effect sizes for the k -th category. In the above model, σ_e^2 is the residual variance, and σ_k^2 is the variance component of the k -th category. In this model, the total SNP heritability is defined as :

$$h_{SNP}^2 = \frac{\sum_{k=1}^K \sigma_k^2}{(\sum_{k=1}^K \sigma_k^2) + \sigma_e^2} \quad (4)$$

The SNP heritability of category k is defined as:

$$h_k^2 = \frac{\sigma_k^2}{(\sum_{k=1}^K \sigma_k^2) + \sigma_e^2}, k \in \{1, \dots, K\} \quad (5)$$

The model in Equation 3 has K partitions such that the variance of effect sizes can differ among the partitions. We can partition SNPs into bins based on minor allele frequency (MAF) and local LD patterns thereby allowing the distribution of effect sizes to vary as a function of MAF and LD tend to be robust to variation in the underlying genetic models and achieve unbiased estimation of SNP heritability. However, in this setting we have several challenges. First, can we estimate the parameters of the model efficiently? Second, for a given genotype, what is the minimum value of K required for an accurate estimation of heritability?

2.2 Method-of-moments for estimating multiple variance components

To estimate the variance components of the multi-component LMM, we use a Method-of-Moments (MoM) estimator that searches for parameter values so that the population moments are close to the sample moments. Since $\mathbb{E}[\mathbf{y}] = 0$, we derived the MoM estimates by equating the population covariance to the empirical covariance. The population covariance is given by:

$$\text{cov}(\mathbf{y}) = E[\mathbf{y}\mathbf{y}^T] - E[\mathbf{y}]E[\mathbf{y}^T] = \sum_k \sigma_k^2 \mathbf{K}_k + \sigma_e^2 \mathbf{I}_N \quad (6)$$

Here $\mathbf{K}_k = \frac{\mathbf{X}_k \mathbf{X}_k^T}{M_k}$ is the genetic relatedness matrix (GRM) computed from all SNPs of k -th category. Using $\mathbf{y}\mathbf{y}^T$ as our estimate of the empirical covariance, we need to solve the following least squares problem to find the variance components.

$$(\tilde{\sigma}_1^2, \dots, \tilde{\sigma}_K^2, \tilde{\sigma}_e^2) = \underset{(\sigma_1^2, \dots, \sigma_K^2, \sigma_e^2)}{\operatorname{argmin}} \|\mathbf{y}\mathbf{y}^T - \sum_k \sigma_k^2 \mathbf{K}_k + \sigma_e^2 \mathbf{I}\|_F^2 \quad (7)$$

It is not hard to see that the MoM estimator satisfies the following normal equations:

$$\begin{bmatrix} \mathbf{T} & \mathbf{b} \\ \mathbf{b}^T & N \end{bmatrix} \begin{bmatrix} \tilde{\sigma}_1^2 \\ \vdots \\ \tilde{\sigma}_k^2 \\ \vdots \\ \tilde{\sigma}_e^2 \end{bmatrix} = \begin{bmatrix} \mathbf{c} \\ \mathbf{y}^T \mathbf{y} \end{bmatrix} \quad (8)$$

Here \mathbf{T} is a $K \times K$ matrix with entries $T_{k,l} = \operatorname{tr}(\mathbf{K}_k \mathbf{K}_l)$, $k, l \in \{1, \dots, K\}$, \mathbf{b} is a K -vector with entries $b_k = \operatorname{tr}(\mathbf{K}_k) = N$ (because \mathbf{X}_k s is standardized), and \mathbf{c} is a K -vector with entries $c_k = \mathbf{y}^T \mathbf{K}_k \mathbf{y}$.

Every GRM \mathbf{K}_k can be computed in time $\mathcal{O}(N^2 M_k)$ and $\mathcal{O}(N^2)$ memory. Given K GRMs, the quantities $T_{k,l}$, c_k , $k, l \in \{1, \dots, K\}$, can be computed in $\mathcal{O}(NM)$. Given the quantities $T_{k,l}$, c_k , the normal equation 8 can be solved in $\mathcal{O}(K^3)$. Therefore, the total time complexity for estimating the variance components is $\mathcal{O}(N^2 M + K^3)$.

2.2.1 RHE-reg-mc: Randomized estimator of multiple variance components

The key bottleneck in solving the normal equation 8 is the computation of $T_{k,l}$, $k, l \in \{1, \dots, K\}$ which takes $\mathcal{O}(N^2 M)$. Instead of computing the exact value of $T_{k,l}$, we use Hutchinson's estimator of the trace [9]. This estimator uses the fact that for a given $N \times N$ matrix \mathbf{C} , $\mathbf{z}^T \mathbf{C} \mathbf{z}$ is an unbiased estimator of $\operatorname{tr}(\mathbf{C})$ ($E[\mathbf{z}^T \mathbf{C} \mathbf{z}] = \operatorname{tr}[\mathbf{C}]$) where \mathbf{z} be a random vector with mean zero and covariance \mathbf{I}_N . Hence, we can estimate the values $T_{k,l}$, $k, l \in \{1, \dots, K\}$ as follows:

$$T_{k,l} = \operatorname{tr}(\mathbf{K}_k \mathbf{K}_l) \approx \widehat{T}_{k,l} = \frac{1}{B} \frac{1}{M_k M_l} \sum_b \mathbf{z}_b^T \mathbf{X}_k \mathbf{X}_k^T \mathbf{X}_l \mathbf{X}_l^T \mathbf{z}_b \quad (9)$$

Here $\mathbf{z}_1, \dots, \mathbf{z}_B$ are B independent random vectors with zero mean and covariance \mathbf{I}_N . In our method, we draw these random vectors independently from a standard normal distribution. Note that computing $T_{k,l}$ by using the unbiased estimator involves four matrix-vector multiplications which is repeated B times. Therefore, the total running time for estimating the values $T_{k,l}$ is $\mathcal{O}(NMB)$.

Moreover, we can leverage the structure of the genotype matrix which only contains entries in $\{0, 1, 2\}$. For a fixed genotype matrix \mathbf{X}_k , we can improve the per iteration time complexity of matrix-vector multiplication from $\mathcal{O}(NM)$ to $\mathcal{O}(\frac{NM}{\max(\log_3(N), \log_3(M))})$ by using the Mailman algorithm [12]. Solving the normal equations takes $\mathcal{O}(K^3)$ time so that the overall time complexity of our algorithm is $\mathcal{O}(\frac{NMB}{\max(\log_3(N), \log_3(M))} + K^3)$.

2.3 Computing the Standard Errors of the estimates

We obtain standard errors for RHE-reg-mc using a block jackknife [10]. The jackknife is a useful resampling technique for estimating the variance and bias of an estimator [21]. A jackknife subsample is created by leaving out a subset of observations from a dataset. The jackknife estimate of a parameter can be found by estimating the parameter for each subsample omitting the i -th jackknife block. A naive way to compute jackknife estimate requires computing the estimator of the parameters for every subsample. For instance, in our problem, if we define J jackknife blocks, then we need to run RHE-reg-mc for every subsample which takes $\mathcal{O}(J(\frac{NMB}{\max(\log_3(N), \log_3(M))} + K^3))$. We propose an efficient way to compute the jackknife estimate in time $\mathcal{O}(\frac{NMB}{\max(\log_3(N), \log_3(M))} + JK^3)$.

Let \mathbf{X} be a $N \times M$ matrix of standardized genotypes where N and M are the number of individuals and SNPs respectively. To generate J jackknife subsamples, we partition \mathbf{X} into J non-overlapping blocks $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(J)}$ such that $\mathbf{X} = [\mathbf{X}^{(1)}, \mathbf{X}^{(2)}, \dots, \mathbf{X}^{(J)}]$. Note that for every j , $\mathbf{X}^{(j)}$ is a $N \times M_j$ matrix where M_j is the number of SNPs in the j -th block.

We create the j -th jackknife subsample by removing the j -th block $\mathbf{X}^{(j)}$ from \mathbf{X} . To estimate the variance components of the j -th jackknife subsample, we need to compute the corresponding quantities of

the j th subsample in the normal equations 8. Let $\mathbf{K}_k^{(-j)}$ be the GRM of the k -th partition which is created by removing the j -th block $\mathbf{X}^{(j)}$ from \mathbf{X} where $k \in \{1, \dots, K\}$, $j \in \{1, \dots, J\}$. In Appendix A we show that we can compute $tr(\widehat{\mathbf{K}_k^{(-j)} \mathbf{K}_l^{(-j)}})$ and $\mathbf{y}^T \mathbf{K}_i^{(-j)} \mathbf{y}$, for all $k, l \in \{1, \dots, K\}$, $j \in \{1, \dots, J\}$, in time $\mathcal{O}(\frac{NMB}{\max(\log_3(N), \log_3(M))})$ which does not effect the total running time of the algorithm.

Therefore, for every jackknife subsample, we can estimate the corresponding variance components in $\mathcal{O}(K^3)$ by solving the corresponding normal equations 8. Given the estimates of variance components for each jackknife subsamples, we can compute jackknife estimate of the variance, estimate the bias, and bias-corrected jackknife estimate of the variance components.

2.4 Including covariates

We can extend our model in Equation 2 to include covariates as follows:

$$\mathbf{y} | \epsilon, \beta_1, \dots, \beta_k = \mathbf{W}\alpha + \sum_k \mathbf{X}_k \beta_k + \epsilon \quad (10)$$

Here \mathbf{W} is a $N \times C$ matrix of covariates while α is a C -vector of fixed effects. In Appendix B, we show that in this setting, we need to solve the following normal equations to estimate the variance components.

$$\begin{bmatrix} \mathbf{T} & \mathbf{b} \\ \mathbf{b}^T & N - C \end{bmatrix} \begin{bmatrix} \sigma_1^2 \\ \vdots \\ \sigma_k^2 \\ \sigma_e^2 \end{bmatrix} = \begin{bmatrix} \mathbf{c} \\ \mathbf{y}^T \mathbf{V} \mathbf{y} \end{bmatrix} \quad (11)$$

Here $\mathbf{V} = \mathbf{I}_N - \mathbf{W}(\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T$ and \mathbf{T} is a $K \times K$ matrix where $T_{k,l} = tr(\mathbf{K}_k \mathbf{V} \mathbf{K}_l \mathbf{V})$, and \mathbf{b} is a K -vector where $b_k = tr(\mathbf{V} \mathbf{K}_k)$, and \mathbf{c} is a K -vector where $c_k = \mathbf{y}^T \mathbf{V} \mathbf{K}_k \mathbf{V} \mathbf{y}$. Commonly, the number of covariates C is small (tens to hundreds) so that including covariates does not significantly affect the computational cost. The cost of computing of the elements of the normal equations 11 includes the cost of inverting $\mathbf{W}^T \mathbf{W}$ which is a $C \times C$ matrix and multiplying \mathbf{W} by a real-valued N -vector which can be done in $\mathcal{O}(C^3 + NC)$.

3 Results

3.1 Estimating total heritability by partitioning based on MAF and LD

3.1.1 Simulations

We performed simulations to compare the performance of RHE-reg-mc with several state-of-the-art methods for heritability estimation that cover the spectrum of methods that have been proposed. These methods include single-component as well as multi-component LMMs which use likelihood maximization for parameter estimation. These methods require access to individual-level genotypes and phenotypes. We also compared to methods that only require access to summary statistics and are typically quite scalable. GCTA and BOLT-REML estimate heritability by maximizing the restricted maximum likelihood (REML). BOLT-REML is a computationally efficient approximate method to compute the REML estimator. GCTA-ldms is the extension of GCTA to a multi-component LMM where the variance components are typically defined by binning SNPs according to their MAF as well as local LD. LDAK is similar to GCTA, except that it assumes allelic effects are a function of LD scores. Among the summary statistic methods, LD score regression (LDSC) uses the slope from the GWAS χ^2 statistics regressed on the LD scores to estimate the h_{SNP}^2 . Stratified LD score method (S-LDSC) is an extension of LDSC for partitioning heritability from summary statistics. SumHer is the summary statistic analog of LDAK [17].

We considered two simulation settings. The small-scale setting was designed so that we could compare the accuracies of our method to state-of-the-art methods on the same data. In this setting, we simulated phenotypes from a subsampled set of genotypes from the UK Biobank [20]. Specifically, we chose $M = 14,000$ SNPs from chromosome 1 of the UK Biobank Axiom array and a subset of $N = 9,000$ individuals. In the

large-scale simulation setting, we simulated phenotypes for the full set of UK Biobank genotypes consisting of $M = 590,000$ array SNPs and $N = 337,000$ unrelated white British individuals.

We simulated phenotypes from genotypes using the following model for the genetic architecture:

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \\ \boldsymbol{\epsilon} | \sigma_e^2 &\sim \mathcal{N}(\mathbf{0}, \sigma_e^2 \mathbf{I}_N) \\ \boldsymbol{\beta}_m &\sim \mathcal{N}(\mathbf{0}, c[f_m(1 - f_m)]^a w_m^b), m \in \{1, \dots, M\} \end{aligned} \quad (12)$$

where $a \in \{0, 0.75\}$, $b \in \{0, 1\}$, c is a constant, β_m , f_m and w_m are the effect size, the minor allele frequency and LDK score of m^{th} SNP respectively. The LD score of a SNP is defined to be the sum of the squared correlation of the SNP with all other SNPs that lie within a specific distance, and the LDK score of a SNP is computed based on local levels of LD such that the LDK score tends to be higher for SNPs in regions of low LD[19]. The above models relating genotype to phenotype are commonly used in methods for estimating SNP heritability: the GCTA Model (when $a = b = 0$ in Equation 12), which is used by the software GCTA [26] and LD Score regression (LDSC) [2], and the LDK Model (where $a = 0.75, b = 1$ in Equation 12) used by software LDK [19]. Moreover, under each model, we varied the proportion and minor allele frequency (MAF) of causal variants (CVs). Proportion of causal variants set to be either 100% or 1%, and MAF of causal variants drawn uniformly from $[0, 1]$ or $[0.01, 0.05]$ to consider genetic architectures that are either infinitesimal or sparse as well genetic architectures that include a mixture of common and rare SNPs as well as one that includes only common SNPs. The GCTA Model assumes that heritability is independent of LD, while the LDK Model assumes that heritability varies according to local levels of LD[18].

We generated 100 sets of simulated phenotypes for each setting of parameters and report accuracies averaged over these 100 sets.

3.1.2 Accuracy

We compared the accuracy of RHE-reg-mc with the most popular single component methods such as BOLT-REML [14], GCTA [26], LDK [19], LDSC [2], SumHer [17], as well as multi-components methods such as GCTA-ldms and S-LDSC (BOLT-REML can also be estimate multi-component variance components but we did not explore this option). In the small scale setting, we compared RHE-reg-mc with BOLT-REML, GCTA, GCTA-ldms, and LDK, each of which use individual-level phenotypes and genotypes[19, 2, 26, 17]. In the large scale setting, we only compared RHE-reg-mc with methods that use summary statistics like SumHer, LDSC, and S-LDSC because the existing methods based on individual-level phenotypes and genotype are not scalable to large data sets. For the multi-component methods, we applied GCTA-ldms by binning SNPs into 8 bins based on 2 bins for MAF ($\text{MAF} < 0.05, \text{MAF} > 0.05$), MAF refers to the frequency at which the second most common allele occurs in a SNP, and 4 bins based on quartiles of the LD score of a SNP. The LD score of a SNP is defined to be the sum of the squared correlation of the SNP with all other SNPs that lie within a specific distance. We used the default LD score computation that is used by GCTA. For RHE-reg-mc, we used 16 bins formed by the combination of 4 bins based on MAF ($\text{MAF} < 0.009, 0.009 < \text{MAF} < 0.011, 0.011 < \text{MAF} < 0.05, 0.05 < \text{MAF} < 0.5$) as well as 4 bins based on quartiles of the LDK score of a SNP. The LDK score of a SNP is computed based on local levels of LD such that the LDK score tends to be higher for SNPs in regions of low LD[19].

Figure 1 and Table 1 show that both multi-component methods, *i.e.*, GCTA-ldms and RHE-reg-mc, have the least bias across the settings considered. RHE-reg-mc tends to have larger standard errors relative to GCTA-ldms consistent with the lower statistical efficiency of method-of-moments estimators relative to REML estimators. In the large-scale setting (Figure 2 and Table 2), RHE-reg-mc has reduced bias as well as low standard errors compared to other methods.

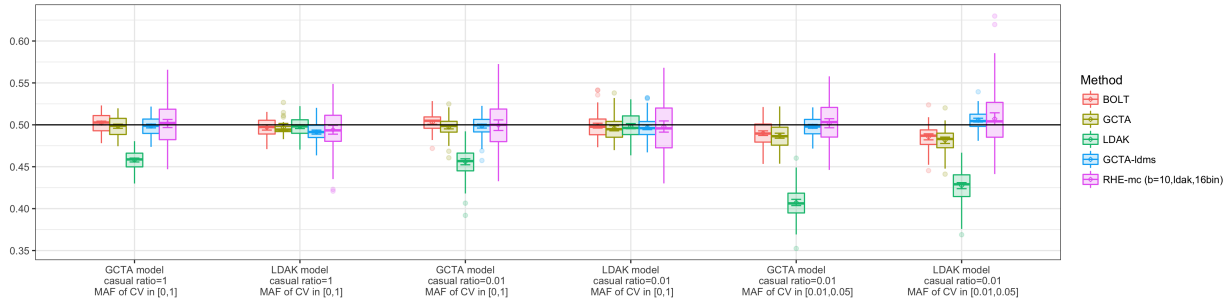


Figure 1: **RHE-reg-mc is relatively unbiased in small scale simulated data** ($M = 14,000$ array SNPs and $N = 9,000$ individuals): We ran 100 replicates where the true heritability of the phenotype is 0.5. We compared methods for heritability estimation under both GCTA and LDKA models with different proportions of causal variants as well as MAF distributions. BOLT-REML (run with single component), GCTA, and LDKA are single component methods. GCTA-ldms is a multi-component method that is applied by binning SNPs to 8 bins based on 2 MAF bins and 4 bins based on quartiles of LD scores. RHE-reg-mc is our proposed multi-component method that is applied by binning SNPs to 16 bins based on 4 MAF bins and 4 bins based on LDKA weights, and $b = 10$ random vector. GCTA-ldms and RHE-reg-mc are relatively unbiased.

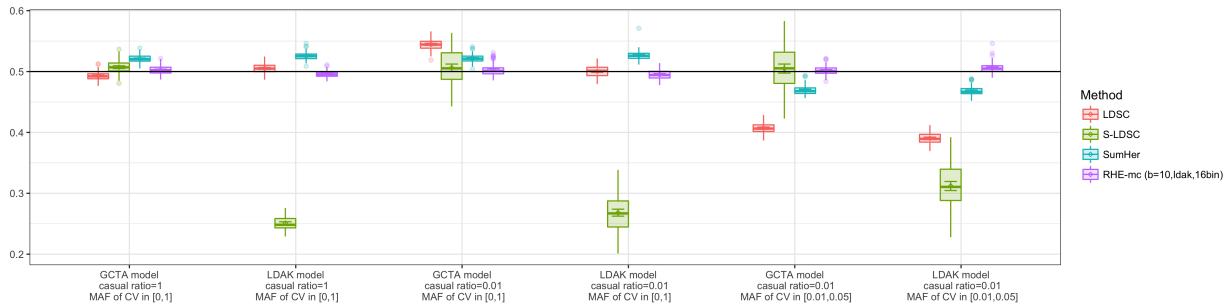


Figure 2: **RHE-reg-mc is relatively unbiased in large scale simulated data** ($M=590K$ array SNPs and $N=337K$ individuals): We ran 100 replicates where the true heritability of the phenotype is 0.5. We compared methods for heritability estimation under both GCTA and LDKA models with different proportions of causal variants and MAF. LDSC, S-LDSC and SumHer methods are based on summary statistics. RHE-reg-mc is an individual-level method that is applied by binning SNPs to 16 bins based on 4 MAF bins and 4 LDKA bins. RHE-reg-mc is relatively unbiased across the settings.

Method	ave(bias),GCTA model	ave(sd),GCTA model	ave(bias),LDKA model	ave(sd),LDKA model
RHE-reg-mc($b=10,16bin$)	0.001	0.027	0.005	0.031
GCTA-ldms	0.001	0.011	0.005	0.012
LDKA	0.059	0.015	0.025	0.014
GCTA	0.005	0.012	0.008	0.013
BOLT	0.004	0.011	0.006	0.013

Table 1: The average of standard errors and the average of absolute values of bias across all settings are reported for each method for small scale data.

Method	ave(bias),GCTA model	ave(sd),GCTA model	ave(bias),LDAK model	ave(sd),LDAK model
RHE-reg-mc(b=10,16bin)	0.002	0.007	0.004	0.007
SumHer	0.024	0.006	0.027	0.007
S-LDSC	0.006	0.026	0.222	0.02
LDSC	0.048	0.008	0.038	0.008

Table 2: The average of standard errors and the average of absolute values of bias across all settings are reported for each method for large scale data.

3.2 Computational Efficiency

We compared the running time and memory usage of RHE-reg-mc method with GCTA [26] and BOLT-REML [14]. In this comparison, we used the UK Biobank genotypes consisting of around 500,000 SNPs over different sample sizes. For each data set, we ran RHE-reg-mc with $B = 10$ random vectors and 22 bins (SNPs partitioned based on chromosomes). We ran GCTA and BOLT-REML for a single variance component. All computations were restricted to a single core on a standard compute machine.

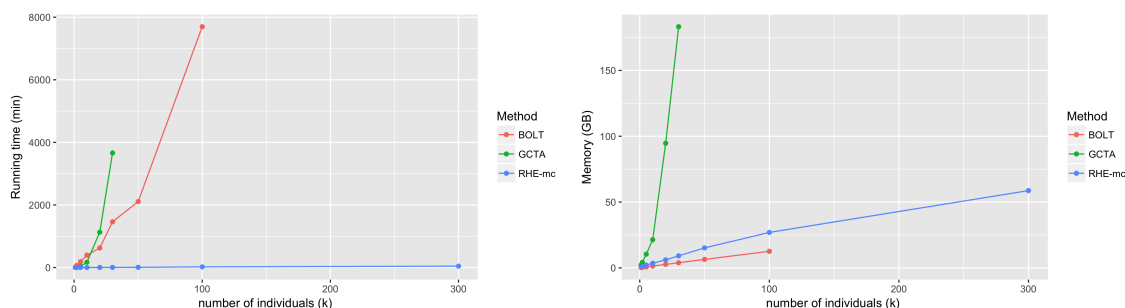


Figure 3: **RHE-reg-mc is efficient:** We compare the runtime and memory usage of methods for heritability estimation. We vary the number of samples while fixing the number of SNPs to 500,000. We run RHE-reg-mc with 22 bins defined based on chromosomes. We ran GCTA and BOLT-REML assuming a single variance component as multiple variance component versions tend to be more computationally intensive.

Figure 3 shows that we could not run single-component GCTA to sample sizes beyond 50,000 due to memory constraints. Single-component BOLT-REML took about 5 days to run on 100,000 individuals while the computation of RHE-reg-mc was about 20 minutes. Our bench-marking experiments show that RHE-reg-mc is around 400 times faster than other state-of-the-art methods such as BOLT-REML on a dataset of 100k individuals and 500k SNPs. Extrapolating this result, we expect that RHE-reg-mc could run on large datasets such as the full UK Biobank which contains half a million individuals genotyped at tens of millions of SNPs efficiently. The memory usage of RHE-reg-mc is linear with respect to sample size. The running time and accuracy of RHE-reg-mc relies on the choice of the number of random vectors B . In practice, it turns out that the estimator is highly accurate with a small $B \approx 10$ across all datasets analyzed.

3.3 Partitioning heritability

To examine the ability of RHE-reg-mc to partition the heritability explained across genomic regions, we partition SNPs into 1 MB regions so that the first 50 bins each explains 1% of the variance ($h_{SNP}^2 = 0.5$) while the rest of the bins have zero heritability. Figure 4 shows that RHE-reg-mc accurately estimates heritability in each partition on 100,000 individuals and 460,000 SNPs. Further, RHE-reg-mc computed the partitioned heritability in 20 minutes.

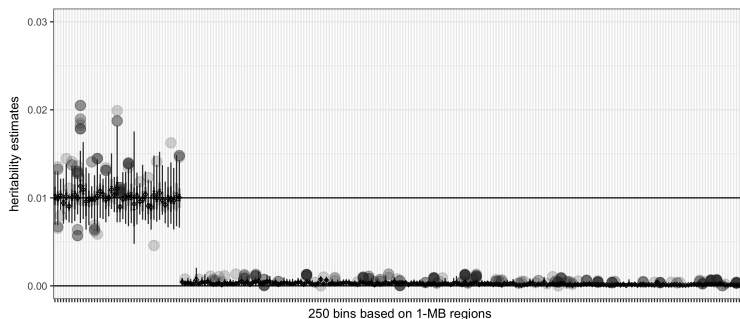


Figure 4: **RHE-reg-mc accurately partitions heritability:** We simulated phenotypes from 460,000 SNPs of UK Biobank genotype over 100,000 individuals. Each of the first 50 bins has a heritability of 1% ($h_{SNP}^2 = 0.5$) while the rest of the bins have zero heritability.

3.4 Application to phenotypes in the UK Biobank

Finally, we used RHE-reg-mc to partition the heritability of three phenotypes (trunk fat percentage, diastolic blood pressure, and systolic blood pressure) on the UK Biobank dataset consisting of around 500,000 SNPs and 300,000 individuals. Partitioning SNPs according to 22 chromosomes, we see that longer chromosomes length have higher heritability consistent with studies in other traits such as height [25] (Figure 5).

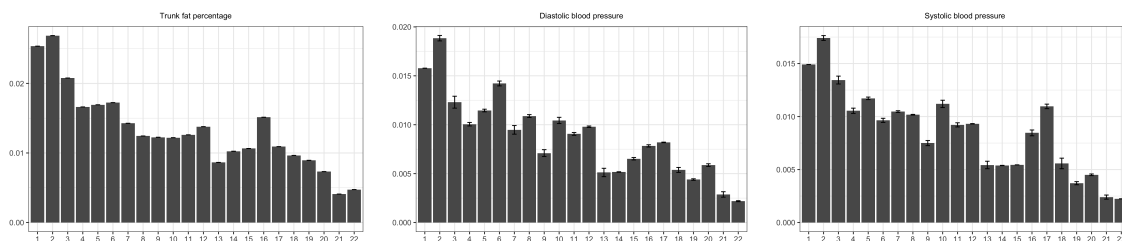


Figure 5: **RHE-reg-mc heritability estimates for trunk fat percentage, diastolic blood pressure, and systolic blood pressure in the UK Biobank**

4 Discussion

We have described RHE-reg-mc, a scalable estimator of multiple variance components in linear mixed models. RHE-reg-mc uses a randomized Method-of-Moments estimator to estimate a large number of variance components on datasets with hundreds of thousands of individuals and SNPs in less than an hour. The ability to estimate multiple variance components efficiently is useful both in obtaining unbiased genome-wide SNP heritability estimates as well as in heritability partitioning analyses.

There are several ways to further improve the runtime and memory usage of RHE-reg-mc method. First, in the context of multiple phenotypes measured for a fixed genotype matrix, the matrix T in normal equation 8 is the same for every given phenotype. Hence, we just need to compute the sufficient statistics of matrix T once in the multiple phenotype setting. Second, RHE-reg-mc can perform its computation in a streaming version which can lead to a highly memory efficient implementation.

5 Acknowledgments

This research was conducted using the UK Biobank Resource under applications 33127 and 33297. We thank the participants of UK Biobank for making this work possible. We thank Rob Brown for feedback on this manuscript. This work was funded by NIH grants R01HG009120 (B.P. and K.S.B.), R35GM125055 (S.S.), an Alfred P. Sloan Research Fellowship (S.S.), and a NSF grant III-1705121 (Y.W. and S.S.).

References

- [1] Laura Almasy and John Blangero. Multipoint quantitative-trait linkage analysis in general pedigrees. *The American Journal of Human Genetics*, 62(5):1198–1211, 1998.
- [2] Brendan K Bulik-Sullivan, Po-Ru Loh, Hilary K Finucane, Stephan Ripke, Jian Yang, Nick Patterson, Mark J Daly, Alkes L Price, Benjamin M Neale, Schizophrenia Working Group of the Psychiatric Genomics Consortium, et al. Ld score regression distinguishes confounding from polygenicity in genome-wide association studies. *Nature genetics*, 47(3):291, 2015.
- [3] Luke M Evans, Rasool Tahmasbi, Scott I Vrieze, Gonçalo R Abecasis, Sayantan Das, Steven Gazal, Douglas W Bjelland, Teresa R Candia, Michael E Goddard, Benjamin M Neale, et al. Comparison of methods that use whole genome data to estimate the heritability and genetic architecture of complex traits. *Nature genetics*, 50(5):737, 2018.
- [4] Hilary K Finucane, Brendan Bulik-Sullivan, Alexander Gusev, Gosia Trynka, Yakir Reshef, Po-Ru Loh, Verneri Anttila, Han Xu, Chongzhi Zang, Kyle Farh, et al. Partitioning heritability by functional annotation using genome-wide association summary statistics. *Nature genetics*, 47(11):1228, 2015.
- [5] David Golan, Eric S Lander, and Saharon Rosset. Measuring missing heritability: inferring the contribution of common variants. *Proceedings of the National Academy of Sciences*, 111(49):E5272–E5281, 2014.
- [6] Alexander Gusev, Gaurav Bhatia, Noah Zaitlen, Bjarni J Vilhjálmsson, Dorothee Diogo, Eli A Stahl, Peter K Gregersen, Jane Worthington, Lars Klareskog, Soumya Raychaudhuri, et al. Quantifying missing heritability at known gwas loci. *PLoS genetics*, 9(12):e1003993, 2013.
- [7] Alexander Gusev, S Hong Lee, Gosia Trynka, Hilary Finucane, Bjarni J Vilhjálmsson, Han Xu, Chongzhi Zang, Stephan Ripke, Brendan Bulik-Sullivan, Eli Stahl, et al. Partitioning heritability of regulatory and cell-type-specific variants across 11 common diseases. *The American Journal of Human Genetics*, 95(5):535–552, 2014.
- [8] JK Haseman and RC Elston. The investigation of linkage between a quantitative trait and a marker locus. *Behavior genetics*, 2(1):3–19, 1972.
- [9] MF Hutchinson. A stochastic estimator of the trace of the influence matrix for laplacian smoothing splines. *Communications in Statistics-Simulation and Computation*, 18(3):1059–1076, 1989.
- [10] Hans R Kunsch. The jackknife and the bootstrap for general stationary observations. *The Annals of Statistics*, 17(3):1217–1241, 1989.
- [11] Sang Hong Lee, Naomi R Wray, Michael E Goddard, and Peter M Visscher. Estimating missing heritability for disease from genome-wide association studies. *The American Journal of Human Genetics*, 88(3):294–305, 2011.
- [12] Edo Liberty and Steven W Zucker. The mailman algorithm: A note on matrix–vector multiplication. *Information Processing Letters*, 109(3):179–182, 2009.
- [13] Christoph Lippert, Gerald Quon, Eun Yong Kang, Carl M Kadie, Jennifer Listgarten, and David Heckerman. The benefits of selecting phenotype-specific variants for applications of mixed models in genomics. *Scientific reports*, 3:1815, 2013.
- [14] Po-Ru Loh, Gaurav Bhatia, Alexander Gusev, Hilary K Finucane, Brendan K Bulik-Sullivan, Samuela J Pollack, Teresa R de Candia, Sang Hong Lee, Naomi R Wray, Kenneth S Kendler, et al. Contrasting genetic architectures of schizophrenia and other complex diseases using fast variance-components analysis. *Nature genetics*, 47(12):1385, 2015.
- [15] Charles E McCulloch and Shayle R Searle. *Generalized, linear, and mixed models*. John Wiley & Sons, 2004.

- [16] H Desmond Patterson and Robin Thompson. Recovery of inter-block information when block sizes are unequal. *Biometrika*, 58(3):545–554, 1971.
- [17] Doug Speed and David Balding. Better estimation of snp heritability from summary statistics provides a new understanding of the genetic architecture of complex traits. *bioRxiv*, page 284976, 2018.
- [18] Doug Speed, Na Cai, Michael R Johnson, Sergey Nejentsev, David J Balding, UCLEB Consortium, et al. Reevaluation of snp heritability in complex human traits. *Nature genetics*, 49(7):986, 2017.
- [19] Doug Speed, Gibran Hemani, Michael R Johnson, and David J Balding. Improved heritability estimation from genome-wide snps. *The American Journal of Human Genetics*, 91(6):1011–1021, 2012.
- [20] Cathie Sudlow, John Gallacher, Naomi Allen, Valerie Beral, Paul Burton, John Danesh, Paul Downey, Paul Elliott, Jane Green, Martin Landray, et al. Uk biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS medicine*, 12(3):e1001779, 2015.
- [21] John Tukey. Bias and confidence in not quite large samples. *Ann. Math. Statist.*, 29:614, 1958.
- [22] Peter M Visscher, William G Hill, and Naomi R Wray. Heritability in the genomics era: concepts and misconceptions. *Nature reviews genetics*, 9(4):255, 2008.
- [23] Yue Wu and Sriram Sankararaman. A scalable estimator of snp heritability for biobank-scale data. *bioRxiv*, page 294470, 2018.
- [24] Jian Yang, Andrew Bakshi, Zhihong Zhu, Gibran Hemani, Anna AE Vinkhuyzen, Sang Hong Lee, Matthew R Robinson, John RB Perry, Ilja M Nolte, Jana V van Vliet-Ostaptchouk, et al. Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body mass index. *Nature genetics*, 47(10):1114, 2015.
- [25] Jian Yang, Beben Benyamin, Brian P McEvoy, Scott Gordon, Anjali K Henders, Dale R Nyholt, Pamela A Madden, Andrew C Heath, Nicholas G Martin, Grant W Montgomery, et al. Common snps explain a large proportion of the heritability for human height. *Nature genetics*, 42(7):565, 2010.
- [26] Jian Yang, S Hong Lee, Michael E Goddard, and Peter M Visscher. Gcta: a tool for genome-wide complex trait analysis. *The American Journal of Human Genetics*, 88(1):76–82, 2011.
- [27] Xiang Zhou. A unified framework for variance component estimation with summary statistics in genome-wide association studies. *The annals of applied statistics*, 11(4):2027, 2017.
- [28] Xiang Zhou and Matthew Stephens. Genome-wide efficient mixed-model analysis for association studies. *Nature genetics*, 44(7):821, 2012.

A Computing jackknife standard errors

For simplicity, here we explain how we can estimate the corresponding parameters of jackknife subsamples when we have only one variance component (one partition). It is easy to see how this approach can be generalized to the multiple variance components setting.

Let \mathbf{X} be a $N \times M$ matrix of standardized genotypes where N and M are the number of individuals and SNPs respectively. We partition \mathbf{X} to J non-overlapping blocks $\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(J)}$ such that $\mathbf{X} = [\mathbf{X}^{(1)}, \dots, \mathbf{X}^{(J)}]$. Note that for every j , \mathbf{X}_j is a $N \times M_j$ matrix where M_j is the number of SNPs in the j -th block.

Let $\mathbf{K} = \frac{1}{M} \mathbf{X} \mathbf{X}^T$ be the GRM. Moreover, Let $\mathbf{K}^{(-j)}$ be the GRM of the subsample created by removing the j -th block from \mathbf{X} . We will show how the values $\mathbf{y}^T \mathbf{K}^{(-j)} \mathbf{y}$ and $\text{tr}(\mathbf{K}^{(-j)^2})$ can be computed for every $j \in \{1, \dots, J\}$ efficiently.

$$\mathbf{X}^T \mathbf{y} = [\mathbf{X}^{(1)T} \mathbf{y}, \dots, \mathbf{X}^{(J)T} \mathbf{y}]^T \quad (13)$$

$$(\mathbf{X}^T \mathbf{y})^T (\mathbf{X}^T \mathbf{y}) = (\mathbf{X}^{(1)T} \mathbf{y})^2 + \dots + (\mathbf{X}^{(J)T} \mathbf{y})^2 \quad (14)$$

Therefore, for every j we have :

$$\mathbf{y}^T \mathbf{K}^{(-j)} \mathbf{y} = \frac{1}{M_j} [(\mathbf{X}^T \mathbf{y})^T (\mathbf{X}^T \mathbf{y}) - (\mathbf{X}^{(j)T} \mathbf{y})^2] \quad (15)$$

$$\mathbf{y}^T \mathbf{K} \mathbf{y} = \frac{1}{M} \mathbf{y}^T \mathbf{X} \mathbf{X}^T \mathbf{y} = \frac{1}{M} (\mathbf{X}^T \mathbf{y})^T (\mathbf{X}^T \mathbf{y}) \quad (16)$$

Therefore, we can compute $\mathbf{y}^T \mathbf{K}^{(j)} \mathbf{y}$, for every j , and $\mathbf{y}^T \mathbf{K} \mathbf{y}$ in $\mathcal{O}(\frac{BNM}{\max(\log_3 N, \log_3 M)})$.

We rewrite $\widehat{\text{tr}(\mathbf{K}^2)}$ as:

$$\widehat{\text{tr}(\mathbf{K}^2)} = \frac{1}{BM^2} \sum_b (\mathbf{X} \mathbf{X}^T \mathbf{z}_b)^T (\mathbf{X} \mathbf{X}^T \mathbf{z}_b) \quad (17)$$

Assume that $\mathbf{v}_b = \mathbf{X}^T \mathbf{z}_b$, then we have:

$$\mathbf{X} \mathbf{X}^T \mathbf{z}_b = \mathbf{X} \mathbf{v}_b = (\mathbf{v}_b^T \mathbf{X}^T)^T \quad (18)$$

Recall that \mathbf{v}_b is a vector of length M . We partition \mathbf{v}_b to J sub-vectors corresponding to the partitions of \mathbf{X} : $\mathbf{v}_{b1}, \dots, \mathbf{v}_{bJ}$ such that $\mathbf{v}_b = [\mathbf{v}_{b1}, \dots, \mathbf{v}_{bJ}]^T$ and \mathbf{v}_{bj} is a sub-vector of length M_j , for every $j \in \{1, \dots, J\}$.

$$\mathbf{X} \mathbf{v}_b = \mathbf{X}^{(1)} \mathbf{v}_{b1} + \dots + \mathbf{X}^{(J)} \mathbf{v}_{bJ} \quad (19)$$

Assume that $\mathbf{c}_{bj} = \mathbf{X}^{(j)} \mathbf{v}_{bj}$ for every $j \in \{1, \dots, J\}$ and $b \in \{1, \dots, B\}$. Let $C_{total} = \sum_b (\mathbf{c}_{b1}^T + \dots + \mathbf{c}_{bJ}^T) (\mathbf{c}_{b1} + \dots + \mathbf{c}_{bJ})$. We then have:

$$\widehat{\text{tr}(\mathbf{K}^2)} = \frac{1}{BM^2} C_{total} \quad (20)$$

Let $\mathbf{d}_b = \mathbf{c}_{b1} + \dots + \mathbf{c}_{bJ}$ for every $b \in B$. Now we can compute $\widehat{\text{tr}(\mathbf{K}^{(-j)^2})}$ for every $j \in \{1, \dots, J\}$ based on the values of C_{total} , \mathbf{d}_b , \mathbf{c}_{bj} in a way that does not effect the total running time :

$$\begin{aligned} \widehat{\text{tr}(\mathbf{K}^{(-j)^2})} &= \frac{1}{BM_j^2} [C_{total} - \sum_b (\mathbf{c}_{bj}^T \mathbf{d}_b + (\mathbf{d}_b - \mathbf{c}_{bj})^T \mathbf{c}_{bj})] \\ &= \frac{1}{BM_j^2} [C_{total} - \sum_b (\mathbf{c}_{bj}^T \mathbf{d}_b + \mathbf{d}_b^T \mathbf{c}_{bj} - \mathbf{c}_{bj}^T \mathbf{c}_{bj})] \\ &= \frac{1}{BM_j^2} [C_{total} - \sum_b (2\mathbf{c}_{bj}^T \mathbf{d}_b - \mathbf{c}_{bj}^T \mathbf{c}_{bj})] \\ &= \frac{1}{BM_j^2} [C_{total} - 2 \sum_b \mathbf{c}_{bj}^T \mathbf{d}_b + \sum_b \mathbf{c}_{bj}^T \mathbf{c}_{bj}] \end{aligned} \quad (21)$$

It is not hard to see that we can compute C_{total} in $\mathcal{O}(BNM)$. Moreover, we can compute c_{bj} for every $b = 1, \dots, B$ and $j = 1, \dots, J$ in $\mathcal{O}(NM_j)$, so all c_{bj} can be computed in $\mathcal{O}(BNM)$. We can again use the Mailman algorithm to speed up the vector-matrix multiplications. Therefore, the total running time for computing the variance component of all jackknife subsamples is $\mathcal{O}(\frac{BNM}{\max(\log_3 N, \log_3 M)})$.

B Including Covariates

$$\mathbf{y}|\epsilon, \beta_1, \dots, \beta_k = \mathbf{W}\alpha + \sum_k \mathbf{X}_k \beta_k + \epsilon \quad (22)$$

Here \mathbf{W} is a $N \times C$ matrix of covariates while α is a C -vector of coefficients. It is easy to see that the matrix $\mathbf{V} = \mathbf{I}_N - \mathbf{W}(\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T$ is symmetric and idempotent ($\mathbf{V}^2 = \mathbf{V}$) of rank $N - C$. Therefore, we consider the eigendecomposition of $\mathbf{V} = \mathbf{E} \mathbf{D} \mathbf{E}^T$, where \mathbf{D} is a diagonal matrix with $N - C$ ones and C zeros on the diagonal (we can assume that first $N - C$ elements are one). Now let the matrix $\mathbf{U}_{N \times (N-C)}$ represent the first $N - C$ columns of \mathbf{E} . It is not hard to see that \mathbf{U} satisfies $\mathbf{U}^T \mathbf{U} = \mathbf{I}_{N-C}$, $\mathbf{U} \mathbf{U}^T = \mathbf{V}$, $\mathbf{U}^T \mathbf{W} = 0$. Now we multiply by \mathbf{U}^T on both sides of the above equation:

$$\mathbf{U}^T \mathbf{y} = \mathbf{U}^T \sum_k \mathbf{X}_k \beta_k + \mathbf{U}^T \epsilon \quad (23)$$

$$\text{cov}(\mathbf{U}^T \mathbf{y}) = E[\mathbf{U}^T \mathbf{y} (\mathbf{U}^T \mathbf{y})^T] - E[\mathbf{U}^T \mathbf{y}] E[\mathbf{U}^T \mathbf{y}] \quad (24)$$

The matrix \mathbf{U}^T is constant and the vector \mathbf{y} is random. Therefore, we have $E[\mathbf{U}^T \mathbf{y}] = \mathbf{U}^T E[\mathbf{y}]$.

$$\begin{aligned} \mathbf{U}^T \mathbf{y} (\mathbf{U}^T \mathbf{y})^T &= (\mathbf{U}^T \sum_k \mathbf{X}_k \beta_k + \mathbf{U}^T \epsilon) (\mathbf{U}^T \sum_k \mathbf{X}_k \beta_k + \mathbf{U}^T \epsilon)^T = \\ &= \sum_i \sum_j \mathbf{U}^T \mathbf{X}_i \beta_i (\mathbf{U}^T \mathbf{X}_j \beta_j)^T + (\mathbf{U}^T \epsilon) \sum_i (\mathbf{U}^T \mathbf{X}_i \beta_i)^T + \sum_i \mathbf{U}^T \mathbf{X}_i \beta_i (\mathbf{U}^T \epsilon)^T + \mathbf{U}^T \epsilon (\mathbf{U}^T \epsilon)^T \end{aligned} \quad (25)$$

Hence

$$E[\mathbf{U}^T \mathbf{y} (\mathbf{U}^T \mathbf{y})^T] = \sum_k \frac{\sigma_{gk}^2}{M_k} (\mathbf{U}^T \mathbf{X}_k) (\mathbf{U}^T \mathbf{X}_k)^T + \sigma_\epsilon^2 \mathbf{U}^T \mathbf{U} \quad (26)$$

Using $\mathbf{K}_k = \frac{\mathbf{X}_k \mathbf{X}_k^T}{M_k}$, we have:

$$\text{cov}(\mathbf{U}^T \mathbf{y}) = \mathbf{U}^T \left(\sum_k \sigma_{gk}^2 \mathbf{K}_k \right) \mathbf{U} + \sigma_\epsilon^2 \mathbf{I}_{N-C} \quad (27)$$

The MoM estimator is obtained by solving the following ordinary least squares problem:

$$(\tilde{\sigma}_1^2, \dots, \tilde{\sigma}_K^2, \tilde{\sigma}_\epsilon^2) = \underset{(\sigma_1^2, \dots, \sigma_K^2, \sigma_\epsilon^2)}{\text{argmin}} \left\| \mathbf{U}^T \mathbf{y} (\mathbf{U}^T \mathbf{y})^T - \mathbf{U}^T \left(\sum_k \sigma_k^2 \mathbf{K}_k \right) \mathbf{U} - \sigma_\epsilon^2 \mathbf{I}_{N-C} \right\|_F^2 \quad (28)$$