

Heritability in the genome-wide association era

Noah Zaitlen · Peter Kraft

Received: 29 February 2012 / Accepted: 29 June 2012 / Published online: 21 July 2012
© Springer-Verlag 2012

Abstract Heritability, the fraction of phenotypic variation explained by genetic variation, has been estimated for many phenotypes in a range of populations, organisms, and time points. The recent development of efficient genotyping and sequencing technology has led researchers to attempt to identify the genetic variants responsible for the genetic component of phenotype directly via GWAS. The gap between the phenotypic variance explained by GWAS results and those estimated from classical heritability methods has been termed the “missing heritability problem”. In this work, we examine modern methods for estimating heritability, which use the genotype and sequence data directly. We discuss them in the context of classical heritability methods, the missing heritability problem, and

describe their implications for understanding the genetic architecture of complex phenotypes.

Introduction

Since their debut in 2005 genome-wide associations studies (GWAS) have identified thousands of single nucleotide polymorphisms (SNPs) associated with hundreds of different phenotypes (Hindorff et al. 2009). Despite this success, the total fraction of the phenotypic variation explained for most phenotypes remains small relative to the published heritability estimates, which are estimated using the trait covariance among relatives (Eichler et al. 2010; Maher 2008; Manolio et al. 2009). This “missing heritability problem” raises questions about the methods used to estimate heritability as well as the genetic architecture of complex phenotypes.

Many explanations for the sources of missing heritability have been proposed including structural variations, gene–environment interactions, epistatic interactions, parent of origin effects, and errors in narrow-sense heritability estimates (Eichler et al. 2010; Manolio et al. 2009; Zuk et al. 2012). Of particular interest is the distribution of causal variants along the genome, their number, and their frequency spectrum. GWAS are particularly suited to capture common variants and so violation of the common disease common variant model may lead to missing heritability. In Fisher’s infinitesimal model, there are expected to be a large number of rare variants associated with disease. The rare-allele model proposes that rare variants of large effect account for a significant fraction of phenotypic variation, and it has been proposed that these can give rise to synthetic association in common variants (Dickson et al. 2010; Gibson 2011).

Determining which combination of these hypotheses is correct and where the majority of phenotypic variation lays

N. Zaitlen · P. Kraft (✉)
Department of Epidemiology, Harvard School of Public Health,
Boston, MA 02115, USA
e-mail: pkraft@hsph.harvard.edu

N. Zaitlen · P. Kraft
Department of Biostatistics, Harvard School of Public Health,
Boston, MA 02115, USA

N. Zaitlen · P. Kraft
Broad Institute of Harvard and Massachusetts Institute
of Technology, Cambridge, MA 02142, USA

N. Zaitlen (✉)
Program in Molecular and Genetic Epidemiology,
Harvard School of Public Health, Boston, MA 02115, USA
e-mail: nzaitlen@hsph.harvard.edu

P. Kraft
Program in Molecular and Genetic Epidemiology,
Harvard School of Public Health, 665 Huntington Avenue,
Building 2 Room 209, Boston, MA 02115, USA

has significant implications for the future success of association studies as well as the clinical utility of genetic risk prediction. It is possible to decouple some of these proposed genetic architectures without directly identifying the causal variants themselves. For example, Wray et al. (2011) show the potential for comparing heritability and sibling relative risk estimates to determine the validity of a rare-variant model (Gibson 2011; Wray and Goddard 2010).

Recently, Yang et al. (2010) proposed using linear-mixed models (LMMs) to estimate a lower bound on the total narrow-sense heritability estimation from GWAS data as well determining how much of the phenotypic variation is due to SNPs in LD with those on genotyping platforms. The results of this approach have broad implications for the genetic architecture of phenotypes as well as the future success of GWAS.

In this work, we examine the problem heritability estimation in the GWAS era and how it relates to the missing heritability problem. We briefly review the classical methods of heritability estimation and contrast them with relatively recent use of genotype data to estimate the component of heritability explained by common SNPs via the LMM approach. We discuss the relative merits of the different methods in terms of potential confounding factors as well as what they tell us about the distribution of causal variants and the potential returns of future GWAS. Finally, we discuss the prospects for using LMM to predict human traits, including disease risk.

Background

Heritability is a measure of the contribution of genetics to phenotype. Wright and Fisher formalized the concept by writing phenotypic variance as the sum of genetic variance and environmental variance, $\sigma_p^2 = \sigma_G^2 + \sigma_e^2$. Broad sense heritability H^2 is the ratio of total genetic variance to phenotypic variance $H^2 = \frac{\sigma_G^2}{\sigma_p^2}$. This measure includes the effects of gene–gene interactions (epistatic effects) σ_I^2 , dominance effects σ_D^2 , and additive effects σ_g^2 such that $\sigma_G^2 = \sigma_g^2 + \sigma_D^2 + \sigma_I^2$. Narrow-sense heritability h^2 measures just the additive contribution of genetic variation to phenotype $h^2 = \frac{\sigma_g^2}{\sigma_p^2}$ (Falconer 1989; Lynch and Walsh 1998).

In this work, we discuss estimates of narrow-sense heritability h^2 unless stated otherwise. This is done because we focus on GWAS and the missing heritability problem. Most traditional estimates of heritability using the correlations among related individuals are presumed to estimate h^2 , although these estimates can be biased. For example, the classical estimate involving the regression of offspring trait values on the mean parental values does not include

the dominance component of variance, but the epistatic component does contribute to the estimate. The epistatic component is typically (and perhaps incorrectly) assumed to be 0 for identifiability purposes (Falconer 1989; Zuk et al. 2012). GWAS estimates of individual-marker effect sizes are generally measured marginally, ignoring dominance and interaction effects, so the “bottom-up” heritability estimates from GWAS (defined below) are narrow-sense estimates.

The additive model

In a GWAS, we are given a set of N_s SNPs $S = \{s_1, s_2, \dots, s_{N_s}\}$ genotyped on N_i individuals with phenotypes $Y = y_1, y_2, \dots, y_{N_i}$. Each genotype has value (0, 1, 2) and the genotypes of the j th individual are $G_j = g_{1j}, g_{2j}, \dots, g_{N_s j}$ with minor allele frequencies p_1, p_2, \dots, p_{N_s} . Let C be the set of N_c causal SNPs, which along with environmental factors determine the phenotype of each individual. The ability of GWAS to identify the genetic contribution to trait variance will depend on the proportion of SNPs in C that are in S or in linkage disequilibrium with one or more SNPs in S .

In an additive model, the phenotype of each individual is defined by a sum of linear effects

$$y_j = m + \sum_{i \in C} z_{ij} \alpha_i + \varepsilon_j \quad (1)$$

where $z_{ij} = \frac{g_{ij} - 2p_i}{\sqrt{2p_i(1-p_i)}}$ are the normalized genotypes, α_i is the effect size of SNP s_i , ε_j is the environmental contribution, and Y is normalized to have variance 1. The environmental contribution is assumed to be normally distributed $\varepsilon_j \sim N(0, \sigma_e^2)$, and ε_j and ε_k are independently distributed for $j \neq k$.

Marginal GWAS “bottom-up” heritability estimation

The genetic variance in an additive model is computed by the sum of the squared-effect sizes of the normalized genotypes $\sigma_g^2 = \sum_i \alpha_i^2$ and the heritability is the ratio of the genetic variance to the total phenotypic variance $h^2 = \frac{\sigma_g^2}{\sigma_g^2 + \sigma_e^2} = \sigma_g^2$, where σ_e^2 is environmental contribution to phenotype and $\sigma_g^2 + \sigma_e^2 = \sigma_Y^2 = 1$.

Given a GWAS, one can compute an estimate of the genetic variance $\hat{\sigma}_g^2$ using the effect size estimates from the markers with a pre-specified genome-wide significance level. This can be used to compute an estimate of the heritability $h_{\text{GWAS}}^2 = \frac{\hat{\sigma}_g^2}{\sigma_Y^2}$, which is defined as “bottom-up” heritability estimation by Zuk et al. (2012).

Unfortunately, the full set of causal variants and their effect sizes are not known, so h_{GWAS}^2 will typically

underestimate the total heritability. The winner’s curse (Ioannidis 2007, 2008; Kraft 2008) and the inclusion of false-positive markers in the bottom-up estimate of genetic variance could in principle lead to an overestimate of heritability. The difference between h^2 and h^2_{GWAS} is known as the “missing heritability”. It is the additive genetic variance not yet captured with GWAS or other methods of identifying associated variants.

Classical “top-down” heritability estimation

The classical methods of heritability estimation are based on an intuitive concept. Phenotypes that are highly correlated among relatives in patterns consistent with Mendelian inheritance are more heritable than those that are weakly correlated among relatives. The formalization of this idea by Fisher (1918) and Wright (1921) is the foundation of heritability estimation.

Consider the correlation between the phenotype of two individuals in the additive model above:

$$\begin{aligned} \text{cor}(y_j, y_k) &= \text{cov}(y_j, y_k) \\ &= \text{cov}\left(\sum_{i \in C} z_{ij} \alpha_i, \sum_{i \in C} z_{ik} \alpha_i\right) = \frac{\sigma_g^2}{N_C} \sum_{i \in C} \frac{\text{cov}(z_{ij}, z_{ik})}{\text{var}(z_i)} \\ &= \sigma_g^2 K_{\text{Causal}}[j, k]. \end{aligned}$$

K_{Causal} is the genetic covariance matrix (Kang et al. 2010; Price et al. 2006; Yang et al. 2010) defined at the causal SNPs. The entry for element j, k in the matrix is:

$$K_{\text{Causal},jk} = \frac{1}{N_C} \sum_{i \in C} \frac{(g_{ij} - 2\hat{p}_i)(g_{ik} - 2\hat{p}_i)}{2\hat{p}_i(1 - \hat{p}_i)}.$$

Until recently, the genotypes of individuals were unavailable and even now the set of causal variants is unknown, so alternative means of estimating K_{Causal} are required. The classical and still widely used approach is to collect sets of related individuals from known pedigrees. The estimate of $K_{\text{Causal},jk}$ is twice the kinship coefficient or $2\Phi_{jk}$. Here Φ_{jk} is the probability that an allele drawn at random from j is identical by descent to a randomly drawn allele from k , and can be calculated from the known pedigree structure (Lange 2002). Many of the familiar values for Φ_{jk} such as $\Phi_{jk} = 1/4$ for full siblings assume that founders share no alleles identical by descent, which may not be true in the presence of inbreeding or population substructure (Lange 2002; Powell et al. 2010). We call the matrix estimated from these pedigree-based estimates K_{Ped} , and it serves as an estimate of K_{Causal} . Given this matrix, the problem of heritability estimation is reduced to estimating $\hat{\sigma}_g^2$ from the observed covariance of the phenotypes of the related individuals.

It is worth stressing that the entries in K_{Ped} are the sums of the expected cross products $E[Z_{ij} \times Z_{ik}]$, while the actual covariance K_{Causal} depends on the observed cross products $z_{ij} \times z_{ik}$. The actual covariance will vary around its expected value for most relative pairs. Visscher et al. (2006) proposed using an estimate of K_{Causal} based on observed genotype data as a more accurate method for estimating heritability using related individuals. For a sample of unrelated sibling pairs, the values of the entries in K_{Ped} are $1/2$ for siblings and 0 otherwise. It follows (making the questionable assumption that the dominance, epistatic, and shared environmental components of variance are 0) that $\hat{\sigma}_g^2$ is twice the average correlation of the phenotype across the sib-pairs (Falconer 1989). If the average correlation among the normalized height of siblings in a population is 0.4 then the heritability estimate for height is 0.8.

Pedigree-based linear-mixed model estimates of heritability

When multiple classes of relationship are measured, as is the case in extended pedigrees, one can take advantage of all the relationships simultaneously via a LMM (Lange 2002; Shaw 1987), where the $1 \times N_{\text{subjects}}$ phenotype vector Y is distributed as a multivariate normal random variable with mean M and variance–covariance matrix Σ . The mean vector M captures the fixed effects of observed covariates (e.g. sex, age, or principal components of genetic variation). The variance–covariance matrix is:

$$\Sigma = \text{var}\left(\sum_{i \in C} z_i \alpha_i\right) + \text{var}(\varepsilon) = K_{\text{Causal}} \sigma_g^2 + I \sigma_\varepsilon^2.$$

To estimate heritability via a LMM, the restricted maximum likelihood (REML) estimate of $\hat{\sigma}_g^2$ is computed and the heritability estimate is $\hat{h}^2 = \frac{\hat{\sigma}_g^2}{\hat{\sigma}_y^2}$. REML is used to estimate the components of variance instead of maximum likelihood to avoid a bias introduced by the fixed effects (Shaw 1987). Since K_{Causal} is not known, K_{Ped} serves as an estimate. There are several algorithms for REML estimation (Kang et al. 2010; Lange 2002; Shaw 1987), but most are computationally expensive due to the cost of matrix inversion. Lippert et al. (2011) recently developed a fast method when the number of individuals exceeds the number of markers. When only one type of relationship is available (e.g. only sibs) then the REML estimator will give the same estimate as the covariance-based approach described above.

There are many extensions to this LMM approach that allow estimation of different components of heritability. These include dominance effects (Lynch and Walsh 1998), gene–gene interaction (Yang et al. 2011a), the shared

genetic basis of multiple phenotypes (cross heritability) (Boehnke et al. 1986; Deary et al. 2012; Lange and Boehnke 1983; Macgregor et al. 2006; Price et al. 2011), heritability from different genomic regions (Yang et al. 2011b), and the effects of shared environment (Lynch and Walsh 1998).

Top-down heritability estimates are susceptible to a range of confounding factors, which can bias estimates. These include gene–environment correlations, selection, non-random mating, and inbreeding (Lynch and Walsh 1998; Visscher et al. 2008). Recently, Zuk et al. (2012) showed that certain types of epistatic interactions can inflate estimates of narrow-sense heritability.

Heritability in the GWAS era

The availability of genotype data over large collections of individuals has opened-up new approaches to estimating heritability. These methods apply the same LMM method described above, but replace the K_{Ped} estimate of K_{Causal} with estimates based on genotype data. We examine only the simple additive estimate of heritability, but each of the extensions listed above may be utilized for each estimator of K_{Causal} .

Heritability using realized IBD

When genetic data are collected over the set of individuals in the study, it is possible to estimate the total fraction of the genome shared identical by descent (IBD). Siblings for example do not share exactly 50 % of their genome with each other (Visscher et al. 2006). Using the genetic data to estimate the fraction of genome shared, IBD gives another means of estimating K_{Causal} , which we call \hat{K}_{IBD} . Provided that the IBD estimates are accurate, this matrix will be a better estimate of K_{Causal} than K_{Ped} and therefore require fewer individuals to achieve a robust estimate of the heritability.

To illustrate this approach, Visscher et al. (2007) used the software package Merlin (Abecasis et al. 2002) to estimate IBD for a collection of twins to generate \hat{K}_{IBD} from 791 autosomal markers and estimate several components of the heritability of height.

Heritability of common variants using observed genetic covariance

Recently, LMMs have been applied to GWAS data in an attempt to partition the “missing” heritability into variants tagged by GWAS SNPs (mostly common) and those that are not (mostly rare) (51). This use of the LMM links modern statistical approaches for high-dimensional data

analysis (penalized regression) with classical models in statistical genetics (de los Campos et al. 2010).

This LMM approach uses the same REML-based estimate of $\hat{\sigma}_g^2$ given above, but the matrix used is an empirical estimate of the genetic covariance (K_{GCV}) instead of \hat{K}_{IBD} or K_{Ped} (Yang et al. 2010, 2011a). This is similar to the “pseudo-heritability” estimate proposed by Kang et al. (2010). The relationship matrix K_{GCV} is computed in a nearly identical way to K_{Causal} , but because the set of causal variants C is unknown, the full set of genotyped SNPs in the GWAS is used directly as a proxy for K_{Causal} .

This approach—which we refer to as the Yang–Visscher or LMM- K_{GCV} approach—relies on the equivalence between the LMM,

$$y_j = \alpha + g_j + \varepsilon_j,$$

with $\text{cov}(g_j, g_k) = K_{\text{GCV},jk} \sigma_g^2$ and $\text{cov}(\varepsilon_j, \varepsilon_k) = 0$, and the random effects model,

$$y_j = \alpha + \sum_{i \in S} \beta_i z_{ij} + \varepsilon_j,$$

with the β_i i.i.d. $N(0, \sigma_g^2/N_s)$. This equivalency provides the motivation for the claim that the LMM using K_{GCV} estimates the proportion of additive genetic variance tagged by the GWAS markers. For unrelated individuals, causal variants that are not correlated with any of the z_{ij} for $i \in S$ (e.g. many rare variants) do not contribute to the estimate of σ_g^2 .

Note that hidden relatedness between individuals would bias σ_g^2 since untagged causal variants would still tend to have the same correlation structure (related to K_{IBD}) as the causal variants that are tagged thereby inflating the estimate of the portion of variability explained by the measured SNPs. Moreover, epistatic effects may also confound estimates of the additive genetic component σ_g^2 if the sample contains closely related individuals. For distantly related individuals, such as those in most GWAS, the effects of dominance and epistasis on heritability estimates are much more attenuated.

The LMM just described is a special case of a general class of regression models defined by any similarity matrix K calculated from the GWAS data, with $\text{cov}(g_j, g_k) = K_{jk} \sigma_g^2$ (de los Campos et al. 2010). Other choices of K may yield improved trait prediction, as they implicitly include non-additive effects (dominance, epistasis). However, precisely because the genetic variance includes non-additive components, the heritability estimates from these more general models can be difficult to interpret.

One of the advantages of this LMM approach using K_{GCV} is that individuals may be selected randomly with respect to their environmental exposures preventing

confounding from shared environments that can affect pedigree-based estimates. In addition, they can inform researchers about the potential success of future GWAS conducted on the phenotype of interest. The K_{GCV} -based estimates provide an upper bound on the total fraction of phenotypic variance explained by future GWAS on the same set of SNPs. They also provide a lower bound on the total narrow-sense heritability of the phenotype.

The application of K_{GCV} to heritability estimation was proposed by Hayes et al. (2009) in the context of related individuals. In this case, K_{KGV} will serve as an estimate of K_{IBD} and therefore give an estimate of the total narrow-sense heritability. The crucial difference in the Yang–Visscher approach is the assumption that when the individuals are distantly related, the K_{GCV} matrix provides no information apart from that contained in the genotyped SNPs. Thus the heritability estimate in this situation will be the narrow-sense heritability due exclusively to the SNPs in LD with those on the genotyping platform.

Violations of model assumptions

Each of the heritability estimation methods described above make different assumptions about the model generating phenotype. The estimates of heritability may be biased when these assumptions are broken.

While pedigree-based estimates of heritability have been examined for decades, the Yang–Visscher approach is a very recent development and there are many open questions about the factors that can affect these estimates of heritability. Here we give several examples of such factors and perform some simple experiments to examine their effects. These are in no way meant to be exhaustive or conclusive, but rather to inform the reader of potential issues.

Violations of additivity

Zuk et al. (2012) show that when certain types of epistatic (gene–gene) interactions exist the estimates of heritability found from pedigree estimates, such as MZ versus DZ twins, will be upwardly biased. In this situation, bottom-up estimates will never reach the top-down estimate of heritability. They propose that this is a possible element of the “missing heritability problem”, and that the true narrow-sense heritability maybe substantially lower than current estimates for certain phenotypes (Zuk et al. 2012).

To examine this problem in the context of Yang–Visscher heritability estimates, we simulated data sets using the epistatic “limiting pathway” models of Zuk et al. (2012), LP(1), LP(3), and LP(4). We simulated case–control genotypes and phenotypes of 2,000 randomly

Table 1 Yang–Visscher and bottom-up estimates of heritability (and their standard error over 1,000 replications) under three limiting pathway models of phenotype

Model	Yang–Visscher	Bottom-up
LP(1)	0.630 (0.012)	0.631 (0.007)
LP(3)	0.398 (0.024)	0.394 (0.024)
LP(4)	0.333 (0.025)	0.333 (0.024)

For $K > 1$ the pedigree-based top-down estimates of heritability will be inflated. An LP(4) model with narrow-sense heritability of 36.4 % will have an estimated heritability of 61.8 % by parent–offspring regression (Zuk et al. 2012). The Yang–Visscher estimate is not affected by this bias

ascertained unrelated individuals with 200 causal variants in each pathway, an effect size of 0.1, a minor allele frequency of 0.5, and prevalence of 50 %. We computed a bottom-up adjusted h^2 estimate via linear regression as well as Yang–Visscher estimate of heritability, using all causal variants to estimate K_{GCV} . The results are shown in Table 1 and demonstrate that the Yang–Visscher approach is not susceptible to confounding from epistatic interaction under the LP model of interaction. If closely related individuals were used then the Yang–Visscher estimate would be upwardly biased from the epistatic component of variance.

Thus, the LMM estimates of heritability from unrelated individuals provide a benchmark to assess how much of the total narrow-sense heritability currently known GWAS-identified trait markers explain—a benchmark that is not influenced by “phantom heritability” due to epistatic interactions. The ratio of the bottom-up additive genetic variance estimated using GWAS-identified markers to the LMM estimate of the additive genetic variance estimates the proportion of GWAS-identifiable markers that have been identified to date.

Violations of exchangeability

The Yang–Visscher approach assumes a polygenic model of disease in which many markers of small effect contribute to variance in genetic risk. Specifically, it assumes marker effect sizes are all drawn from the same normal distribution, $\beta \sim N(0, \sigma_g^2/N_s)$. There are, however, many diseases where there are outlier markers with strikingly different effects. For example, GWAS have identified dozens of markers associated with type 1 diabetes and rheumatoid arthritis, most of which have very small effects relative to the long-established risk variants in the MHC; for both of these diseases, the variants in the MHC have per-allele relative risks roughly three times larger than the relative risks for the GWAS-identified risk variants (Barrett et al. 2009; Stahl et al. 2010).

To examine the effect of such extreme variants, we simulated 1,000 GWAS of 1,500 individuals with a single causal variant. The genotypes at 1,000 marker loci (including the causal locus) were generated by random binomials with minor allele frequencies drawn uniformly between 0.05 and 0.5. The true heritability of the phenotype was 0.5 and the average estimate over the 1,000 GWAS was 0.50, suggesting that violations of the infinitesimal model do not strongly effect estimates of heritability.

Addition of non-causal variants

For many phenotypes, K_{GCV} will contain a large number of variants unlinked to any causal variants. To examine the effect of these variants on the estimates of heritability, we repeated the experiment above with 10 causal variants and 10^2 , 10^3 , and 5×10^3 additional independent (i.e. non-causal) variants. The true heritability of the phenotype was 0.5 and the mean heritability across the 1,000 simulated GWAS was 0.50 in all studies. However, the standard deviations were 0.018, 0.025, and 0.067, showing that the effect of additional variants is to increase standard error of the heritability estimates. The results did not change qualitatively for other values of h^2 . Other factors that affect the standard error of heritability estimates include the study sample size as well as the true heritability. Alternative disease models, such as mixtures of infinitesimals, described by Park et al. (2011) have not yet been investigated in this context; the possibility that they lead to biased heritability estimates remains open.

Sample size considerations

To investigate the precision of LMM estimates of h^2 using K_{GCV} in real-world situations, we used GWAS data on 10,503 individuals from two European-ancestry cohorts, the Nurses' Health Study and Health Professionals Follow-up Study. We simulated continuous phenotypes as a function of 500 SNPs, according to Eq. 1, constraining the SNP effects so that the resulting phenotype had the desired heritability ($h^2 = 0.50, 0.25, 0.10$). We also simulated a binary phenotype using the liability threshold model, with liability given by Eq. 1, and prevalence 10%. We estimated h^2 using the LMM approach as implemented in GCTA (Yang et al. 2011a), applied to a set of 151,019 markers (including the 500 causal variants) chosen to have low linkage disequilibrium ($r^2 < 0.2$), varying the sample size from 1,000 to 10,503.

Results from single replicates are shown in Table 2. Precision increases roughly linearly with increasing log sample size. For sample sizes under 2,000, the 95% confidence intervals are wide (>0.40), and, for modest

Table 2 LMM estimates of narrow-sense heritability using K_{GCV} and their standard errors for phenotypes simulated conditional on empirical GWAS data (described in text, under “Sample size considerations”)

Sample size	True h^2		
	0.10	0.25	0.50
Continuous phenotype			
1,000	0.000 (0.167)	0.228 (0.229)	0.773 (0.124)
1,999	0.141 (0.112)	0.207 (0.118)	0.567 (0.104)
3,993	0.079 (0.059)	0.302 (0.061)	0.631 (0.050)
7,989	0.104 (0.031)	0.0297 (0.031)	0.594 (0.027)
10,503	0.136 (0.024)	0.321 (0.025)	0.583 (0.021)
Binary phenotype			
2,099	0.111 (0.111)	0.314 (0.125)	0.414 (0.108)
10,503	0.125 (0.069)	0.224 (0.070)	0.648 (0.075)

For binary phenotypes, heritabilities are on the liability scale, calculated using the transformation described in the section “Ascertainment and case-control phenotypes”

heritabilities (under 25%, consistent with the observed heritabilities for many complex traits), they include 0. This suggests that accurate estimation of narrow-sense heritabilities will require large sample sizes, on the order of 5,000–10,000 or more, at least as big as those needed to identify individual markers with modest effects. Published studies using the LMM- K_{GCV} approach to estimate the narrow-sense heritability due to GWAS markers for continuous traits like height and body mass index used between 4,000 and 11,500 subjects (Yang et al. 2010, 2011b). Care must be taken when combining studies to reach such large sample sizes, as this may introduce population substructure and corresponding environmental variation of non-genetic risk factors, potentially biasing estimates of heritability.

Addition of markers in LD with the causal variants

The additive model assumes that all of the tested variants are independent. In reality, there is extensive LD between causal and non-causal variants in the genome. To examine the potential for LD to affect heritability estimates, we repeated the experiment above with 4 causal variants, and 1 additional causal variant repeated 100 times simulating extensive LD for a particular SNP, and 10^4 non-causal variants. The true heritability was 0.5 and the average estimated heritability was 0.40 showing that LD patterns can significantly affect heritability estimates. We note that this is an extreme example meant to demonstrate the potential for bias. Yang et al. (2010) simulated phenotypes over real GWAS data (i.e. with real LD patterns) and found estimates within two standard errors of the true heritability.

Table 3 The genetic covariance between pairs of individuals with a range of IBDs, estimate from N_s SNPs

$N_s \backslash IBD$	0.025	0.05	0.1	0.25	0.5
1,000	0.025 (0.032)	0.05 (0.032)	0.10 (0.033)	0.25 (0.034)	0.50 (0.039)
10,000	0.025 (0.010)	0.05 (0.010)	0.10 (0.011)	0.25 (0.011)	0.50 (0.012)
100,000	0.025 (0.001)	0.05 (0.001)	0.10 (0.001)	0.25 (0.001)	0.50 (0.001)

K_{GCV} is an unbiased estimate of IBD and the variance of the estimate in parenthesis is function of the number of available SNPs

Distant/cryptic relatedness in the study

Provided that the individuals in a GWAS are unrelated, the matrix K_{GCV} contains no information about SNPs out of LD with the genotyped SNPs. If the study contains related individuals, however, the LMM estimate of heritability will contain some additional genetic variance due to variants not tagged by the GWAS SNPs. This is because K_{GCV} is an unbiased estimate of K_{IBD} , as we illustrate below. Since there are no truly “unrelated” individuals, any GWAS will contain a range of distantly related individuals. Yang et al. (2011a) suggest removing individuals with $K_{GCV} > 0.025$ in the case of quantitative phenotypes and 0.05 for dichotomous phenotypes.

We simulated 1,000 pairs of individuals that shared 0.5, 0.1, 0.05 and 0.025 of their genome IBD and compute K_{GCV} for each pair. We repeated this experiment using 10^4 , 10^5 , and 10^6 SNPs. The results are presented in Table 3. In each case, the mean estimate of IBD is close to the true IBD showing that the K_{GCV} is a good estimate of K_{IBD} . The standard error is independent of the true IBD and decreases as a function of the number of independent SNPs.

For distantly related individuals, the signal from IBD will typically be small relative to the signal from the causal variants. Here, a concern is confounding due to cryptic relatedness, where more closely related individuals tend to have similar trait values for non-genetic reasons (Kang et al. 2010). The influence of low levels of IBD in the Yang–Visscher approach remains an open question. It is possible to test explicitly for inflation due to relatedness, by simulating phenotypes over odd chromosomes and estimating heritability over even chromosomes (Visscher et al. 2010).

Population substructure

Individuals from different populations have different minor allele frequencies as well different environmental exposures. In a case–control study, this can lead to significant confounding if there is a difference in the phenotypic mean between the populations, and is usually corrected with a principal component adjustment. Browning and Browning (2011) show that under certain extreme population differences, this can lead to biases in heritability estimates. Yang et al. (2011b) show that using PC adjustment will mitigate

this inflation. They also propose to estimate the effects of population stratification and cryptic relatedness by performing heritability estimation over each chromosome. This procedure has not yet been examined in detail in the published literature.

Another type of population stratification arises when there is a difference in the phenotypic variance (but not necessarily mean phenotype) between the populations. In this case, PCA will not adequately adjust for population substructure leading to inflation in standard GWAS (McPeck and Abney 2008). Furthermore, the interpretation of heritability may be ambiguous in this scenario, since each of the sub populations will likely have different heritability estimates.

Heritability is defined with respect to a population at a particular time. The heritability of lung cancer will be dramatically different between a population where some people smoke and a population of only non-smokers. Thus bottom-up GWAS heritability estimates and those from published heritability studies can only be compared if they come from the same population and are conducted at similar times.

Imputation and rare variants

Currently the Yang–Visscher approach has been performed using observed SNPs, genotyped using the same platform (Yang et al. 2010, 2011b). Given the success of imputation in the GWAS community, one of the open questions is the possibility of leveraging external reference panels such as the HapMap to determine if additional signal lies within the additional SNPs genotypes in the panel. High-throughput sequencing data are available with a large number of rare variants. The proper way to include dense maps of common markers and rare variants in heritability estimation—notably in light of the discussion of the impact of linkage disequilibrium patterns, above—is an area of current research.

Ascertainment and case–control phenotypes

For binary traits, the percent of trait variance captured by K_{GCV} when analyzing a discontinuous 1–0 case–control phenotype in the LMM framework is not directly

comparable to commonly quoted heritabilities from some family studies (e.g. MZ–DZ twin comparisons), which are measures of the percentage of the underlying liability captured by inherited factors (Dempster and Lerner 1949; Visscher et al. 2008). Nor is there a simple link between individual-locus odds ratios and bottom-up estimates of the genetic variance and liability-scale heritability. There is a simple relationship between the familial recurrence risk and the additive genetic component of variance on the log relative risk scale (Pharoah et al. 2002). Moreover, individual-marker allele frequencies and relative risk estimates from GWAS can be directly related to heritability on the log relative risk scale (Pharoah et al. 2008). See Wray and Goddard (2010) for a thorough discussion of the relationship between individual-marker relative risks and heritability measured on different scales.

For the LMM, the phenotypic variance captured by K_{GCV} depends on disease prevalence and sampling scheme. By construction, the heritability of liability is independent of prevalence. When estimating the heritability of case–control phenotypes, the ascertainment strategy and prevalence of disease will affect the final heritability estimate. To address this issue, it is possible to transform the disease scale heritability estimate to a liability scale heritability estimate, which accounts for both ascertainment and prevalence (Dempster and Lerner 1949; Lee et al. 2011):

$$h_{\text{liability}}^2 = h_{\text{Obs}}^2 \frac{F(1-F)}{\phi(\Phi^{-1}(F))^2} \frac{F(1-F)}{P(1-P)}$$

F is the prevalence, ϕ is the normal pdf, Φ is the normal cdf and P is the proportion of cases in the sample. The justification for this elegant adjustment depends on a rather simple model for ascertainment, namely, that selection for inclusion is independent of all other covariates conditional on disease status. This will not be the case in many practical situations (e.g. matched case–control studies), where ascertainment depends on other factors that are usually associated with disease risk and may also be associated with genotype. The impact of violations of this assumption is unclear.

Phenotypic prediction

The LMM using K_{GCV} also offers a means of phenotypic prediction using the best linear unbiased predictors or BLUPs (Lynch and Walsh 1998). The expected trait value for a new individual (who did not contribute to the data set used to fit the LMM) is given by:

$$\hat{y} = \alpha + \hat{g} = \alpha + \sum_{i \in S} \hat{\beta}_i z_i$$

This is similar to the “polygenic” models proposed by Purcell et al. (2009) and Evans et al. (2009), in that the

predictor uses information contained in SNPs that do not reach the genome-wide significance threshold. But where the “polygenic model” performs feature selection only building predictors using markers with single-SNP (marginal) p values below some threshold (often much larger than the stringent GWAS threshold), the LMM approach builds predictors using all available SNPs simultaneously. The LMM predictor is closely related to ridge regression, a penalized regression procedure that often outperforms variable selection procedures in terms of minimizing prediction error in new data sets (Harrell 2001; Hastie et al. 2001).

The accuracy of the LMM predictor is a function of narrow-sense heritability, the number of markers included in the LMM, the true genetic architecture, and the sample size in the data set used to fit the LMM. The sample size determines the accuracy with which β_i can be estimated. The squared correlation between the LMM predictor and trait values in new observations is typically far smaller than the heritability estimate from the LMM (the theoretical maximum of the squared correlation); this is because of the variability in the estimated β_i s (Daetwyler et al. 2008; Visscher et al. 2010).

Conclusion

The Yang–Visscher approach to heritability estimation provides a means of estimating the contribution of SNPs in LD with those on genotyping platforms to the total phenotypic variation. In the context of GWAS, these estimates answer questions about the genetic architecture of complex phenotypes. The growing number of GWAS identified loci, as well as their small effect sizes, has led to speculation about genetic models of disease.

There has been significant recent debate about the success or failure of GWAS (Eichler et al. 2010; Gibson 2011; Visscher et al. 2012). This has in turn reinvigorated the debate about the distribution of causal variants. Goldstein demonstrated the possibility for rare variants to induce synthetic associations (Dickson et al. 2010), and there have been several recent works discussing the common disease common variant, strong and weak rare variants, the infinitesimal, and other disease models (Gibson 2011).

There has also been speculation about the location of the “missing heritability” with discussions of parent of origin effects, epistatic interactions, gene–environment interactions, structural variation, and other cache’s of genetic variation not well captured by current GWAS or their analysis methods (Eichler et al. 2010; Visscher et al. 2012; Zuk et al. 2012).

The work of Yang and Visscher discussed here as well as other GWAS-based approaches (Lango Allen et al.

2010; So et al. 2011a, 2011b; Yang et al. 2011c) provide insights relevant to these questions. They estimate heritability restricted to a certain class of SNPs (i.e. those in LD with genotyped SNPs), are not confounded by many of the factors biasing traditional methods of heritability estimation, and are fundamentally different than bottom-up methods. In principle, these procedures could also be used to build phenotype prediction algorithms incorporating markers beyond the small number identified at genome-wide significance levels. However, very large sample sizes will be needed to obtain accurate estimates and precise prediction algorithms.

Acknowledgments The authors thank Alkes Price, Eli Stahl and Dan Stram for helpful discussions, and Poorva Mudgal for programming support. NZ was supported by NIH fellowship 5T32ES007142-27, PK by NIH grant R21 DK084529.

References

- Abecasis GR, Cherny SS, Cookson WO, Cardon LR (2002) Merlin—rapid analysis of dense genetic maps using sparse gene flow trees. *Nat Genet* 30:97–101
- Barrett JC, Clayton DG, Concannon P, Akolkar B, Cooper JD, Erlich HA, Julier C, Morahan G, Nerup J, Nierras C, Plagnol V, Pociot F, Schuilenburg H, Smyth DJ, Stevens H, Todd JA, Walker NM, Rich SS (2009) Genome-wide association study and meta-analysis find that over 40 loci affect risk of type 1 diabetes. *Nat Genet* 41:703–707
- Boehnke M, Moll PP, Lange K, Weidman WH, Kottke BA (1986) Univariate and bivariate analyses of cholesterol and triglyceride levels in pedigrees. *Am J Med Genet* 23:775–792
- Browning SR, Browning BL (2011) Population structure can inflate SNP-based heritability estimates. *Am J Hum Genet* 89:191–193 (author reply 193–195)
- Daetwyler HD, Villanueva B, Woolliams JA (2008) Accuracy of predicting the genetic risk of disease using a genome-wide approach. *PLoS One* 3:e3395
- De los Campos G, Gianola D, Allison DB (2010) Predicting genetic predisposition in humans: the promise of whole-genome markers. *Nat Rev Genet* 11:880–886
- Deary IJ, Yang J, Davies G, Harris SE, Tenesa A, Liewald D, Luciano M, Lopez LM, Gow AJ, Corley J, Redmond P, Fox HC, Rowe SJ, Haggarty P, McNeill G, Goddard ME, Porteous DJ, Whalley LJ, Starr JM, Visscher PM (2012) Genetic contributions to stability and change in intelligence from childhood to old age. *Nature* 482(7384):212–215. doi:10.1038/nature10781
- Dempster E, Lerner I (1949) Heritability of threshold characters. *Genetics* 35:212–236
- Dickson SP, Wang K, Krantz I, Hakonarson H, Goldstein DB (2010) Rare variants create synthetic genome-wide associations. *PLoS Biol* 8:e1000294
- Eichler EE, Flint J, Gibson G, Kong A, Leal SM, Moore JH, Nadeau JH (2010) Missing heritability and strategies for finding the underlying causes of complex disease. *Nat Rev Genet* 11:446–450
- Evans DM, Visscher PM, Wray NR (2009) Harnessing the information contained within genome-wide association studies to improve individual prediction of complex disease risk. *Hum Mol Genet* 18:3525–3531
- Falconer DS (1989) Introduction to quantitative genetics, 3rd edn. Wiley, Burnt Mill
- Fisher R (1918) The correlation among relatives on the supposition of Mendelian inheritance. *Trans Roy Soc Edinburgh* 52:399–433
- Gibson G (2011) Rare and common variants: twenty arguments. *Nat Rev Genet* 13:135–145
- Harrell F (2001) Regression modeling strategies. Springer, New York
- Hastie T, Tibshirani R, Friedman J (2001) The elements of statistical learning. Springer, New York
- Hayes BJ, Visscher PM, Goddard ME (2009) Increased accuracy of artificial selection by using the realized relationship matrix. *Genet Res (Camb)* 91:47–60
- Hindorf LA, Sethupathy P, Junkins HA, Ramos EM, Mehta JP, Collins FS, Manolio TA (2009) Potential etiologic and functional implications of genome-wide association loci for human diseases and traits. *Proc Natl Acad Sci USA* 106:9362–9367
- Ioannidis JP (2007) Non-replication and inconsistency in the genome-wide association setting. *Hum Hered* 64:203–213
- Ioannidis JP (2008) Why most discovered true associations are inflated. *Epidemiology* 19:640–648
- Kang HM, Sul JH, Service SK, Zaitlen NA, Kong SY, Freimer NB, Sabatti C, Eskin E (2010) Variance component model to account for sample structure in genome-wide association studies. *Nat Genet* 42:348–354
- Kraft P (2008) Curses—winner’s and otherwise—in genetic epidemiology. *Epidemiology* 19:649–651 (discussion 657–658)
- Lange K (2002) Mathematical and statistical methods for genetic analysis. Springer, New York
- Lange K, Boehnke M (1983) Extensions to pedigree analysis. IV. Covariance components models for multivariate traits. *Am J Med Genet* 14:513–524
- Lango Allen H, Estrada K, Lettre G, Berndt SI, Weedon MN, Rivadeneira F, Willer CJ, Jackson AU, Vedantam S, Raychaudhuri S, Ferreira T, Wood AR, Weyant RJ, Segre AV, Speliotes EK, Wheeler E, Soranzo N, Park JH, Yang J, Gudbjartsson D, Heard-Costa NL, Randall JC, Qi L, Vernon Smith A, Magi R, Pastinen T, Liang L, Heid IM, Luan J, Thorleifsson G, Winkler TW, Goddard ME, Sin Lo K, Palmer C, Workalemahu T, Aulchenko YS, Johansson A, Zillikens MC, Feitosa MF, Esko T, Johnson T, Ketkar S, Kraft P, Mangino M, Prokopenko I, Absher D, Albrecht E, Ernst F, Glazer NL, Hayward C, Hottenga JJ, Jacobs KB, Knowles JW, Kutalik Z, Monda KL, Polasek O, Preuss M, Rayner NW, Robertson NR, Steinthorsdottir V, Tyrer JP, Voight BF, Wiklund F, Xu J, Hua Zhao J, Nyholt DR, Pellikka N, Perola M, Perry JR, Surakka I, Tammesoo ML, Altmaier EL, Amin N, Aspelund T, Bhangale T, Boucher G, Chasman DI, Chen C, Coin L, Cooper MN, Dixon AL, Gibson Q, Grundberg E, Hao K, Juhani Juntila M, Kaplan LM, Kettunen J, König IR, Kwan T, Lawrence RW, Levinson DF, Lorentzon M, McKnight B, Morris AP, Muller M, Suh Ngwa J, Purcell S, Rafelt S, Salem RM, Salvi E et al (2010) Hundreds of variants clustered in genomic loci and biological pathways affect human height. *Nature* 467:832–838
- Lee SH, Wray NR, Goddard ME, Visscher PM (2011) Estimating missing heritability for disease from genome-wide association studies. *Am J Hum Genet* 88:294–305
- Lippert C, Listgarten J, Liu Y, Kadie CM, Davidson RI, Heckerman D (2011) FaST linear mixed models for genome-wide association studies. *Nat Methods* 8:833–835
- Lynch M, Walsh B (1998) Genetics and analysis of quantitative traits. Sinauer, Sunderland
- Macgregor S, Cornes BK, Martin NG, Visscher PM (2006) Bias, precision and heritability of self-reported and clinically measured height in Australian twins. *Hum Genet* 120:571–580
- Maher B (2008) Personal genomes: the case of the missing heritability. *Nature* 456:18–21

- Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorf LA, Hunter DJ, McCarthy MI, Ramos EM, Cardon LR, Chakravarti A, Cho JH, Guttmacher AE, Kong A, Kruglyak L, Mardis E, Rotimi CN, Slatkin M, Valle D, Whittemore AS, Boehnke M, Clark AG, Eichler EE, Gibson G, Haines JL, Mackay TF, McCarroll SA, Visscher PM (2009) Finding the missing heritability of complex diseases. *Nature* 461:747–753
- Park JH, Gail MH, Weinberg CR, Carroll RJ, Chung CC, Wang Z, Chanock SJ, Fraumeni JF Jr, Chatterjee N (2011) Distribution of allele frequencies and effect sizes and their interrelationships for common genetic susceptibility variants. *Proc Natl Acad Sci USA* 108:18026–18031
- Pharoah PD, Antoniou A, Bobrow M, Zimmern RL, Easton DF, Ponder BA (2002) Polygenic susceptibility to breast cancer and implications for prevention. *Nat Genet* 31:33–36
- Pharoah PD, Antoniou AC, Easton DF, Ponder BA (2008) Polygenes, risk prediction, and targeted prevention of breast cancer. *N Engl J Med* 358:2796–2803
- Powell JE, Visscher PM, Goddard ME (2010) Reconciling the analysis of IBD and IBS in complex trait studies. *Nat Rev Genet* 11(11):800–805
- Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D (2006) Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet* 38:904–909
- Price AL, Helgason A, Thorleifsson G, McCarroll SA, Kong A, Stefansson K (2011) Single-tissue and cross-tissue heritability of gene expression via identity-by-descent in related or unrelated individuals. *PLoS Genet* 7:e1001317
- Purcell SM, Wray NR, Stone JL, Visscher PM, O'Donovan MC, Sullivan PF, Sklar P (2009) Common polygenic variation contributes to risk of schizophrenia and bipolar disorder. *Nature* 460:748–752
- Shaw R (1987) Maximum-likelihood approaches applied to quantitative genetics of natural populations. *Evolution* 41:812–826
- So HC, Gui AH, Cherny SS, Sham PC (2011a) Evaluating the heritability explained by known susceptibility variants: a survey of ten complex diseases. *Genet Epidemiol* 35:310–317
- So HC, Li M, Sham PC (2011b) Uncovering the total heritability explained by all true susceptibility variants in a genome-wide association study. *Genet Epidemiol* 35:447–456
- Stahl EA, Raychaudhuri S, Remmers EF, Xie G, Eyre S, Thomson BP, Li Y, Kurreeman FA, Zhernakova A, Hinks A, Guiducci C, Chen R, Alfredsson L, Amos CI, Ardlie KG, Barton A, Bowes J, Brouwer E, Burt NP, Catanese JJ, Coblyn J, Coenen MJ, Costenbader KH, Criswell LA, Crusius JB, Cui J, de Bakker PI, De Jager PL, Ding B, Emery P, Flynn E, Harrison P, Hocking LJ, Huizinga TW, Kastner DL, Ke X, Lee AT, Liu X, Martin P, Morgan AW, Padyukov L, Posthumus MD, Radstake TR, Reid DM, Seielstad M, Seldin MF, Shadick NA, Steer S, Tak PP, Thomson W, van der Helm-van Mil AH, van der Horst-Bruinsma IE, van der Schoot CE, van Riel PL, Weinblatt ME, Wilson AG, Wolbink GJ, Wordsworth BP, Wijmenga C, Karlson EW, Toes RE, de Vries N, Begovich AB, Worthington J, Siminovitch KA, Gregersen PK, Klareskog L, Plenge RM (2010) Genome-wide association study meta-analysis identifies seven new rheumatoid arthritis risk loci. *Nat Genet* 42:508–514
- Visscher PM, Medland SE, Ferreira MA, Morley KI, Zhu G, Cornes BK, Montgomery GW, Martin NG (2006) Assumption-free estimation of heritability from genome-wide identity-by-descent sharing between full siblings. *PLoS Genet* 2:e41
- Visscher PM, Macgregor S, Benyamin B, Zhu G, Gordon S, Medland S, Hill WG, Hottenga JJ, Willemsen G, Boomsma DI, Liu YZ, Deng HW, Montgomery GW, Martin NG (2007) Genome partitioning of genetic variation for height from 11,214 sibling pairs. *Am J Hum Genet* 81:1104–1110
- Visscher PM, Hill WG, Wray NR (2008) Heritability in the genomics era—concepts and misconceptions. *Nat Rev Genet* 9:255–266
- Visscher PM, Yang J, Goddard ME (2010) A commentary on ‘common SNPs explain a large proportion of the heritability for human height’ by Yang et al. (2010). *Twin Res Hum Genet* 13:517–524
- Visscher PM, Brown MA, McCarthy MI, Yang J (2012) Five years of GWAS discovery. *Am J Hum Genet* 90:7–24
- Wray NR, Goddard ME (2010) Multi-locus models of genetic risk of disease. *Genome Med* 2:10
- Wray NR, Purcell SM, Visscher PM (2011) Synthetic associations created by rare variants do not explain most GWAS results. *PLoS Biol* 9:e1000579
- Wright S (1921) Systems of mating. *Genetics* 6:111–178
- Yang J, Benyamin B, McEvoy BP, Gordon S, Henders AK, Nyholt DR, Madden PA, Heath AC, Martin NG, Montgomery GW, Goddard ME, Visscher PM (2010) Common SNPs explain a large proportion of the heritability for human height. *Nat Genet* 42:565–569
- Yang J, Lee SH, Goddard ME, Visscher PM (2011a) GCTA: a tool for genome-wide complex trait analysis. *Am J Hum Genet* 88:76–82
- Yang J, Manolio TA, Pasquale LR, Boerwinkle E, Caporaso N, Cunningham JM, de Andrade M, Feenstra B, Feingold E, Hayes MG, Hill WG, Landi MT, Alonso A, Lettre G, Lin P, Ling H, Lowe W, Mathias RA, Melbye M, Pugh E, Cornelis MC, Weir BS, Goddard ME, Visscher PM (2011b) Genome partitioning of genetic variation for complex traits using common SNPs. *Nat Genet* 43:519–525
- Yang J, Weedon MN, Purcell S, Lettre G, Estrada K, Willer CJ, Smith AV, Ingelsson E, O'Connell JR, Mangino M, Magi R, Madden PA, Heath AC, Nyholt DR, Martin NG, Montgomery GW, Frayling TM, Hirschhorn JN, McCarthy MI, Goddard ME, Visscher PM (2011c) Genomic inflation factors under polygenic inheritance. *Eur J Hum Genet* 19:807–812
- Zuk O, Hechter E, Sunyaev SR, Lander ES (2012) The mystery of missing heritability: genetic interactions create phantom heritability. *Proc Natl Acad Sci USA* 109(4):1193–1198