

Supplementary Information
Confidence intervals for heritability via Haseman-Elston
regression

Tamar Sofer

Contents

1	Mathematical derivation	1
2	Computation	6
2.1	Variance component estimators	6
2.2	Confidence intervals for the variance components	6
2.3	Computing heritability estimates and their confidence intervals	7
2.3.1	A faster algorithm when the kinship matrix is the only source of correlation in the model	8
2.3.2	Meta-analysis of across studies when kinship is the only source of correlation	8
3	The Hispanic Community Health Study/Study of Latinos	10
3.1	Genotyping, imputation and quality control	11
3.2	Heritability estimation in the HCHS/SOL	11
4	Additional simulations	14
4.1	Results	15

1 Mathematical derivation

Suppose that a quantitative trait Y , measured on n individuals, follows the regression model

$$y_i = \mathbf{w}_i^T \boldsymbol{\beta} + \epsilon_i, \quad i = 1, \dots, n$$

where

$$E[\boldsymbol{\epsilon}] = \mathbf{0} \quad (1)$$

$$\text{var}[\boldsymbol{\epsilon}] = \sigma_e^2 \mathbf{I}_{n \times n} + \sigma_a^2 \mathbf{A}_1 + \dots + \sigma_k^2 \mathbf{K} = \boldsymbol{\Sigma} \quad (2)$$

and $\mathbf{A}, \dots, \mathbf{K}$ are $n \times n$ matrices modeling correlations between individuals. Let $a_{i,j}, \dots, k_{i,j}$ denote the i, j entries of the matrices $\mathbf{A}, \dots, \mathbf{K}$. Assuming that random effects due to $\mathbf{A}, \dots, \mathbf{K}$ are independent, we have that:

$$E[\epsilon_i \epsilon_j] = \text{cov}(\epsilon_i, \epsilon_j) = \sigma_e^2 \mathcal{I}_{(i=j)} + \sigma_a^2 a_{i,j} + \dots + \sigma_k^2 k_{i,j}.$$

Let $\widehat{\boldsymbol{\beta}}$ be an unbiased estimator of $\boldsymbol{\beta}$, and let $\hat{\epsilon}_i = y_i - \mathbf{w}_i^T \widehat{\boldsymbol{\beta}}$ be an estimator $\epsilon_i, i = 1, \dots, n$. We estimate the variance components in a residual regression, i.e. by taking the vector all unique pairs of residuals $\hat{\epsilon}_i \hat{\epsilon}_j, i \leq j$ (we can do it by taking the upper diagonal sub-matrix of $\widehat{\boldsymbol{\epsilon}} \widehat{\boldsymbol{\epsilon}}^T$ that includes the diagonal), denoted by $\tilde{\boldsymbol{\epsilon}}^d$ and regressing it according to the above model. The regression design matrix is given by:

$$\mathbf{X} = \begin{pmatrix} 1 & a_{1,1} & \dots & k_{1,1} \\ 0 & a_{1,2} & \dots & k_{1,2} \\ \vdots & \vdots & \vdots & \vdots \\ 0 & a_{1,n} & \dots & k_{1,n} \\ 1 & a_{2,2} & \dots & k_{2,2} \\ 0 & a_{2,3} & \dots & k_{2,3} \\ \vdots & \vdots & \vdots & \vdots \\ 0 & a_{2,n} & \dots & k_{2,n} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & a_{n-1,n-1} & \dots & k_{n-1,n-1} \\ 0 & a_{n-1,n} & \dots & k_{n-1,n} \\ 1 & a_{n,n} & \dots & k_{n,n} \end{pmatrix} = \begin{pmatrix} 1 & 1 & 1 & 1 \\ 0 & a_{1,2} & \dots & k_{1,2} \\ \vdots & \vdots & \vdots & \vdots \\ 0 & a_{1,n} & \dots & k_{1,n} \\ 1 & 1 & 1 & 1 \\ 0 & a_{2,3} & \dots & k_{2,3} \\ \vdots & \vdots & \vdots & \vdots \\ 0 & a_{2,n} & \dots & k_{2,n} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 1 & 1 \\ 0 & a_{n-1,n} & \dots & k_{n-1,n} \\ 1 & 1 & 1 & 1 \end{pmatrix}$$

(because $a_{i,i}, \dots, k_{i,i} = 1$ for all i). Denote, for simplicity of presentation, the vector of off-diagonal elements of $\mathbf{A}, \dots, \mathbf{K}$ by $\mathbf{l} = (l_{1,2}, l_{1,3}, \dots, l_{1,n}, l_{2,3}, \dots, l_{n-1,n})^T, l = 1, \dots, k$, and the vector of off-diagonal elements of $\hat{\boldsymbol{\epsilon}}\hat{\boldsymbol{\epsilon}}^T$ by $\tilde{\boldsymbol{\epsilon}}$. Then the least squares estimator of $(\sigma_e^2, \sigma_a^2, \dots, \sigma_k^2)$ is given by $(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\tilde{\boldsymbol{\epsilon}}^d$. Clearly, we have that

$$(\mathbf{X}^T\mathbf{X}) = \begin{pmatrix} n & n & n & \dots & n \\ n & n + \mathbf{a}^T\mathbf{a} & n + \mathbf{a}^T\mathbf{b} & \dots & n + \mathbf{a}^T\mathbf{k} \\ \vdots & & & & \vdots \\ n & n + \mathbf{k}^T\mathbf{a} & n + \mathbf{k}^T\mathbf{b} & \dots & n + \mathbf{k}^T\mathbf{k} \end{pmatrix}.$$

This is most likely a positive definite matrix as (we assume that) the matrices $\mathbf{A}, \dots, \mathbf{K}$ are not highly correlated. In addition, we have that

$$\mathbf{X}^T\tilde{\boldsymbol{\epsilon}} = \begin{pmatrix} \sum_{i=1}^n \hat{\epsilon}_i^2 \\ \sum_{i=1}^n \hat{\epsilon}_i^2 + \mathbf{a}^T\tilde{\boldsymbol{\epsilon}} \\ \vdots \\ \sum_{i=1}^n \hat{\epsilon}_i^2 + \mathbf{k}^T\tilde{\boldsymbol{\epsilon}} \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^n \hat{\epsilon}_i^2 \\ \sum_{i=1}^n \hat{\epsilon}_i^2 \\ \vdots \\ \sum_{i=1}^n \hat{\epsilon}_i^2 \end{pmatrix} + \begin{pmatrix} 0 \\ \mathbf{a}^T\tilde{\boldsymbol{\epsilon}} \\ \vdots \\ \mathbf{k}^T\tilde{\boldsymbol{\epsilon}} \end{pmatrix}.$$

Lemma 1:

$$(\mathbf{X}^T\mathbf{X})^{-1} \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{1}{n} \\ 0 \\ \vdots \\ 0 \end{pmatrix} \Leftrightarrow (\mathbf{X}^T\mathbf{X}) \begin{pmatrix} \frac{1}{n} \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$$

Proof: Because $(\mathbf{X}^T\mathbf{X})$ is non-singular, and from the properties of the inverse matrix, we have that $(\mathbf{X}^T\mathbf{X})^{-1}(\mathbf{X}^T\mathbf{X})\mathbf{v} = \mathbf{v}$ for every \mathbf{v} . ■

Lemma 2: Variance component estimators corresponding to the matrices $\mathbf{A}, \dots, \mathbf{K}$ depend only on the between-observation residuals of the form $\epsilon_i\epsilon_j$ for $i \neq j$ and do not depend on $\epsilon_i^2, i = 1, \dots, n$.

Proof: By noting that

$$(\mathbf{X}^T \mathbf{X}) \begin{pmatrix} \frac{1}{n} \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix},$$

we get from Lemma 1 that

$$(\mathbf{X}^T \mathbf{X})^{-1} \begin{pmatrix} \sum_{i=1}^n \hat{\epsilon}_i^2 \\ \sum_{i=1}^n \hat{\epsilon}_i^2 \\ \vdots \\ \sum_{i=1}^n \hat{\epsilon}_i^2 \end{pmatrix} = \begin{pmatrix} \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$$

which proves that the term $\sum_{i=1}^n \hat{\epsilon}_i^2$ contributes only to the estimator $\hat{\sigma}_e^2$. ■

Lemma 3: Denote by $\sigma_T^2 = \sigma_e^2 + \sigma_a^2 + \dots + \sigma_k^2$. Then $\hat{\sigma}_T^2 = \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2$.

Proof: We show that $\hat{\sigma}_e^2 + \hat{\sigma}_a^2 + \dots + \hat{\sigma}_k^2 = \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2$. In the proof of Lemma 2 we saw that

$$(\mathbf{X}^T \mathbf{X})^{-1} \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix} = \begin{pmatrix} \frac{1}{n} \\ 0 \\ \vdots \\ 0 \end{pmatrix}.$$

Since this $(\mathbf{X}^T \mathbf{X})^{-1}$ is symmetric, it follows that

$$\begin{pmatrix} 1 & 1 & \dots & 1 \end{pmatrix} (\mathbf{X}^T \mathbf{X})^{-1} = \begin{pmatrix} \frac{1}{n} & 0 & \dots & 0 \end{pmatrix}$$

Therefore

$$\begin{aligned} \hat{\sigma}_e^2 + \hat{\sigma}_a^2 + \dots + \hat{\sigma}_k^2 &\equiv \begin{pmatrix} 1 & 1 & \dots & 1 \end{pmatrix} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \tilde{\epsilon}) \\ &= \begin{pmatrix} \frac{1}{n} & 0 & \dots & 0 \end{pmatrix} (\mathbf{X}^T \tilde{\epsilon}) = \frac{1}{n} \sum_{i=1}^n \hat{\epsilon}_i^2, \end{aligned}$$

which completes the proof. ■

Lemma 4: *An estimator of the ratio between any variance component (or sum of variance components) and the total variance is a ratio between two quadratic forms.*

Proof: For $\mathbf{L} = \mathbf{A}, \dots, \mathbf{K}$, a quantity of the form $\mathbf{l}^T \tilde{\boldsymbol{\epsilon}} = \hat{\boldsymbol{\epsilon}}^T \mathbf{L}^- \hat{\boldsymbol{\epsilon}}/2$, where the matrix \mathbf{L}^- is the matrix \mathbf{L} with all diagonal values set to 0. An estimator a variance component σ_l^2 is a linear sum of the quadratic forms $\hat{\boldsymbol{\epsilon}}^T \mathbf{A}^- \hat{\boldsymbol{\epsilon}}, \dots, \hat{\boldsymbol{\epsilon}}^T \mathbf{K}^- \hat{\boldsymbol{\epsilon}}$, with coefficients the entries of the corresponding row of $(\mathbf{X}^T \mathbf{X})^{-1}$. Since a weighted sum of quadratic forms is a quadratic form, any variance component (and a sum of variance components) is also a quadratic form. Similarly, the total variance estimators is the quadratic form $\frac{1}{n} \hat{\boldsymbol{\epsilon}}^T \hat{\boldsymbol{\epsilon}}$. ■

Theorem: *We say that two matrices \mathbf{C}_1 and \mathbf{C}_2 are orthogonal in the trace inner product, or “trace orthogonal” if $\text{tr}(\mathbf{C}_1 \mathbf{C}_2) = 0$. If a matrix \mathbf{L}^- is trace orthogonal to all other matrices in the set $\{\mathbf{A}^-, \dots, \mathbf{K}^-\}$, then*

$$\hat{\sigma}_l^2 = \frac{1}{\sum_{j>i} l_{i,j}^2} \sum_{j>i} l_{i,j} \hat{\epsilon}_i \hat{\epsilon}_j.$$

Proof: Without loss of generality, assume that \mathbf{A} is trace orthogonal to $\mathbf{B}, \dots, \mathbf{K}$. First note that for the symmetric matrices with diagonal values set to zero $\mathbf{A}^-, \dots, \mathbf{K}^-$, $\text{tr}(\mathbf{A}^- \mathbf{L}^-) = 0$ if and only if $\mathbf{a}^T \mathbf{l} = 0$. Then

$$(\mathbf{X}^T \mathbf{X}) = \begin{pmatrix} n & n & n & \dots & n \\ n & n + \mathbf{a}^T \mathbf{a} & n & \dots & n \\ n & n & n + \mathbf{b}^T \mathbf{b} & \dots & n + \mathbf{b}^T \mathbf{k} \\ \vdots & \vdots & \vdots & & \vdots \\ n & n & n + \mathbf{k}^T \mathbf{b} & \dots & n + \mathbf{k}^T \mathbf{k}. \end{pmatrix}$$

Denote by $(\mathbf{X}^T \mathbf{X})_{[i,j]}^{-1}$ the i, j element in the matrix $(\mathbf{X}^T \mathbf{X})^{-1}$. First, we notice that the entries $(\mathbf{X}^T \mathbf{X})_{[2,j]}^{-1}, j = 3, \dots, k+1$ are all 0, because the $(\mathbf{X}^T \mathbf{X})_{[i,j]}^{-1}$ entry is a constant times the i, j minor of $(\mathbf{X}^T \mathbf{X})$, which has two identical columns (corresponding to the 1st and 2nd columns of $(\mathbf{X}^T \mathbf{X})$ when removing its 2nd row). Since the sum of the 2nd row of $(\mathbf{X}^T \mathbf{X})^{-1}$ is equal to 0, as we saw before, we get that $(\mathbf{X}^T \mathbf{X})_{[2,1]}^{-1} = -(\mathbf{X}^T \mathbf{X})_{[2,2]}^{-1}$.

We now argue that

$$\begin{aligned}
\hat{\sigma}_a^2 &\equiv (\mathbf{X}^T \mathbf{X})_{[2,1]}^{-1} \sum_{i=1}^n \epsilon_i^2 + (\mathbf{X}^T \mathbf{X})_{[2,2]}^{-1} \left(\sum_{i=1}^n \epsilon_i^2 + \mathbf{a}^T \tilde{\boldsymbol{\epsilon}} \right) \\
&= (\mathbf{X}^T \mathbf{X})_{[2,1]}^{-1} \sum_{i=1}^n \epsilon_i^2 - (\mathbf{X}^T \mathbf{X})_{[2,1]}^{-1} \left(\sum_{i=1}^n \epsilon_i^2 + \mathbf{a}^T \tilde{\boldsymbol{\epsilon}} \right) \\
&= -(\mathbf{X}^T \mathbf{X})_{[2,1]}^{-1} \mathbf{a}^T \tilde{\boldsymbol{\epsilon}} \stackrel{(4)}{=} \frac{1}{\sum_{j>i} a_{i,j}^2} \sum_{j>i} a_{i,j} \epsilon_i \epsilon_j,
\end{aligned}$$

where we need to show that equality (4) holds to complete the proof. We need to show that $-(\mathbf{X}^T \mathbf{X})_{[2,1]}^{-1} = \frac{1}{\sum_{j>i} a_{i,j}^2}$. Consider now the matrix $\mathbf{X}^T \mathbf{X}$. One can derive its determinant from its second row, as:

$$\begin{aligned}
|\mathbf{X}^T \mathbf{X}| &= (\mathbf{X}^T \mathbf{X})_{[2,1]} M_{2,1} - (\mathbf{X}^T \mathbf{X})_{[2,2]} M_{2,2} + \dots + (-1)^{k+1} (\mathbf{X}^T \mathbf{X})_{[2,1k]} M_{2,k} \\
&= n M_{2,1} - (n + \mathbf{a}^T \mathbf{a}) M_{2,2} + 0 + \dots + 0 \\
&= n |\mathbf{X}^T \mathbf{X}| (\mathbf{X}^T \mathbf{X})_{[2,1]}^{-1} - (n + \mathbf{a}^T \mathbf{a}) |\mathbf{X}^T \mathbf{X}| (\mathbf{X}^T \mathbf{X})_{[2,2]}^{-1} \\
&= n |\mathbf{X}^T \mathbf{X}| (\mathbf{X}^T \mathbf{X})_{[2,1]}^{-1} - (n + \mathbf{a}^T \mathbf{a}) |\mathbf{X}^T \mathbf{X}| (\mathbf{X}^T \mathbf{X})_{[2,1]}^{-1} \\
&= -\mathbf{a}^T \mathbf{a} |\mathbf{X}^T \mathbf{X}| (\mathbf{X}^T \mathbf{X})_{[2,1]}^{-1}
\end{aligned}$$

Therefore, we get that $-(\mathbf{X}^T \mathbf{X})_{[2,1]}^{-1} = \frac{1}{\mathbf{a}^T \mathbf{a}}$, which completes the proof. \blacksquare

2 Computation

2.1 Variance component estimators

While any unbiased estimator of $\hat{\boldsymbol{\beta}}$ suffices to generate residuals $\hat{\boldsymbol{\epsilon}}$ and use them to obtain variance component estimators as $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \tilde{\boldsymbol{\epsilon}}^d$, a more efficient estimator iterates between estimating $\boldsymbol{\beta}$ and the variance component estimator as follows:

1. Initialization step: set $\hat{\boldsymbol{\beta}}^{(0)} = (\mathbf{W}^T \mathbf{W})^{-1} \mathbf{W}^T \mathbf{y}$.

2. Iteration step:

(a) Given the k th estimator of $\boldsymbol{\beta}$, $\hat{\boldsymbol{\beta}}^{(k)}$, set $\hat{\boldsymbol{\epsilon}} = \mathbf{y} - \mathbf{W} \boldsymbol{\beta}^{(k)}$ and $\tilde{\boldsymbol{\epsilon}}$ is the vector of upper diagonal matrix (including the diagonal) of $\hat{\boldsymbol{\epsilon}} \hat{\boldsymbol{\epsilon}}^T$. Set $\hat{\boldsymbol{\sigma}}^{2,(k)} = (\hat{\sigma}_e^{2,(k)}, \hat{\sigma}_a^{2,(k)}, \dots, \hat{\sigma}_k^{2,(k)}) =$

$$(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \tilde{\boldsymbol{\epsilon}}^d.$$

- (b) Given the k th estimator of σ^2 , $\hat{\boldsymbol{\sigma}}^{2,(k)}$, let $\hat{\boldsymbol{\Sigma}}^{(k)} = \hat{\sigma}_e^{2,(k)} \mathbf{I}_{n \times n} + \hat{\sigma}_a^{2,(k)} \mathbf{A} + \dots + \hat{\sigma}_k^{2,(k)} \mathbf{K}$ with inverse $\hat{\boldsymbol{\Sigma}}^{-1,(k)}$. Set $\hat{\boldsymbol{\beta}}^{(k+1)} = (\mathbf{W}^T \hat{\boldsymbol{\Sigma}}^{-1,(k)} \mathbf{W})^{-1} \mathbf{W}^T \hat{\boldsymbol{\Sigma}}^{-1,(k)} \mathbf{y}$

The iteration step repeats until convergence.

2.2 Confidence intervals for the variance components

From Lemma 4, any variance components (or sum of variance components) is given as a quadratic form. Let \mathbf{Q} be the quadratic form corresponding to a variance component estimate $\hat{\sigma}_l^2$, such that $\hat{\sigma}_l^2 = \hat{\boldsymbol{\epsilon}}^T \mathbf{Q} \hat{\boldsymbol{\epsilon}}$. Then this $\hat{\sigma}_l^2$ is distributed as the sum of independent $\chi_{(1)}^2$ variables in $\sum_{i=1}^n \lambda_i \chi_{(1)}^2$, where $\lambda_1, \dots, \lambda_n$ are the eigenvalues of $\mathbf{Q} \text{cov}(\hat{\boldsymbol{\epsilon}})$. In practice, for $\text{cov}(\hat{\boldsymbol{\epsilon}})$ we use the estimated $\hat{\boldsymbol{\Sigma}}(\hat{\sigma}_e^2, \dots, \hat{\sigma}_k^2)$. Functions in the package `CompQuadForm` calculate the probability function (or survival function) of this quadratic form based on $\lambda_1, \dots, \lambda_n$. While it takes times to compute the eigenvalues, once they are computed, a calculating the probabilities associated with the quadratic form over a grid is simple and quick. We can test the hypothesis $H_0 : \sigma_l^2 = 0$ by calculating the probability

$$\Pr(\hat{\boldsymbol{\epsilon}}^T \mathbf{Q} \hat{\boldsymbol{\epsilon}} = 0) = 1 - \Pr(\hat{\boldsymbol{\epsilon}}^T \mathbf{Q} \hat{\boldsymbol{\epsilon}} > 0),$$

and calculate two-sided confidence intervals for $\hat{\sigma}_l^2$ by calculating the survival probabilities over a grid, and taking the appropriate quantiles. For example, for a 95% confidence interval we take the values (c_1, c_2) for which

$$\begin{aligned} c_1 &= u : \Pr(\boldsymbol{\epsilon}^T \mathbf{Q} \boldsymbol{\epsilon} > u) = 0.025 \\ c_2 &= u : \Pr(\boldsymbol{\epsilon}^T \mathbf{Q} \boldsymbol{\epsilon} > u) = 0.975. \end{aligned}$$

We find these values using a binary search on the interval $[0, \hat{\sigma}_T^2]$.

2.3 Computing heritability estimates and their confidence intervals

Suppose that the variance component corresponding to the kinship matrix is σ_k^2 , which quadratic form denoted by \mathbf{Q}_k . We estimate heritability as $\hat{h}_k = \hat{\sigma}_k^2 / \hat{\sigma}_T^2$. However, we cannot use the confidence intervals for σ_k^2 to construct confidence intervals for h_k . Instead, we note that the point

estimate is \hat{h}_k is given by:

$$\hat{h}_k = \frac{\hat{\boldsymbol{\epsilon}}^T \mathbf{Q}_k \hat{\boldsymbol{\epsilon}}}{\frac{1}{n} \hat{\boldsymbol{\epsilon}}^T \mathbf{I} \hat{\boldsymbol{\epsilon}}} = \frac{\mathbf{x}^T \hat{\boldsymbol{\Sigma}}^{1/2} \mathbf{Q}_k \hat{\boldsymbol{\Sigma}}^{1/2} \mathbf{x}}{\frac{1}{n} \mathbf{x}^T \hat{\boldsymbol{\Sigma}} \mathbf{x}} = \frac{\mathbf{x}^T \mathbf{F} \mathbf{x}}{\mathbf{x}^T \mathbf{G} \mathbf{x}}$$

where $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, for $\mathbf{F} = \hat{\boldsymbol{\Sigma}}^{1/2} \mathbf{Q}_k \hat{\boldsymbol{\Sigma}}^{1/2}$ and $\mathbf{G} = \hat{\boldsymbol{\Sigma}}/n$. Thus, it is a ratio between two quadratic forms in (what we assume are) normal variables. For the squared root $\hat{\boldsymbol{\Sigma}}^{1/2}$, we use the Cholesky decomposition of $\hat{\boldsymbol{\Sigma}}$.

Now, we use the saddlepoint approximation for the distribution of a ratio of quadratic forms in normal variables, proposed by Lieberman (1994). For a given potential value of h_k , say h_k^* , we can calculate the survival probability

$$Pr(h_k \geq h_k^*) \cong 1 - \Phi(\hat{\xi}) + \phi(\hat{\xi}) \left[\frac{1}{\hat{z}} - \frac{1}{\hat{\xi}} \right]$$

where Φ and ϕ are the standard normal cdf and pdf, and

$$\begin{aligned} \hat{z} &= \hat{\omega} \left\{ 2 \sum_{i=1}^n \frac{d_i^{*2}}{(1 - 2\hat{\omega}d_i^*)^2} \right\}^{1/2} \\ \hat{\xi} &= \left\{ \sum_{i=1}^n \ln(1 - 2\hat{\omega}d_i^*) \right\}^{1/2} \text{sgn}(\hat{\omega}) \end{aligned}$$

and d_1^*, \dots, d_n^* are the eigenvalues of the matrix $\mathbf{D}^* = \mathbf{F} - h_k^* \mathbf{G}$, and $\hat{\omega}$ is the corresponding saddlepoint satisfying

$$\sum_{i=1}^n \frac{d_i^*}{1 - 2\hat{\omega}d_i^*} = 0.$$

Confidence intervals are then built, as before, using a binary search to find the values satisfying the required probabilities at the tails.

2.3.1 A faster algorithm when the kinship matrix is the only source of correlation in the model

Computing the eigenvalues $d_1^*(h_k^*), \dots, d_n^*(h_k^*)$ takes time. However, in the case where we only have a single kinship matrix, denoted by \mathbf{K} we can compute the eigen decomposition of the matrix \mathbf{K}^- once to obtain eigenvalues $\lambda_1, \dots, \lambda_n$, and then transform these eigenvalues to obtain the

eigenvalues $d_1^*(h_k^*), \dots, d_n^*(h_k^*)$ for each value h_k^* . To see this, suppose that \mathbf{u} is an eigenvector of \mathbf{K}^- with eigenvalues λ . Then, by definition:

$$\mathbf{K}^- \mathbf{u} = \lambda \mathbf{u}.$$

Since $\boldsymbol{\Sigma} = \sigma_k^2(\mathbf{K}^- + \mathbf{I}) + \sigma_e^2\mathbf{I}$, it is straightforward to see that \mathbf{u} is also an eigenvector of $\boldsymbol{\Sigma}$:

$$\boldsymbol{\Sigma} \mathbf{u} = [\sigma_k^2(\mathbf{K}^- + \mathbf{I}) + \sigma_e^2\mathbf{I}] \mathbf{u} = (\sigma_k^2\lambda + \sigma_k^2 + \sigma_e^2) \mathbf{u}.$$

Similarly, \mathbf{u} is an eigenvector of $\boldsymbol{\Sigma}^{1/2}$ with eigenvalue $\sqrt{\sigma_k^2\lambda + \sigma_k^2 + \sigma_e^2}$, which finally leads us to the transformation between an eigenvalue λ of \mathbf{A} to an eigenvalue of $\mathbf{D}^* = \mathbf{F} - h_k^*\mathbf{G}$ given by:

$$d_i^*(h_k^*, \lambda_i) = \frac{1}{2 \sum_{i < j} v_{ij}^2} \lambda_i (\lambda_i \sigma_k^2 + \sigma_k^2 + \sigma_e^2) - h_k^* (\lambda_i \sigma_k^2 + \sigma_k^2 + \sigma_e^2) / n.$$

As before, we use the estimated $\hat{\sigma}_k^2, \hat{\sigma}_e^2$ instead of the true unknown quantities.

2.3.2 Meta-analysis of across studies when kinship is the only source of correlation

A meta-analytic estimator. Suppose that there are S studies that we wanted to combined in meta-analysis. We assume that kinship is the only source of correlation. Each study has a vector of residuals $\hat{\boldsymbol{\epsilon}}_s = (\hat{\epsilon}_{s,1}, \dots, \hat{\epsilon}_{s,n_s})^T, s = 1, \dots, S$. Consider the Haseman-Elston regression, but incomplete, so that only the pairs of multiplied residuals within study are used (i.e. only $\hat{\epsilon}_{s,i}\hat{\epsilon}_{s,j}$ are regressed against entries of the kinship covariance matrix, but not $\hat{\epsilon}_{s,i}\hat{\epsilon}_{t,j}$). Therefore, cross-study kinship estimates are not used in the regression, however no assumption is made on them. In other words, we do not need to assume that participant in one study is genetically independent (no alleles shared IBD) of a participant in another study. It is straightforward to show that the meta-analytic estimator of σ_e^2 is given by $\hat{\sigma}_e^2 = \sum_{s=1}^S \sum_{i=1}^{n_s} \hat{\epsilon}_{s,i}^2$. Let $\hat{\boldsymbol{\epsilon}} = (\hat{\boldsymbol{\epsilon}}_1^T, \dots, \hat{\boldsymbol{\epsilon}}_S^T)^T$. Then the meta-analysis kinship variance component estimator is given by

$$\hat{\sigma}_k^2 = \frac{1}{\text{tr}(\mathbf{K}_s^- \mathbf{K}_s^-)} \hat{\boldsymbol{\epsilon}}^T \mathbf{K}_s^- \hat{\boldsymbol{\epsilon}}$$

where \mathbf{K}_s^- is the block diagonal matrix that have all the study-specific kinship matrix (without their diagonal values) arranged diagonally, as

$$\mathbf{K}^s = \begin{pmatrix} \mathbf{K}_1^- & \mathbf{0} & \dots & \dots \\ \mathbf{0} & \mathbf{K}_2^- & & \mathbf{0} \\ \vdots & & \ddots & \\ \vdots & \mathbf{0} & \mathbf{0} & \mathbf{K}_S^- \end{pmatrix}$$

To see that this meta-analytic estimator of σ_k^2 is unbiased, note first that $\text{cov}(\hat{\epsilon}) = (\sigma_e^2 + \sigma_k^2)\mathbf{I} + \sigma_k^2\mathbf{K}^-$, where now \mathbf{K}^- is the kinship matrix with kinship coefficients between the individuals across studies. Now, from characteristics of quadratic forms, we have that

$$\begin{aligned} E[\hat{\sigma}_k^2] &= E\left[\frac{1}{\text{tr}(\mathbf{K}_s^- \mathbf{K}_s^-)} \hat{\epsilon}^T \mathbf{K}_s^- \hat{\epsilon}\right] = \frac{1}{\text{tr}(\mathbf{K}_s^- \mathbf{K}_s^-)} \text{tr}(\mathbf{K}_s^- \text{cov}(\hat{\epsilon})) \\ &= \frac{1}{\text{tr}(\mathbf{K}_s^- \mathbf{K}_s^-)} \text{tr}(\mathbf{K}_s^- (\sigma_e^2 + \sigma_k^2)\mathbf{I} + \sigma_k^2\mathbf{K}^-) \\ &= \frac{1}{\text{tr}(\mathbf{K}_s^- \mathbf{K}_s^-)} \text{tr}(\mathbf{K}_s^- \sigma_k^2 \mathbf{K}_s^-) = \sigma_k^2. \end{aligned}$$

Computing the meta-analytic heritability estimator and confidence intervals. The eigenvalues result shows that all we need to calculate heritability estimates and confidence intervals are eigenvalues of the matrix \mathbf{K}^- (the kinship matrix without the diagonal), estimated σ_e^2 , σ_k^2 , and the sum of the entries of \mathbf{K}^- ($2\sum_{i<j} k_{ij}^2$). This result could be used to extend our methods to meta-analysis of information from multiple studies. Suppose that each of m independent studies calculated the residuals from a “null model” (i.e. a regression model without genetic fixed effects other than PCs). Then, each study s reports:

1. $\mathcal{K}^s = 2\sum_{i<j} k_{ij}^2$,
2. $\hat{\sigma}_{k,s}^2$,
3. $\hat{\sigma}_{e,s}^2$,
4. The number of participants in the study n_s ,
5. The eigenvalues $\lambda_1^s, \dots, \lambda_{n_s}^s$ of the matrix \mathbf{K}_s^- .

Then, the meta-analysis estimates of the kinship and error variance components, and \mathcal{K}^S are given by:

$$\begin{aligned}\hat{\sigma}_k^2 &= \frac{\sum_{s=1}^S \mathcal{K}^s \hat{\sigma}_{k,s}^2}{\sum_{s=1}^S \mathcal{K}^s} \\ \hat{\sigma}_e^2 &= \frac{\sum_{s=1}^S n_s \hat{\sigma}_{e,s}^2}{\sum_{s=1}^S n_s}, \\ \mathcal{K}^S &= \sum_{s=1}^S \mathcal{K}^s,\end{aligned}$$

and the eigenvalues of the cross-study \mathbf{K}^- matrix are taken to be $\lambda_1^1, \dots, \lambda_{n_1}^1, \dots, \lambda_1^S, \dots, \lambda_{n_S}^S$. Using these, the central location that can calculate heritability estimates and confidence intervals.

3 The Hispanic Community Health Study/Study of Latinos

The HCHS/SOL, (LaVange et al., 2010; Sorlie et al., 2010)) is a community based cohort study, following self-identified Hispanic individuals from four field centers (Chicago, IL; Miami, FL; Bronx, NY; and San Diego, CA). Individuals were sampled via a two-stage sampling scheme, in which households were randomly sampled from sampled community block units. Almost 13,000 study participants consented for genotyping. HCHS/SOL individuals are classified into ‘genetic analysis groups’, classes that are based on self reported ethnicities and genetic similarity (Conomos et al., 2016). The genetic analysis groups are Central American, Cuban, Dominican, Mexican, Puerto Rican, and South American. This study was approved by the institutional review boards at each field center, where all subjects gave written informed consent.

3.1 Genotyping, imputation and quality control

Blood samples from HCHS/SOL individuals were genotyped on a custom array consisting of Illumina Omni 2.5M content plus $\sim 150,000$ custom markers selected to include ancestry-informative markers, variants characteristic of Amerindian populations, known GWAS hits and other candidate gene polymorphisms. Quality control was similar to the procedure described in Laurie et al. (2010) and included checks for sample identity, batch effects, missing call rate, chromosomal anomalies (Laurie et al., 2012), deviation from Hardy-Weinberg equilibrium, Mendelian errors, and duplicate

sample discordance. 12,803 samples passed quality control, and 2,232,944 SNPs passed quality filters. Pairwise kinship coefficients and principal components reflecting ancestry were estimated in an iterative procedure which accounts for admixture (Conomos et al., 2016). All common variants were used to estimate kinship coefficients.

3.2 Heritability estimation in the HCHS/SOL

In each group of interest, including all individuals (‘pooled’ analysis), or specific genetic analysis groups, we randomly removed related individuals, to generate a set of individuals without any pair having kinship coefficient higher than 2^{-11} . Due to the sampling structure of the HCHS/SOL, the correlation between individuals is modeled in a kinship matrix, and also via matrices corresponding to community block units, and households. We estimated variance components via the procedure described here, with the three correlation matrices. We utilized the availability of environmental correlation to also estimate the contribution of modeled environmental factors (block unit and household) to the phenotypic variance. Finally, we also demonstrate the use of our method for meta-analysis by removing individuals from shared household to generate a restricted set in which none of the individuals live in the same house, and used the proposed procedure for calculating heritability in meta-analysis. Note that for this purpose we neglected block unit correlation and assume that there is no correlation due to block unit sharing.

We estimated heritability for the FEV1 (a measure of lung function), standing height, depression score (CESD10, a sum of ten questionnaire items related to depression in the past few weeks of filling the form), SBP (systolic blood pressure), and dental caries, a count of tooth decays and cavities across all teeth of a participant. Finally, all regression models were adjusted (via the design matrix \mathbf{W}) to the 5 first principal components, study center, age, sex, and genetic analysis group (in the pooled models). For some traits we used additional covariates. Table 1 provides the various estimates and confidence intervals.

Analysis	n	Height	Depression score	SBP	Dental caries	FEV1
Full: Including environmentally correlated individuals						
Heritability	10,255	0.58 (0.47,0.69)	0.04 (0.00,0.10)	0.20 (0.12,0.27)	0.17 (0.09,0.25)	0.25 (0.17,0.33)
Environmental variance	10,255	0.10 (0.06,0.14)	0.06 (0.02,0.10)	0.02 (0.00,0.06)	0.06 (0.02,0.10)	0.04 (0.00,0.09)
Restricted: without environmentally correlated individuals (heritability only)						
Pooled	7,848	0.57 (0.45,0.69)	0.03 (0.00,0.11)	0.19 (0.11,0.28)	0.19 (0.09,0.28)	0.25 (0.16,0.34)
Central American	867	0.26 (0.00,0.88)	0.00 (0.00,0.62)	0.00 (0.00,0.56)	0.17 (0.00,0.75)	0.16 (0.00,0.75)
South American	544	0.48 (0.00,1.00)	0.00 (0.00,0.88)	0.93 (0.06,1.00)	0.00 (0.00,0.88)	0.57 (0.00,1.00)
Mexican	2,862	0.48 (0.28,0.72)	0.04 (0.00,0.22)	0.20 (0.02,0.41)	0.24 (0.06,0.44)	0.27 (0.06,0.50)
Puerto Rican	1,479	0.47 (0.12,0.88)	0.00 (0.00,0.38)	0.27 (0.00,0.62)	0.25 (0.00,0.62)	0.49 (0.12,0.88)
Cuban	1,370	0.83 (0.31,1.00)	0.00 (0.00,0.50)	0.23 (0.00,0.75)	0.06 (0.00,0.62)	0.09 (0.00,0.62)
Dominican	726	0.84 (0.06,1.00)	0.02 (0.00,0.75)	0.71 (0.00,1.00)	0.00 (0.00,0.75)	0.61 (0.00,1.00)
Meta-analysis	7,848	0.53 (0.38,0.72)	0.02 (0.00,0.16)	0.24 (0.09,0.41)	0.19 (0.03,0.34)	0.30 (0.16,0.47)

Table 1: Comparison of heritability estimates obtained in analysis of various subgroups of the HCHS/SOL. The top and bottom parts consider data sets with and without environmentally correlated individuals. When environmental correlation was present, the analysis included all participants. The heritability is $\hat{\sigma}_k^2/\hat{\sigma}_T^2$, and the environmental variance proportion is $(\hat{\sigma}_h^2 + \hat{\sigma}_c^2)/\hat{\sigma}_T^2$, for σ_h^2, σ_c^2 variance component corresponding to household and community sharing matrices. The models without environmentally correlated individuals had smaller sample sizes (because individuals who shared household were removed), and we compared heritability estimates from the pooled analysis of all individuals, the various ethnic subgroups, and their meta-analysis. SBP is systolic blood pressure. Dental caries is a measure of teeth damage. FEV1 is a measure of lung function. The distribution of the ratio between the appropriate quadratic forms was approximation using the saddlepoint approximation of Lieberman (1994).

4 Additional simulations

We performed additional simulations in simple settings in which we expected that confidence intervals that are based on asymptotic normality will be limited: when the kinship values are small, and when the kinship and household matrices are somewhat “correlated”. For the following two settings, we ran 1,000 simulations for each sample size of $n = 1,500, 3,000, 5,000,$ and $12,784,$ with the largest value selected to match the largest HCHS/SOL sample size used. In both settings, the kinship and household matrices were block diagonal with 3×3 matrices, representing sets of three individuals who live in the same house. Therefore, the block corresponding to three people living in the same house in the household matrix was a 3×3 with all values equal to 1. The matrices describing the correlations for the entire simulated sample were thus:

$$\mathbf{H} = \begin{pmatrix} \mathbf{H}_h & \mathbf{0} & \mathbf{0} & \cdots \\ \mathbf{0} & \mathbf{H}_h & \mathbf{0} & \cdots \\ \mathbf{0} & \mathbf{0} & \mathbf{H}_h & \cdots \\ \vdots & & & \ddots \end{pmatrix}, \quad \mathbf{K} = \begin{pmatrix} \mathbf{K}_h & \mathbf{0} & \mathbf{0} & \cdots \\ \mathbf{0} & \mathbf{K}_h & \mathbf{0} & \cdots \\ \mathbf{0} & \mathbf{0} & \mathbf{K}_h & \cdots \\ \vdots & & & \ddots \end{pmatrix},$$

with

$$\mathbf{H}_h = \begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}$$

in the two settings, and

$$\mathbf{K}_h^1 = \begin{pmatrix} 1 & 0.4 & 0.5 \\ 0.4 & 1 & 0.6 \\ 0.5 & 0.6 & 1 \end{pmatrix}$$

in the “correlated kinship and household matrices” settings 1, and

$$\mathbf{K}_h^2 = \begin{pmatrix} 1 & 0.05 & 0.05 \\ 0.05 & 1 & 0.1 \\ 0.05 & 0.1 & 1 \end{pmatrix}$$

in the “small kinship values” settings 2. In both settings the residuals were sample to have $\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma_e^2 \mathbf{I} + \sigma_k^2 \mathbf{K} + \sigma_h^2 \mathbf{H})$, with $\sigma_e^2 = 100$, $\sigma_k^2 = 40$, and $\sigma_h^2 = 15$. To simulate ϵ with this distribution we computed the squared root of the matrix $\Sigma = \sigma_e^2 \mathbf{I} + \sigma_k^2 \mathbf{K} + \sigma_h^2 \mathbf{H}$. Then, we used this $\Sigma^{1/2}$ by having $\epsilon = \Sigma^{1/2} \epsilon_{iid}$, where $\epsilon_{iid} \sim \mathcal{N}(0, \mathbf{I})$, which is straight forward to sample.

4.1 Results

Tables 2 and 3 provide the results comparing the proposed HE-based method for calculating confidence intervals with the AI-REML method implemented in the GENESIS R package that uses normal approximation for the asymptotic distribution of the variance components for settings 1 and 2 respectively. The performance measures are the same as in the main manuscript: coverage (the proportion of simulations the true value was contained in the confidence interval), width (average width of the confidence interval), and root-mean-square-error comparing the estimated variance component to its true value across simulations.

One can see that the performance of both methods improve with the sample size, but HE confidence intervals always have coverage at least 0.95, and tending to be larger when the sample sizes are smaller, while the REML (GENESIS) confidence intervals often have poor coverage for small sample size and/or small values of the correlation matrices. The household matrix is the same in the two settings. However, in the first settings, the household and kinship matrices are more similar, and this leads to higher uncertainly - wider confidence intervals and larger RMSEs - of the corresponding variance component. Interestingly, the coverage of the household variance component of GENESIS remains poor in setting 2 even in large sample sizes, despite similar width and RMSE to the HE. This is because the estimated kinship variance component was in fact 0 in many of the simulations. The mean width of the confidence intervals tend to be larger for HE (which also have more coverage) and becomes almost the same as that of GENESIS (REML) as the sample size becomes large.

n	method	Kinship			Household		
		coverage	width	RMSE	coverage	width	RMSE
1500	HE	1.00	0.79	0.21	1.00	0.39	0.10
1500	REML - GENESIS	0.61	0.73	0.18	0.59	0.31	0.09
3000	HE	1.00	0.69	0.19	1.00	0.33	0.09
3000	GENESIS	0.70	0.64	0.16	0.69	0.28	0.08
5000	HE	1.00	0.66	0.18	1.00	0.32	0.09
5000	GENESIS	0.73	0.61	0.16	0.72	0.27	0.08
12784	HE	0.97	0.37	0.09	0.96	0.17	0.05
12784	GENESIS	0.95	0.36	0.09	0.95	0.17	0.05

Table 2: Results from simulation setting 1, with \mathbf{K}_h^1 .

n	method	Kinship			Household		
		coverage	width	RMSE	coverage	width	RMSE
1500	HE	1.00	1.00	0.44	1.00	0.21	0.05
1500	GENESIS	0.38	0.80	0.43	0.89	0.15	0.04
3000	HE	1.00	1.00	0.42	1.00	0.17	0.04
3000	GENESIS	0.46	0.83	0.41	0.83	0.13	0.04
5000	HE	1.00	1.00	0.37	0.96	0.14	0.03
5000	GENESIS	0.56	0.89	0.36	0.84	0.11	0.03
12784	HE	1.00	0.84	0.27	0.97	0.09	0.02
12784	GENESIS	0.76	0.84	0.27	0.88	0.08	0.02

Table 3: Results from simulation setting 2, with \mathbf{K}_h^2 .

References

- CONOMOS, M. P., LAURIE, C. A., STILP, A. M., GOGARTEN, S. M., MCHUGH, C. P., NELSON, S. C., SOFER, T., FERNÁNDEZ-RHODES, L., JUSTICE, A. E., GRAFF, M., YOUNG, K. L., SEYERLE, A., AVERY, C., TAYLOR, K., ROTTER, J., TALAVERA, G., DAVIGLUS, M., WASSERTHEIL-SMOLLER, S., SCHNEIDERMAN, N., HEISS, G., KAPLAN, R., FRANCESCHINI, N., REINER, A., SHAFFER, G., JOHN R AND BARR, KERR, K., BROWNING, S., BROWNING, B., WEIR, B., AVILÉS-SANTA, L., PAPANICOLAOU, G., LUMLEY, T., SZPIRO, A., NORTH, K., RICE, K., THORNTON, T. and LAURIE, C. (2016). Genetic Diversity and Association Studies in US Hispanic/Latino Populations: Applications in the Hispanic Community Health Study/Study of Latinos. *The American Journal of Human Genetics*, **98** 165–184.
- LAURIE, C. ET AL. (2010). Quality control and quality assurance in genotypic data for genome-wide association studies. *Genetic Epidemiology*, **34** 591–602.

- LAURIE, C. C., LAURIE, C. A., RICE, K., DOHENY, K. F., ZELNICK, L. R., MCHUGH, C. P., LING, H., HETRICK, K. N., PUGH, E. W., AMOS, C. ET AL. (2012). Detectable clonal mosaicism from birth to old age and its relationship to cancer. *Nature Genetics*, **44** 642–650.
- LAVANGE, L. M., KALSBECK, W. D., SORLIE, P. D., AVILÉS-SANTA, L. M., KAPLAN, R. C., BARNHART, J., LIU, K., GIACHELLO, A., LEE, D. J., RYAN, J. ET AL. (2010). Sample design and cohort selection in the hispanic community health study/study of latinos. *Annals of epidemiology*, **20** 642–649.
- LIEBERMAN, O. (1994). Saddlepoint approximation for the distribution of a ratio of quadratic forms in normal variables. *Journal of the American Statistical Association*, **89** 924–928.
- SORLIE, P. D., AVILÉS-SANTA, L. M., WASSERTHEIL-SMOLLER, S., KAPLAN, R. C., DAVIGLUS, M. L., GIACHELLO, A. L., SCHNEIDERMAN, N., RAIJ, L., TALAVERA, G., ALLISON, M., LAVANGE, L., CHAMBLESS, L. E. and HEISS, G. (2010). Design and implementation of the hispanic community health study/study of latinos. *Annals of epidemiology*, **20** 629–641.