

ARTICLE

On the reconciliation of missing heritability for genome-wide association studies

Guo-Bo Chen*

The definition of heritability has been unique and clear, but its estimation and estimates vary across studies. Linear mixed model (LMM) and Haseman–Elston (HE) regression analyses are commonly used for estimating heritability from genome-wide association data. This study provides an analytical resolution that can be used to reconcile the differences between LMM and HE in the estimation of heritability given the genetic architecture, which is responsible for these differences. The genetic architecture was classified into three forms via thought experiments: (i) coupling genetic architecture that the quantitative trait loci (QTLs) in the linkage disequilibrium (LD) had a positive covariance; (ii) repulsion genetic architecture that the QTLs in the LD had a negative covariance; (iii) and neutral genetic architecture that the QTLs in the LD had a covariance with a summation of zero. The neutral genetic architecture is so far most embraced, whereas the coupling and the repulsion genetic architecture have not been well investigated. For a quantitative trait under the coupling genetic architecture, HE overestimated the heritability and LMM underestimated the heritability; under the repulsion genetic architecture, HE underestimated but LMM overestimated the heritability for a quantitative trait. These two methods gave identical results under the neutral genetic architecture. A general analytical result for the statistic estimated under HE is given regardless of genetic architecture. In contrast, the performance of LMM remained elusive, such as further depended on the ratio between the sample size and the number of markers, but LMM converged to HE with increased sample size.

European Journal of Human Genetics (2016) 24, 1810–1816; doi:10.1038/ejhg.2016.89; published online 20 July 2016

INTRODUCTION

Heritability has been defined under the context of the multiple regression of an infinitesimal sample size.¹ Recently, various methods have been developed to search for missing heritability in genome-wide association study (GWAS) data, which is typical high-dimensional data ($M \gg N$; ie, the number of markers is larger than the sample size).² Variance component methods, such as the modified Haseman–Elston (HE) regression and the linear mixed model (LMM), estimate heritability as being higher than do single-marker association studies.^{3–5} For a quantitative trait such as height, the estimated heritability was about 0.5 using variance component methods. Compared with the empirical upper bound of the heritability of height,^{6,7} which was thought to be approximately 0.8, the gap of 'missing heritability' has been much narrowed using variance component methods. However, 'true heritability' has not yet been attained. The majority endeavors focused on searching for 'missing heritability' aim to fill the gaps between state-of-the-art (GWAS data; single-marker GWAS provides the lower bound, and variance component GWAS often offers a higher estimate) and old-fashioned designs (such epidemiological data, providing the upper bound of heritability). However, epidemiological and family-based studies are often criticized for potential overestimation of heritability, if not fully justified, due to the shared environment. Similarly, it is unknown whether the variance component will overestimate or underestimate the heritability for GWAS data.

Some researchers showed that LMM provided valid estimate of heritability,⁸ whereas under various genetic architecture the estimated heritability from GWAS data using maximum likelihood framework may be ambiguous.⁹ Recently, Golan *et al.*⁵ and Chen¹⁰ independently

discovered the discrepancy between the estimates from LMM and HE for the estimates of heritability for case-control studies (Chen also found a discrepancy in quantitative traits). Of note, a method by Golan *et al.*⁵ is called phenotype correlation–genotype correlation (PCGC) regression, when without adjustment for covariates PCGC resembles HE. For convenience, we call both of them HE thereby. For more detailed discussion and controversies on the estimation of heritability please refer to Table 1.

Two common issues should be noted. First, the estimation of heritability is most often treated as a statistical procedure: a parameter is estimated and assumed to be heritability as granted. Second, the effect sizes of QTLs are assumed to be from a random distribution. Estimation of heritability can be influenced by the genetic architecture, such as the genomic locations of causal variants/QTLs or the ranges of their effect sizes.¹¹ As heritability is a genetic architecture parameter, it is reasonable to examine how certain forms of genetic architecture, implicitly or explicitly, will influence the estimation of variance component methods. Although little is known about genetic architecture, the estimation of heritability depends on the genetic architecture.¹¹

This study closely scrutinized the genetic architecture without assuming that QTLs were random along the genome, and addressed this implication in the estimation of heritability. As demonstrated in this study, the estimation of heritability depended on the genetic architecture, which can be classified into three forms, underlying a complex trait. Given the various methods proposed for estimating heritability in GWAS data, LMM, which represents a method that is built on maximum likelihood, and HE, which is built on least squares, were studied in detail; they may differ dramatically in

Table 1 A summary of various arguments regarding searching for missing heritability

Author(s)	Methods	Conclusion	
		Quantitative traits	Case-control
Yang <i>et al.</i> ³	LMM	Unbiased	NA
	HE	Unbiased	NA
Lee <i>et al.</i> ⁴	LMM	NA	Unbiased
	HE	NA	NA
Speed <i>et al.</i> ¹³	LMM (weights)	Unbiased	Underestimated
	HE	NA	NA
Lee and Chow ⁸	LMM	Unbiased	NA
		NA	NA
de los Campos <i>et al.</i> ⁹	LMM (a Bayesian version)	Biased	NA
	HE	NA	NA
Golan <i>et al.</i> ⁵	LMM	Unbiased	Underestimate
	HE	Unbiased	Unbiased
Chen. ¹⁰	LMM	Biased (under- or overestimate)	Biased (under- or overestimate)
	HE	Biased (under- or overestimate)	Biased (under- or overestimate)

Abbreviations: HE, modified Haseman–Elston regression; LMM, linear mixed model; NA, if not mentioned in their reports.

their estimations of heritability and reflect the genetic architecture underlying a complex trait. The following conclusions were made: (1) an increased estimate of heritability via variance component methods can be the result of overestimation under the certain forms of genetic architecture; (2) the difference between the estimated heritability in LMM and that in HE may reveal the genetic architecture underlying a complex trait.

MATERIALS AND METHODS

The linear model of a complex trait

For a quantitative trait, under the Hardy–Weinberg equilibrium the additive genetic variance (σ_A^2) is

$$\sigma_A^2 = \sum_{l=1}^L 2p_l q_l \beta_l^2 + \sum_{l_1=1}^L \sum_{l_2 \neq l_1} \rho_{l_1, l_2} \sqrt{2p_{l_1} q_{l_1} 2p_{l_2} q_{l_2}} \beta_{l_1} \beta_{l_2} \quad (1)$$

in which p_l ($l \leq 0.5$) is the allele frequency for the reference allele at the l^{th} QTL, $q_l = 1 - p_l$ the frequency for the alternative allele, ρ_{l_1, l_2} is the correlation measure between the l_1^{th} and the l_2^{th} QTLs, and β_{l_1} is the additive effect for the l_1^{th} locus. Equation 1 is the classic definition of additive genetic variance, referring to page 102 in Lynch and Walsh.¹² For ease of discussion: (i) the phenotype and the genotypes are standardized, and consequently, $h^2 = \sigma_A^2$; (ii) only narrow-sense heritability is discussed here; and (iii) every QTL is perfectly tagged. Due to the context, σ_A^2 and h^2 will be used interchangeably.

The decomposition of genetic architecture

The additive variance component expressed in Equation 1 can be decomposed as

$$\sigma_A^2 = \sigma_{A,w}^2 + \sigma_{A,\bar{w}}^2 = \sum_{l=1}^L \beta_l^2 + \sum_{l_1=1}^L \sum_{l_2 \neq l_1} \rho_{l_1, l_2} \beta_{l_1} \beta_{l_2} \quad (2)$$

$\sum_{l=1}^L \beta_l^2$ is the within-locus variance, denoted as $\sigma_{A,w}^2$, and $\sum_{l_1=1}^L \sum_{l_2 \neq l_1} \rho_{l_1, l_2} \beta_{l_1} \beta_{l_2}$ is the between-locus variance/covariance, denoted as $\sigma_{A,\bar{w}}^2$. Analogously, $h^2 = h_{A,w}^2 + h_{A,\bar{w}}^2$.

The unit of $\sigma_{A,\bar{w}}^2$ is $\rho_{l_1, l_2} \beta_{l_1} \beta_{l_2}$, a three-element product characterizing a pair of QTLs. Given the two possible signs for ρ_{l_1, l_2} , β_{l_1} , and β_{l_2} , it generates eight combinations. Thus, for the ease of discussion of this two-QTL scenario, it is presumed that the reference alleles of two QTLs have been aligned such that $\rho_{l_1, l_2} > 0$ (changing the reference alleles will not change the sign of $\rho_{l_1, l_2} \beta_{l_1} \beta_{l_2}$, Supplementary Note I). If $\rho_{l_1, l_2} \beta_{l_1} \beta_{l_2} > 0$, the genetic unit is called the coupling phase, where QTLs with the same effect sign are clustered together. It is equivalent to argue whether a detected QTL is actually the aggregation of a pair of small-effect QTLs with the same sign. If $\rho_{l_1, l_2} \beta_{l_1} \beta_{l_2} < 0$, it is called the repulsion phase, where QTLs with opposite effect signs are clustered together. It is analogous to argue whether a detected QTL is actually the aggregation of

two QTLs with opposite signs. If $\rho_{l_1, l_2} \beta_{l_1} \beta_{l_2} = 0$, this is the neutral phase, where QTLs are in linkage equilibrium for the two-QTL scenario.

Thus, the smallest genetic architecture in this definition has at least two QTLs, and the total between-locus variance $\sigma_{A,\bar{w}}^2 = \sum_{l_1=1}^L \sum_{l_2 \neq l_1} \rho_{l_1, l_2} \beta_{l_1} \beta_{l_2}$ now can be written as the aggregation of the repulsion and coupling phases along the genome, $\sigma_{A,\bar{w}}^2 = \sum \text{repulsion phase} + \sum \text{coupling phase}$. Depending on if $\sigma_{A,\bar{w}}^2$ is 0, or greater/smaller than 0, the genetic architecture is split into three forms:

- (1) the coupling genetic architecture where $\sigma_{A,\bar{w}}^2 > 0$;
- (2) the repulsion genetic architecture where $\sigma_{A,\bar{w}}^2 < 0$;
- (3) and the neutral genetic architecture where $\sigma_{A,\bar{w}}^2 = 0$.

When the neutral genetic architecture is assumed, the heritability can be simplified as $h_R^2 = h_{A,w}^2 = \sum_{l=1}^L \beta_l^2$. For almost all recent variance component publications,^{3–5,8,9,13} a random distribution of effects, leading to a neutral genetic architecture, along the genome is assumed. Therefore, it is subsequently demonstrated that the coupling/repulsion genetic architecture helps to reconcile the estimation of heritability for GWAS data.

The conventional definition of heritability under the context of multiple regression (h_R^2)

As argued above, the heritability can be split into two components under the multiple regression of infinitesimal sample size

$$h_R^2 = h_{A,w}^2 + h_{A,\bar{w}}^2 \quad (3)$$

Before the availability of GWAS data, h_R^2 is often estimated from epidemiological data via structural equation or linkage analysis.^{6,7} Those estimates are often served as the upper bound in searching 'missing heritability'.

Heritability estimated under LMM (h_{LMM}^2 and $h_{\text{LMM},e}^2$)

In LMM, the variance component of a trait is modeled as

$$\text{var}(y) = A\sigma_A^2 + I\sigma_e^2 = A(\sigma_{A,w}^2 + \sigma_{A,\bar{w}}^2) + I\sigma_e^2$$

Where A is a realized genetic relatedness matrix for samples. Between a pair of individuals i and j , $A_{ij} = \frac{1}{M} \sum_k \frac{(x_{ik} - 2p_k)(x_{jk} - 2p_k)}{2p_k q_k}$, in which M is the number of markers and x_k counts the reference alleles for the k^{th} marker. σ_A^2 can be estimated via the restricted maximum likelihood estimator.^{14,15} It should be noted, as will be shown below, that when $\sigma_{A,\bar{w}}^2 \neq 0$, LMM will give a biased estimate. As discussed by de los Campos *et al.*⁹ the actual estimated statistic, which is often taken as heritability, remains a fundamental question for LMM. h_{LMM}^2 denotes the heritability estimated by LMM, and $h_{\text{LMM},e}^2 = \frac{\sigma_A^2}{\sigma_A^2 + \sigma_e^2}$, in which

$\sigma_A^2 + \sigma_e^2$ is considered a proxy for σ_y^2 . Alternatively, in this study, an *ad hoc* estimate of heritability was also defined as

$$h_{LMM,e}^2 = 1 - \frac{\sigma_e^2}{\sigma_y^2} \quad (4)$$

h_{LMM}^2 and $h_{LMM,e}^2$ might differ across genetic architectures. Other *ad hoc* estimates include introducing weights to A matrix as proposed by Speed *et al.*¹³

Heritability as defined under the HE regression (h_{HE}^2)

Using a modified HE (or PCGC as proposed by Golan *et al.*),^{5,10} the variance component can be modeled as $Y_{ij} = \mu + bA_{ij} + \varepsilon_{ij}$, in which μ is the mean of the model, b is the regression coefficient, $Y_{ij} = (y_i - y_j)^2$ is the squared difference between a pair of individuals, and ε_{ij} is the residual. A is the realized genetic relatedness matrix for samples, as used for LMM above. Chen¹⁰ and Golan *et al.*⁵ both adopted this framework. Golan *et al.* derive the regression coefficient under the assumption that $\sigma_{A,\bar{w}}^2 = 0$; when $\sigma_{A,\bar{w}}^2 \neq 0$, Golan *et al.* did not provide a solution. In Chen's work, the mathematical expectation of the regression coefficient was derived regardless of whether $\sigma_{A,\bar{w}}^2$ was zero or not.¹⁰ After all, Golan *et al.*⁵ took the estimate as heritability directly, whereas Chen found a much broad variation of the estimate and discussed the condition it was, or not was, equal to heritability (h_R^2).

Chen's original work is too long to present here; therefore, only partial results are shown.¹⁰ Without losing generality, the regression coefficient of HE is as follows (see Equation 1 and Table 1 in Chen's original work)

$$b = -2 \frac{\sum_{k=1}^M \sum_{l=1}^L \sum_{kl}^L \rho_{kl} \beta_l \beta_k / M}{\sum_{k=1}^M \sum_{k_2=1}^M \rho_{k_1 k_2}^2 / M^2} \quad (5)$$

The denominator $\sum_{k=1}^M \sum_{k_2=1}^M \rho_{k_1 k_2}^2 / M^2 = \bar{\rho}_{\mathcal{M},\mathcal{M}}^2$ is the averaged linkage disequilibrium (LD) between the markers. The heritability is estimated as $h_{HE}^2 = -\frac{b}{2}$. When there is only one marker in Equation 5, $b = -2\rho_{\mathcal{M},\mathcal{Q}}^2 h_{\mathcal{Q}}^2$, in which $\rho_{\mathcal{M},\mathcal{Q}}$ is the correlation between the marker and the QTL; $h_{\mathcal{Q}}^2$ is the heritability of the QTL. An alternative expression for HE regression coefficient is

$$b = -2 \left\{ \frac{\frac{1}{M} \left\{ \sum_k^M [\sigma_{A,k,w}^2] + \sum_k^M [\sigma_{A,k,\bar{w}}^2] \right\}}{\bar{\rho}_{\mathcal{M},\mathcal{M}}^2} \right\} \quad (6)$$

which decomposes the numerator to the within-locus variance, $\sigma_{A,k,w}^2$, and the between-locus variance, $\sigma_{A,k,\bar{w}}^2$ (Supplementary Note II). Equation 6, which resembles Equation 2, indicates how $\sigma_{A,k,\bar{w}}^2$, the covariance, will influence the estimate. The implications of these two components are not trivial in the inference of genetic architecture for complex traits.

Scenario 1: When it is the neutral genetic architecture, $\sum_k^M [\sigma_{A,k,\bar{w}}^2] = 0$. Equation 5 can be simplified as

$$b = -2\Lambda h_{A,w}^2 \quad (7)$$

in which $\Lambda = \frac{\bar{\rho}_{\mathcal{Q},\mathcal{M}}^2}{\bar{\rho}_{\mathcal{M},\mathcal{M}}^2}$, $\bar{\rho}_{\mathcal{Q},\mathcal{M}}^2$, an unknown parameter, indicates the mean of the LD between a marker and a QTL. When there is no $h_{A,\bar{w}}^2$, $h_{HE}^2 = \frac{b}{-2} = \Lambda h_R^2$. Λ can be 1 when every marker is a QTL, and $h_{HE}^2 = h_R^2$ directly leads to an unbiased estimate of heritability.

Scenario 2: When either the coupling or repulsion genetic architecture is present, the between-locus component contributes to the estimate in Equation 6. $\sum_k^M [\sigma_{A,k,\bar{w}}^2]$ can be positive or negative depending on the underlying genetic architecture. Then there is no simple way to find the heritability (h_R^2) for the trait. As demonstrated in the simulation below, a discrepancy was observed between the respective estimates of LMM and HE.

In addition, weights can be introduced into the HE regression for A_{ij} as proposed by Speed *et al.*¹³ In general, as long as the weights follow a normal distribution, the estimated heritability will be nearly identical to that without weights (Supplementary Note III).

RESULTS

Simulation I: the genetic unit (two QTLs) of the genetic architecture

In order to demonstrate how genetic architecture affects the estimation of heritability, the smallest genetic architecture, which only has two QTLs, was considered first. We simulated 1000 unrelated individuals, and two equally frequent QTLs, which had identical effect sizes, were tagged perfectly on the genome. The heritability was $h_R^2 = 2p_1q_1\beta_1^2 + 2p_2q_2\beta_2^2 + 2\rho_{1,2}\sqrt{2p_1q_1 2p_2q_2}\beta_1\beta_2 = h_{A,w}^2 + h_{A,\bar{w}}^2 = 0.5$. The LD between the pair of two consecutive single-nucleotide polymorphism markers was set to $\rho_{j,j+1} = \{-0.9, -0.8, -0.7, \dots, 0.7, 0.8, 0.9\}$. When $\rho_{j,j+1}$ was positive, it led to the coupling genetic architecture; when $\rho_{j,j+1}$ was negative, it led to the repulsion genetic. The number of genetic markers M was set to 2, 100, 200, 500, 750, 1000, 2000, and 5000, and the allele frequency was 0.5 for each marker. These two QTLs were always located on the $(\frac{M}{2})^{\text{th}}$ and $(\frac{M}{2} + 1)^{\text{th}}$ markers. The genetic relatedness matrix A between the individuals was calculated using all M markers. Heritability was estimated using LMM and HE, respectively. For HE, based on Equation 5, it could be predicted that $E(h_{HE}^2) = \frac{1+\rho_{1,2}}{1+\rho_{1,2}^2} h_R^2$ (Appendix). The analytical results for h_{LMM}^2 was hardly known due to the likelihood, which maximizes to unpredictable maximization if it is not specified correctly.

Figure 1 illustrates the influence of the coupling/repulsion genetic architecture on the estimation of heritability. Under either the coupling or repulsion genetic architecture, neither HE nor LMM gave unbiased estimates of heritability; unbiased estimates were only generated under the neutral genetic architecture (LD = 0), regardless of the number of QTLs and the markers. For HE, the influence of genetic architecture was predictable (Appendix), and the estimated heritability agreed well with $E(h_{HE}^2)$. HE underestimated the true heritability, $\hat{h}_{HE}^2 < 0.5$, under the repulsion genetic architecture; HE overestimated the true heritability, $\hat{h}_{HE}^2 > 0.5$, under the coupling genetic architecture. The whole pattern was consistent for h_{HE}^2 under different M/N ratios. So, the performance of HE should be predictable for the estimation of heritability, at least in the simulated scenarios.

For LMM, the findings were more complicated. The bias of the estimate was not only due to the genetic architecture but also to the ratio of M/N . As observed, when $M/N < 0.5$, \hat{h}_{LMM}^2 overestimated the true heritability under the repulsion genetic architecture, and underestimated the true heritability under the coupling genetic architecture. It seemed that \hat{h}_{LMM}^2 was not influenced by the genetic architecture when $M/N = 0.5$. Nevertheless, \hat{h}_{LMM}^2 changed its response to the genetic architecture when $M/N > 0.5$. However, the number of markers increased and the performance of h_{LMM}^2 converged with $E(h_{HE}^2)$. No known theory can explain the performance of h_{LMM}^2 . The heritability estimated by $h_{LMM,e}^2$ was more precise than that of both h_{LMM}^2 and $E(h_{HE}^2)$ when $M < 500$. When the number of markers was greater than 500, its performance also converged with h_{HE}^2 .

In addition, weights were introduced to generate genetic relatedness,¹³ but a nearly identical patterns for both \hat{h}_{LMM}^2 and \hat{h}_{HE}^2 were observed with weights as without weights (Supplementary Figure S1).

The scenarios for more QTLs were also considered, but the general pattern remained the same for HE and LMM as observed for the two-QTL scenarios (Supplementary Figure S2).

Simulation II: scenarios for case-control data

In previous works by Chen¹⁰ and Golan *et al.*⁵ it was demonstrated that HE (or PCGC) was unbiased in estimating heritability for case-control data. However, that conclusion was incomplete. When the

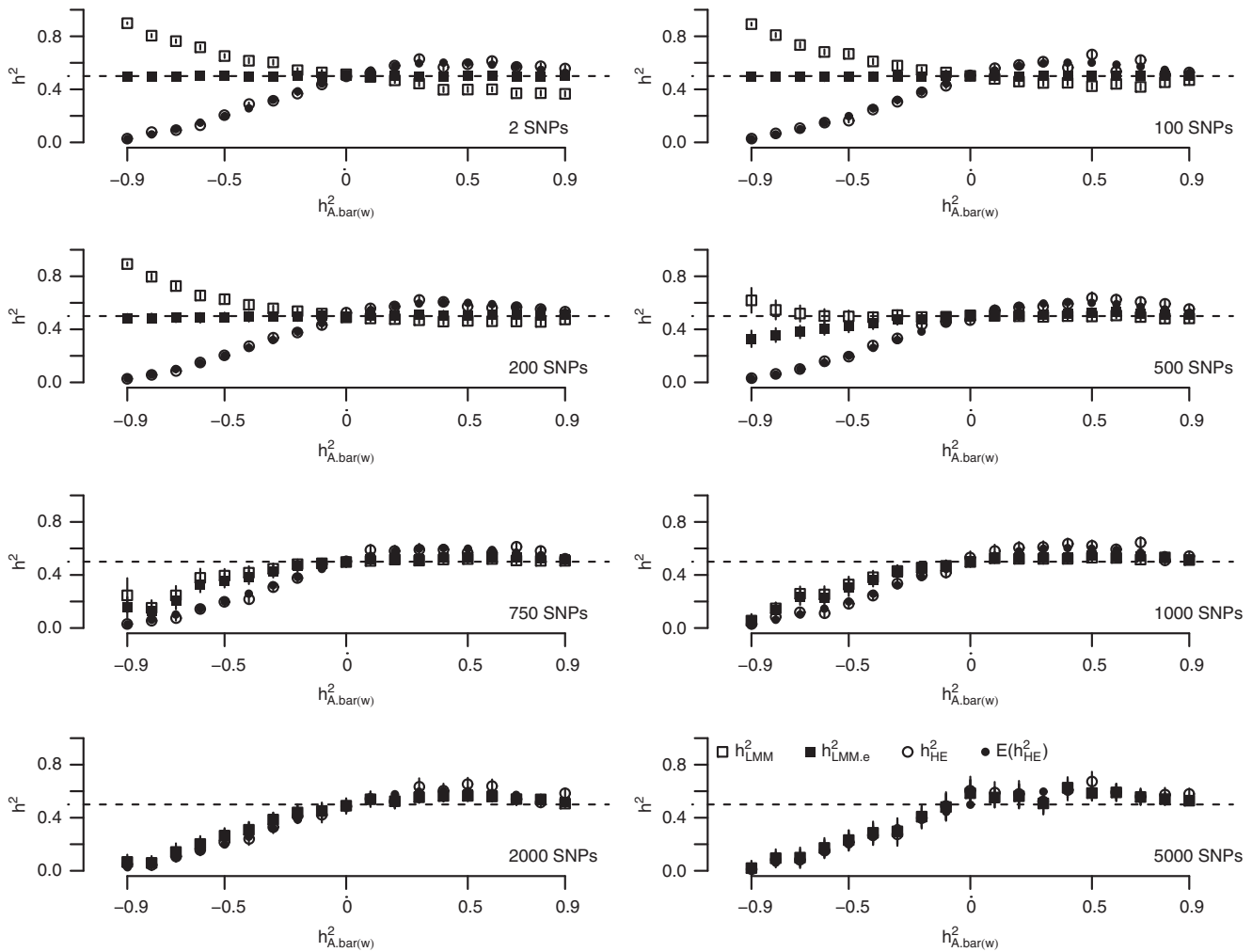


Figure 1 The influence of the genetic architecture on the estimation of heritability from the Haseman–Elston Regression (HE) and the linear mixed model (LMM) – 2 QTLs. The x axis indicates the genetic architecture as quantified by $h^2_{A,\bar{w}}$. $h^2_{A,\bar{w}} < 0$ refers to the repulsion genetic architecture, $h^2_{A,\bar{w}} > 0$ refers to the coupling genetic architecture; and $h^2_{A,\bar{w}} = 0$ refers to the neutral genetic architecture. $E(h^2_{A,\bar{w}})$ can be derived by Equation 5. The SD of each estimate was calculated from 100 replications of the simulations.

base population from which the cases and controls were sampled was characterized by either the coupling or repulsion genetic architecture, HE could also be biased. To demonstrate this phenomenon, 1000 cases and 1000 controls were simulated, $M = \{100, 500, 750, 1000\}$, and equally frequent biallelic QTLs were simulated. To introduce the repulsion and coupling genetic architectures, the effects of the QTLs were sampled from the standard normal distribution, and furthermore, from the first QTL to the last QTL, the effect assumed a quantity of $\Phi^{-1}(P)$. $\Phi^{-1}(P)$ generated a quantity from the normal distribution, given a p -value of P . Here $P = \frac{j}{M}$ for the j^{th} QTL. The LD between two consecutive QTLs was $\rho_{j,j+1} = \{-0.9, -0.8, -0.7, \dots, 0.7, 0.8, 0.9\}$. The total heritability on the liability scale was constrained to 0.5.

As illustrated in Figure 2, after transformation to the liability scale,⁴ a pattern was observed that was similar for quantitative traits: when the base population was under the repulsion genetic architecture, \hat{h}^2_{HE} underestimated the heritability, and when it was under the coupling genetic architecture, \hat{h}^2_{HE} overestimated the heritability. When the base population was under the neutral genetic architecture, \hat{h}^2_{HE} produced an unbiased estimate of the heritability. In contrast, \hat{h}^2_{LMM} depended on both the number of markers and the genetic architecture.

As observed, when $M = 100$ and $K = 0.1$, \hat{h}^2_{LMM} overestimated the heritability when it was under the repulsion genetic architecture, and underestimated the heritability when it was under the coupling genetic architecture. However, the pattern depended upon the number of markers: when the number of markers increased to 1000, \hat{h}^2_{LMM} always underestimated the heritability. $\hat{h}^2_{LMM,e}$ was not as precise in this situation as it was for quantitative traits.

Of note, the crossover between \hat{h}^2_{HE} and \hat{h}^2_{LMM} did not occur at the point where there was neutral genetic architecture, but slightly under the repulsion genetic architecture. This was likely because ascertainment would introduce genetic architecture that resembled the coupling genetic architecture, which is known as Bulmer’s effect in selection studies.¹⁶

Weights were also introduced to the genetic relatedness between individuals,¹³ and the results were nearly identical to those without weighting (Supplementary Figure S3).

Summary: reconciliation of missing heritability in both theory and practice

Table 2 summarizes the theoretical and simulation results presented. The genetic architecture can be classified into coupling, repulsion, and

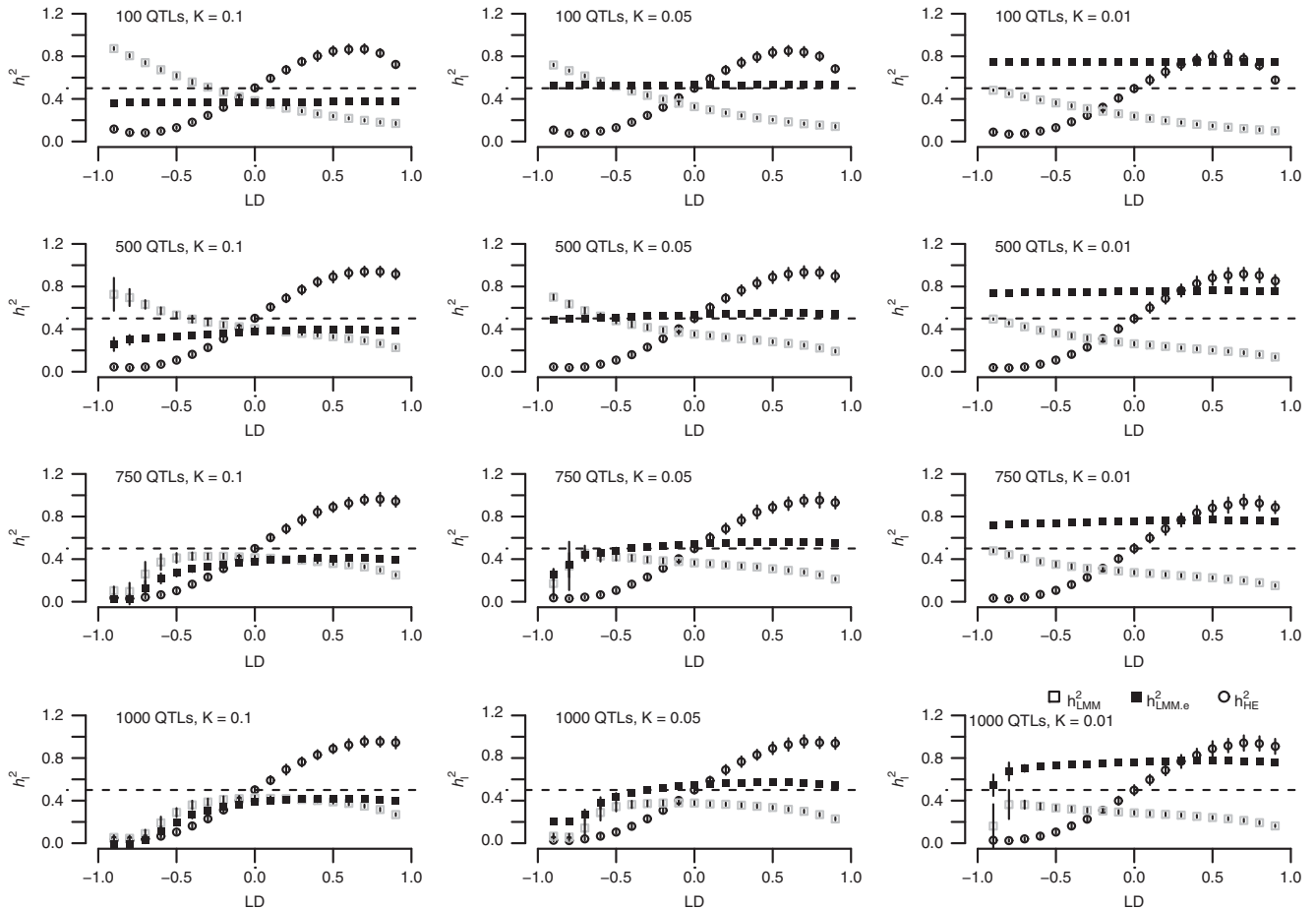


Figure 2 The simulation results for case-control data with varying prevalence (K) and number of QTLs. The x axis quantifies the genetic architecture for the base population, in which the cases and the controls are sampled. The y axis represents the heritability on the liability scale. The x axis reflects the genetic architecture: negative/positive LD indicates the repulsion/coupling genetic architecture; LD=0 indicates the neutral genetic architecture. Of note, when $h_{LMM}^2 = h_{HE}^2$, it was under the repulsion genetic architecture.

neutral genetic architectures. They reflect the physical features of QTLs along the genome, and heritability is a commonly used statistic to summarize this. Depending upon the genetic architecture, multiple regression (h_R^2) – the standard definition for heritability, LMM (h_{LMM}^2), and HE (h_{HE}^2) estimate heritability differently. As heritability is defined under the context of multiple regression, which leads to h_R^2 , it may or may not agree with an alternative heritability estimation, such as h_{LMM}^2 and h_{HE}^2 .

Under the neutral genetic architecture, these three statistics may be closely related. In particular, h_R^2 and h_{HE}^2 are identical, but via different statistical mechanisms, as described in Equation 6. However, under the coupling or repulsion genetic architecture, h_{HE}^2 may under- or overestimate h_R^2 . However, it is not easy to predict the performance of h_{LMM}^2 ; under the neutral genetic architecture, its performance should resemble HE, and very likely converges with the performance of h_{HE}^2 under a wide range of genetic architectures.

In application, the difference between h_{HE}^2 and h_{LMM}^2 , if observed for a trait, may reflect the underlying genetic architecture of that trait. As demonstrated in Simulation I, given the increasing ratio between M/N , h_{LMM}^2 converges with h_{HE}^2 . When these values are close, it does not mean that the estimated heritability is correct; however, this may reflect a condition in which one can presume that the estimated

heritability was likely unbiased. It is unclear whether the convergence is also the case in real data analyses. Further investigation is required to examine how often h_{LMM}^2 converges with h_{HE}^2 , otherwise many reported heritability from LMM remains to be *ad hoc* because of its unwarranted outcomes.

DISCUSSION

As acknowledging the genetic architecture is important for the estimation of heritability, three possible forms of genetic architecture were introduced. Under these three forms, the performance of LMM and HE could be classified. In previous studies, it was suspected that LMM may underestimate the heritability in case-control data,^{5,10} and this study showed that the bias could even occur for quantitative traits. Furthermore, under the coupling genetic architecture, HE overestimated the heritability; under the repulsion genetic architecture, HE underestimated heritability. Under the neutral genetic architecture, HE gave an unbiased estimate. LMM depended on factors other than the genetic architecture, such as the ratio between M and N . Although there was uncertainty in h_{LMM}^2 , an approximation of h_{LMM}^2 could be archived under the neutral genetic architecture. However, as the density of markers can fluctuate the estimation of h_{LMM}^2 , an increased

Table 2 Summary for LMM and HE in the estimation of heritability for genome-wide association data

Genetic architecture ^a	Estimator		
	h_R^2	h_{HE}^2 ^b	h_{LMM}^2 ^c
Neutral	$h_{A,W}^2$	$A h_{A,W}^2$	$h_{LMM}^2 \approx h_{HE}^2 = A h_{A,W}^2$
Positive coupling	$h_{A,W}^2 + h_{A,\bar{W}}^2$	$\frac{\sum_{k=1}^M \sum_{j_1=1}^L \sum_{j_2=1}^L \rho_{k1} \rho_{k2} \beta_{j_1} \beta_{j_2} / M}{\sum_{k_1=1}^M \sum_{k_2=1}^M \rho_{k_1 k_2}^2 / M^2}$	$h_{LMM}^2 \approx h_{HE}^2$
Negative coupling			

Abbreviations: HE, modified Haseman–Elston regression; LMM, linear mixed model. The conclusion was largely based on quantitative traits.
^aNeutral/positive/negative coupling genetic architecture lead to $h_{A,W}^2$ equal/greater/smaller than zero.
^bThe HE regression always has an analytical result regardless of genetic architecture. Under the neutral genetic architecture, it is easy to interpret that $A = \frac{\sigma_{G,W}^2}{\sigma_{A,W}^2}$ (Equation 7), a value between 0 and 1, indicates how well causal variant has been tagged. Under the either positive or negative coupling genetic architecture, the interpretation of the estimate is available in a high-dimension space (Supplementary Note II).
^cLMM does not have close-form result for the estimate. However, when the number of markers was getting larger, its performance was approaching HE.

estimate of heritability may reflect better tagging of QTLs, which may be a good thing, or an overestimation, which is not expected.

These three classes of genetic architecture can be naturally translated into a biological question: how do QTLs emerge on a regional scale and do those nearby QTLs resemble each other or not? A GWAS hit is often an aggregation of much smaller signals, such as those observed in the GIANT height study.¹⁷ In the future, it should be possible to determine whether a region that harbors a GWAS hit actually has more than one signal. The local clustering of QTLs will lead to the repulsion or coupling genetic architecture on the whole-genome scale, but large sample size is required to observe it.

As argued by Bulmer,¹⁶ selection could drive the departure of $\sigma_{A,W}^2$ from zero, and consequently lead to the coupling or repulsion genetic architecture. This raises the question of how likely the repulsion or coupling genetic architecture in real data. A departure from the neutral genetic architecture is indicated when a trait's \hat{h}_{LMM}^2 may differ from its \hat{h}_{HE}^2 . As HE has been under reported in the literature, assessing which genetic architecture is more likely among the three proposed genetic architectures is not possible now. As h_{HE}^2 is easy to implement, testing for genetic architecture forms in various species should be possible,^{5,10,18} particularly among beef cattle or chickens, which are often under strong directional selection and whose traits are also likely under strong selection.

This study used very simple scenarios to demonstrate the genetic architecture and its impact on estimating heritability. Other factors, such as quality control and population structure, may lead to different estimates of heritability using LMM and HE. After all, the current

paradigm used to search for missing heritability favors a higher estimate of heritability; however, one should be careful because a much higher estimate may be an overestimation due to methodological limitations rather than approach the missing heritability.¹⁹

CONFLICT OF INTEREST

The author declares no conflict of interest.

ACKNOWLEDGEMENTS

This article benefited greatly from detailed feedback by the reviewers. I thank Wouter Peyrot and Liyuan Zhou for their helpful discussion. No specific funding was received for this work.

- Falconer DS, Mackay TFC: *Introduction to Quantitative Genetics*. 4th edn. Pearson Education Limited: Harlow, UK, 1996.
- Manolio TA, Collins FS, Cox NJ *et al*: Finding the missing heritability of complex diseases. *Nature* 2009; **461**: 747–753.
- Yang J, Benyamin B, McEvoy BP *et al*: Common SNPs explain a large proportion of the heritability for human height. *Nat Genet* 2010; **42**: 565–569.
- Lee SH, Wray NR, Goddard ME, Visscher PM: Estimating missing heritability for disease from genome-wide association studies. *Am J Hum Genet* 2011; **88**: 294–305.
- Golan D, Lander ES, Rosset S: Measuring missing heritability: Inferring the contribution of common variants. *Proc Natl Acad Sci USA* 2014; **111**: E5272–E5281.
- Visscher PM, Medland SE, Ferreira MAR *et al*: Assumption-free estimation of heritability from genome-wide identity-by-descent sharing between full siblings. *PLoS Genet* 2006; **2**: e41.
- Chen X, Kuja-Halkola R, Rahman I *et al*: Dominant Genetic Variation and Missing Heritability for Human Complex Traits: Insights from Twin versus Genome-wide Common SNP Models. *Am J Hum Genet* 2015; **97**: 708–714.
- Lee JJ, Chow CC: Conditions for the validity of SNP-based heritability estimation. *Hum Genet* 2014; **133**: 1011–1022.
- de los Campos G, Sorensen D, Gianola D: Genomic Heritability: What Is It? *PLoS Genet* 2015; **11**: e1005048.
- Chen G-B: Estimating heritability of complex traits from genome-wide association studies using IBS-based Haseman–Elston regression. *Front Genet* 2014; **5**: 107.
- Moser G, Lee SH, Hayes BJ, Goddard ME, Wray NR, Visscher PM: Simultaneous Discovery, Estimation and Prediction Analysis of Complex Traits Using a Bayesian Mixture Model. *PLoS Genet* 2015; **11**: e1004969.
- Lynch M, Walsh B: *Genetics and Analysis of Quantitative Traits*. Sinauer Associates, Inc.: Sunderland, MA, USA, 1998.
- Speed D, Hemani G, Johnson MR, Balding DJ: Improved Heritability Estimation from Genome-wide SNPs. *Am J Hum Genet* 2012; **91**: 1011–1021.
- Patterson HD, Thompson R: Recovery of inter-block information when block sizes are unequal. *Biometrika* 1971; **58**: 545–554.
- Gilmour AR, Thompson R, Cullis BR: Average information REML: an efficient in linear mixed models variance parameter estimation in linear mixed models. *Biometrics* 1995; **51**: 1440–1450.
- Bulmer MG: The effect of selection on genetic variability. *Am Nat* 1971; **105**: 201–211.
- Yang J, Ferreira T, Morris AP *et al*: Conditional and joint multiple-SNP analysis of GWAS summary statistics identifies additional variants influencing complex traits. *Nat Genet* 2012; **44**: 369–375.
- Hu Z, Yang R-C: Marker-based estimation of genetic parameters in genomics. *PLoS One* 2014; **9**: e102715.
- Yang J, Bakshi A, Zhu Z *et al*: Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body mass index. *Nat Genet* 2015; **47**: 1114–1120.

Supplementary Information accompanies this paper on European Journal of Human Genetics website (<http://www.nature.com/ejhg>)

APPENDIX

Numerical example of a two-QTL scenario

Given $p_1 = p_2 = 0.5$ for both QTLs, which had the same effect size, $\rho_{1,2} \in [-1, 1]$, and the real heritability is $h_R^2 = 2p_1q_1\beta_1^2 + 2p_2q_2\beta_2^2 + 2\rho_{1,2}\sqrt{2p_1q_12p_2q_2}\beta_1\beta_2$. For this case, we set $\beta_1 = \beta_2 = \sqrt{\frac{h_R^2}{(1+\rho_{1,2})}}$.

The numerator of the HE regression coefficient was calculated as $\frac{\sum_{k_1=1}^{\Sigma_1^2} \sum_{k_2=1}^{\Sigma_2^2} \rho_{k_1 k_2} \beta_{k_1} \beta_{k_2}}{2} = \frac{1}{2} \left\{ [\beta_1 \beta_2] \begin{bmatrix} 1 & \rho_{1,2} \\ \rho_{1,2} & \rho_{1,2}^2 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} + [\beta_1 \beta_2] \begin{bmatrix} \rho_{1,2}^2 & \rho_{1,2} \\ \rho_{1,2} & 1 \end{bmatrix} \begin{bmatrix} \beta_1 \\ \beta_2 \end{bmatrix} \right\} = \beta^2 (1 + \rho_{1,2})^2$.

The denominator was calculated as $\bar{\rho}_{\mathcal{M}, \mathcal{M}}^2 = \frac{\sum_{k_1=1}^{\Sigma_1^2} \sum_{k_2=1}^{\Sigma_2^2} \rho_{k_1 k_2}^2}{4} = \frac{1}{4} [1, 1] \begin{bmatrix} 1 & \rho_{1,2}^2 \\ \rho_{1,2}^2 & 1 \end{bmatrix} [1, 1]^T = \frac{1}{2} (1 + \rho_{1,2}^2)$.

Therefore, $b = -\frac{(1+\rho_{1,2})}{\frac{1}{2}(1+\rho_{1,2}^2)} h_R^2 = -2 \frac{1+\rho_{1,2}}{1+\rho_{1,2}^2} h_R^2$. If taking the heritability as the negative half of the regression coefficient, consequently, $E(h_{HE}^2) = \frac{b}{-2} = \frac{1+\rho_{1,2}}{(1+\rho_{1,2}^2)} h_R^2$, which is represented in Figure 1. It means that the statistic estimated from HE may or may not be the heritability. When $\rho_{1,2} = 0$, indicating neutral genetic architecture, HE can provide unbiased estimate of heritability.