

Quasi-regression for heritability

Art B. Owen
Stanford University

March 2012

Abstract

We show in an idealized model that the narrow sense (linear) heritability from d autosomal SNPs can be estimated without bias from n independently sampled individuals, by a method of moments strategy based on quasi-regression. The variance of this estimator is below $C_1/n + C_2d^2/n^3$ for constants $C_j < \infty$, uniformly in $n \geq 2$ and $d \geq 1$. In particular $d \gg n$ is allowed, and heritability can be consistently estimated in some limits where the effects of individual SNPs cannot be consistently estimated due to non-sparsity. Furthermore, when the C_2d^2/n^3 term dominates, then doubling the sample size reduces the variance by almost eight-fold.

The idealization is that the model assumes complete linkage equilibrium, although it does allow for an arbitrary pattern of regression coefficients. In particular the coefficients need not be sparse or Gaussianly distributed, nor independent of the minor allele frequency.

The method and the rate of convergence extend to additive heritability. For a full quadratic model encompassing pairwise interactions with $d^* = 2d + d(d-1)/2$ coefficients the error variance is $O(d^4/n^3)$.

1 Introduction

A missing heritability problem has emerged in genome wide association studies (Manolio et al., 2009). A phenotype known to be largely heritable, for example height, may be found to have only a few significantly predictive SNPs and those few may explain in total only a small portion of the known heritability.

A typical dataset setup has measurements of d SNPs on n subjects where frequently $d \gg n$. Multiple linear regression of the phenotype on SNPs is a useful way to study their effects. Even when $d \gg n$, it is still possible to estimate the needed regression coefficients, if the true coefficient vector is sparse. See for example Zhao and Yu (2006). When that coefficient vector is not sparse, then it is not possible to estimate it accurately (Candès and Davenport, 2011).

2 Notation and models

In this paper, the meaning of $a(n, d) = O(b(n, d))$ is that there is some constant $C < \infty$ for which $|a(n, d)| \leq C \times b(n, d)$ holds for all $n \geq 2$ and all $d \geq 1$. That is, we use non-asymptotic order bounds.

Let y_i be a phenotype for subject i . We assume that the population average of this phenotype has been subtracted from y_i so that $\mathbb{E}(y_i) = 0$. We will suppose that the phenotype is bounded, $|y_i| \leq c < \infty$. Boundedness is not necessary, but we prefer to write c^3 instead of, for example, $\mathbb{E}(y^{12})^{1/4}$. Conversely, we find $\mathbb{E}(y^2)$ more informative than c^2 , so we do not use the bound everywhere.

Let \tilde{x}_{ij} be a genetic measure for SNP j on subject i . For an autosomal SNP, \tilde{x}_{ij} might be 1 if subject i has one or more copies of the minor allele and zero otherwise. Alternatively, $\tilde{x}_{ij} \in \{0, 1, 2\}$ might be the number of copies of minor allele j in subject i . Let $x_{ij} = (\tilde{x}_{ij} - \mathbb{E}(\tilde{x}_{ij}))/\sqrt{\text{Var}(\tilde{x}_{ij})}$ so that $\mathbb{E}(x_{ij}) = 0$ and $\text{Var}(x_{ij}) = 1$ in the relevant population.

We assume additionally that the components x_1, \dots, x_d of \mathbf{x} are independent of each other. This neglects linkage disequilibrium (LD). LD only affects a small fraction of the $O(d^2)$ SNP pairs, but we expect it could be a significant consideration in aggregate because there are so many such pairs.

The k 'th moment of x_j is $\mu_{jk} = \mathbb{E}(x_j^k)$. We will need μ_{j4} for each SNP. If SNP j has minor allele frequency $\epsilon_j > 0$, then $\mu_{4j} = O(\epsilon_j^{-1})$. The variable x_j is bounded. In the motivating examples, $|x_j| \leq c_j$ where $c_j = O(\epsilon_j^{-1/2})$.

In most work, we assume that $\epsilon_j \geq \epsilon_0$ holds for some $\epsilon_0 > 0$. Then μ_{4j} are uniformly bounded. Some intermediate expressions use the average $\bar{\mu}_4 = (1/d) \sum_{j=1}^d \mu_{j4}$. These allow us to consider a limit in which alleles with smaller MAF become eligible for use as n increases. For example, if $\epsilon_j \geq n^{-1/2}$ then $\mu_{4j} \leq n^{1/2}$, and if we only assume $\epsilon_j \geq m/n$ (perhaps the smallest reasonable MAFs to use) for some $m > 0$, then $\mu_{4j} = O(n)$.

The Euclidean norm of \mathbf{x} is denoted $\|\mathbf{x}\|$. We easily find that $\mathbb{E}(\|\mathbf{x}\|^2) = d$ and $\mathbb{E}(\|\mathbf{x}\|^4) = d^2 + d(\bar{\mu}_4 - 1)$, so $\text{Var}(\|\mathbf{x}\|) = d(\bar{\mu}_4 - 1)$. As a result $\|\mathbf{x}\|/d = 1 + O_p(d^{-1/2}\bar{\mu}_4^{1/2})$. In particular

$$\mathbb{E}\left(\sum_{i=1}^n \left(\frac{\|\mathbf{x}_i\|^2}{d} - 1\right)^2\right) = \frac{n}{d}(\bar{\mu}_4 - 1),$$

and so if $d \gg n$, then $\|\mathbf{x}_i\|^2 \approx d$ is a good approximation and one that underlies some of our methods.

We consider several models for heritability. The **linear** model is

$$y_i = \sum_{j=1}^d x_{ij}\beta_j + \varepsilon_i, \quad i = 1, \dots, n \quad (1)$$

where ε_i is a random error uncorrelated with $\mathbf{x}_i = (x_{i1}, \dots, x_{id})$ having mean 0 and variance σ^2 . Effects of the environment, interactions among the SNPs

and interactions between SNPs and the environment are subsumed into ε_i . The pairs (\mathbf{x}_i, y_i) are independent.

Interest naturally centers on the vector $\beta = (\beta_1, \dots, \beta_d)$. Estimation of β is difficult because usually $n \ll d$. If β is sparse, then it can still be estimated using compressed sensing methods. But β might not be sparse and then there are lower bounds (Candès and Davenport, 2011) on its estimation error.

To study linear heritability we are interested in $\sigma_L^2 = \sum_{j=1}^d \beta_j^2 = \|\beta\|^2$. If β is not sparse enough to allow an accurate estimate $\hat{\beta}$ of β , then we cannot rely on $\|\hat{\beta}\|^2$ to be a good estimate of σ_L^2 . Accordingly, we turn our attention to strategies for estimating $\sum_j \beta_j^2$ directly without requiring an accurate estimate of β .

Our starting point is the bias corrected quasi-regression estimator of σ_L^2 from Owen (2000) given below. That algorithm is able to estimate σ_L^2 in an asymptotic regime where $n/d^{2/3} \rightarrow \infty$, without assuming sparsity of β . That is, the sample size n can be sublinear in d . We choose an asymptotic regime with n and d growing together not because more SNPs will be found, but because the usual fixed d and $n \rightarrow \infty$ framework is not a good description for a finite data set with $n \ll d$.

In addition to the linear model, we may also be interested in an additive model which contains x_j^2 terms. If x_j only takes 2 levels, then the centered value $x_j^2 - 1$ is linearly dependent on x_j , and we gain nothing from incorporating it into the linear model. When SNPs are recorded at 3 levels, this additive model allows one to investigate dominance effects. Let

$$z_j = \frac{x_j^2 - \mu_{j3}x_j - 1}{\sqrt{\mu_{j4} - \mu_{j3}^2 - 1}}. \quad (2)$$

Then $\mathbb{E}(z_j) = 0$, $\mathbb{E}(z_j^2) = 1$ and $\mathbb{E}(x_j z_j) = 0$. The **additive** model for heritability is

$$y_i = \sum_{j=1}^d x_{ij} \beta_j + \sum_{j=1}^d z_{ij} \gamma_j + \varepsilon_i, \quad i = 1, \dots, n, \quad (3)$$

where ε_i contains an error uncorrelated with the x_{ij} and the z_{ij} , and does not ordinarily take the same numerical value in models (1), (3) and others that we consider.

We are also interested in the **squares** model

$$y_i = \sum_{j=1}^d z_{ij} \gamma_j + \varepsilon_i, \quad i = 1, \dots, n, \quad (4)$$

the pure **interaction** model

$$y_i = \sum_{j=1}^{d-1} \sum_{k=j+1}^d x_{ij} x_{ik} \beta_{jk} + \varepsilon_i, \quad i = 1, \dots, n, \quad (5)$$

and the **quadratic** model

$$y_i = \sum_{j=1}^d x_{ij} \beta_{jk} + \sum_{j=1}^d z_{ij} \gamma_j + \sum_{j=1}^{d-1} \sum_{k=j+1}^d x_{ij} x_{ik} \beta_{jk} + \varepsilon_i, \quad i = 1, \dots, n. \quad (6)$$

Some of these are of direct biological interest and others play a role in our proofs.

Corresponding to the model parameters above are the sums of squares

$$\sigma_L^2 = \sum_j \beta_j^2, \quad \sigma_S^2 = \sum_j \gamma_j^2, \quad \text{and} \quad \sigma_I^2 = \sum_{j < k} \beta_{jk}^2,$$

along with combinations $\sigma_A^2 = \sigma_L^2 + \sigma_S^2$ and $\sigma_Q^2 = \sigma_A^2 + \sigma_I^2$. These sums of squares take the same value in any model for which all their parts are defined. The following inequality

$$\max(\sigma_L^2, \sigma_S^2, \sigma_I^2) \leq \sigma_Q^2 \leq \mathbb{E}(y^2) \quad (7)$$

will be useful in bounding some error terms. For instance, even though σ_I^2 is the sum of $d(d-1)/2$ nonnegative contributions, that sum is $O(1)$. Also, introducing artificial phenotypes like $y^2 - \mathbb{E}(y^2)$, and applying a version of (7) for them, yields further bounds that we use.

3 Quasi-regression for linear heritability

We begin with estimation of σ_L^2 , for linear heritability. For $d \gg n$ it is infeasible to estimate β_j by least squares. The quasi-regression estimator of β_j is

$$\tilde{\beta}_j = \frac{1}{n} \sum_{i=1}^n x_{ij} y_i.$$

We easily find that $\mathbb{E}(\tilde{\beta}_j) = \beta_j$. The naive estimator of σ_L^2 is

$$\hat{\sigma}_L^2 = \sum_{j=1}^d \tilde{\beta}_j^2.$$

Proposition 1 gives the bias of $\hat{\sigma}_L^2$.

Proposition 1.

$$\mathbb{E}(\hat{\sigma}_L^2) = \frac{n-1}{n} \sigma_L^2 + \frac{1}{n} \mathbb{E}(\|\mathbf{x}\|^2 y^2).$$

Proof. First

$$\begin{aligned} n \mathbb{E}(\tilde{\beta}_j^2) &= \frac{1}{n} \mathbb{E} \left(\sum_i \sum_{i'} x_{ij} x_{i'j} y_i y_{i'} \right) \\ &= \mathbb{E}(x_j^2 y^2) + (n-1) \mathbb{E}(x_j y)^2 \\ &= \mathbb{E}(x_j^2 y^2) + (n-1) \beta_j^2. \end{aligned}$$

Summing over j yields the result. \square

Because $\mathbb{E}(\|\mathbf{x}\|^2) = d$, the bias in $\hat{\sigma}_L^2$ is of order d/n , which is severe when $d \gg n$. But this bias is easily removed. Let

$$\hat{B}_L = \frac{1}{n^2} \sum_{i=1}^n \|\mathbf{x}_i\|^2 y_i^2. \quad (8)$$

We can obtain a bias corrected estimate

$$\hat{\sigma}_{L,BC}^2 = \frac{n}{n-1} (\hat{\sigma}_L^2 - \hat{B}_L), \quad (9)$$

because

$$\mathbb{E}(\hat{\sigma}_{L,BC}^2) = \frac{n}{n-1} \left(\frac{n-1}{n} \sigma_L^2 + \frac{1}{n} \mathbb{E}(\|\mathbf{x}\|^2 y^2) \right) - \frac{1}{n-1} \mathbb{E}(\|\mathbf{x}\|^2 y^2) = \sigma_L^2.$$

A simple approximation to the bias adjustment can be obtained via the concentration of measure approximation $\|\mathbf{x}_i\|^2 \approx d$. Letting,

$$\tilde{B}_L = \frac{d}{n^2} \sum_{i=1}^n y_i^2$$

the resulting estimate is

$$\hat{\sigma}_{L,BC,CM}^2 = \frac{n}{n-1} (\hat{\sigma}_L^2 - \tilde{B}_L). \quad (10)$$

Proposition 2. *Suppose that the phenotype y satisfies the bound $|y| \leq c$. Then*

$$\mathbb{E}((\tilde{B}_L - \hat{B}_L)^2) \leq c^4 \bar{\mu}_4 \frac{d}{n^2}.$$

Proof. First $\tilde{B}_L - \hat{B}_L = n^{-2} \sum_{i=1}^n (d - \|\mathbf{x}_i\|) y_i^2$. Therefore

$$\begin{aligned} \mathbb{E}((\tilde{B}_L - \hat{B}_L)^2) &= \frac{1}{n^3} \mathbb{E}((d - \|\mathbf{x}\|^2)^2 y^4) + \frac{n-1}{n^3} \mathbb{E}((d - \|\mathbf{x}\|^2) y^2)^2 \\ &\leq c^4 \bar{\mu}_4 \frac{d}{n^3} + c^4 \frac{1}{n^2} \mathbb{E}((d - \|\mathbf{x}\|^2)^2) \\ &\leq c^4 \bar{\mu}_4 \frac{d}{n^3} + c^4 \frac{n-1}{n^3} d \bar{\mu}_4. \quad \square \end{aligned}$$

Proposition 2 shows that the concentration of measure approximation changes our estimate of σ_L^2 by an amount with root mean square (RMS) $O(d^{1/2} n^{-1} \bar{\mu}_4^{1/2})$, where c is assumed to be constant in n and d for the phenotype of interest. In an asymptotic regime with $\bar{\mu}_4$ bounded we find that the concentration of measure approximation makes a vanishing RMS effect on our estimate if $d = o(n^2)$. Simulations in Owen (2000) showed very little difference between approximations with and without the concentration of measure approximation.

Theorem 1. Let $\hat{\sigma}_L^2$ be the naive quasi-regression variance estimate and assume that the phenotype satisfies the bound $|y| \leq c$. Then

$$\mathbb{E}((\hat{\sigma}_{L,BC}^2 - \sigma_L^2)^2) = \text{Var}(\hat{\sigma}_{L,BC}^2) = O\left(\frac{1}{n} + \frac{d^2}{n^3}\right).$$

Proof. It suffices to show that $\text{Var}(\hat{\sigma}_L^2 - \hat{B}_L) = O(n^{-1} + d^2n^{-3})$. A lengthy derivation (see Lemma 3 of the appendix) yields

$$\begin{aligned} \text{Var}(\hat{\sigma}_L^2) &\leq \frac{4}{n}c^2\sigma_L^2 + \frac{1}{n^2}\left(4\sigma_Lc^3(d^2 + d\bar{\mu}_4)^{1/2} + 2(d\bar{\mu}_4 + 2)\mathbb{E}(y^4)\right) + \frac{1}{n^3}c^4(d^2 + d\bar{\mu}_4) \\ &= O(n^{-1} + dn^{-2} + d^2n^{-3}) = O(n^{-1} + d^2n^{-3}). \end{aligned}$$

Next

$$\text{Var}(\hat{B}_L) = \frac{1}{n^3}\left(\mathbb{E}(\|\mathbf{x}\|^4y^4) - \mathbb{E}(\|\mathbf{x}\|^2y^2)^2\right) \leq \frac{c^4}{n^3}(d^2 + d\bar{\mu}_4).$$

Thus $\text{Var}(\hat{B}_L) = O(d^2n^{-3})$ and $\text{Var}(\hat{\sigma}_L^2) = O(n^{-1} + d^2n^{-3})$, so $\text{Var}(\hat{\sigma}_L^2 - \hat{B}_L) \leq \text{Var}(\hat{\sigma}_L^2) + \text{Var}(\hat{B}_L) + 2(\text{Var}(\hat{\sigma}_L^2)\text{Var}(\hat{B}_L))^{1/2} = O(n^{-1} + d^2n^{-3})$. \square

When $d \gg n$, the dominant term in the bound for $\text{Var}(\hat{\sigma}_{L,BC}^2)$ is $4c^4d^2n^{-3}$. In that case, doubling n reduces the variance by nearly a factor of eight.

The concentration of measure approximation is accurate to within $O(d^{1/2}n^{-1})$. The standard deviation of $\hat{\sigma}_{L,BC}^2$ is $O(n^{-1/2} + dn^{-3/2})$. As a result, the bias introduced by concentration of measure is of order $\sqrt{n/d} + \sqrt{d/n^3}$ standard errors. It is thus negligible for the values of n and d used in many genome-wide association studies. Both $\hat{\sigma}_{L,BC}^2$ and $\hat{\sigma}_{L,BC,CM}^2$ require $O(nd)$ computation but the latter has a smaller implicit constant.

4 Additive and quadratic heritability

The additive model predicts y by a linear combination of x_j and z_j for $j = 1, \dots, d$. Because x_1, \dots, x_d are independent, it follows that the entire list $x_1, \dots, x_d, z_1, \dots, z_d$ consists of $2d$ mutually uncorrelated predictors. The analysis of the additive model is similar to that of the linear model, except that there are twice as many predictors and the fourth moments of z_j are larger than those of x_j .

The quadratic model incorporates predictors x_jx_k for $j < k$. For distinct pairs $j < k$ and $j' < k'$ with $j \neq j'$ or $k \neq k'$ the predictors x_jx_k and $x_{j'}x_{k'}$ are uncorrelated. Similarly x_jx_k is uncorrelated with all of the $x_{j'}$ and $z_{j'}$ whether or not $j = j'$ or $k = k'$. The big difference in the quadratic model is that there are now more than $d^2/2$ coefficients to estimate instead of $O(d)$ coefficients.

For the additive model, we estimate

$$\tilde{\beta}_j = \frac{1}{n} \sum_{i=1}^n x_{ij}y_i, \quad \text{and} \quad \tilde{\gamma}_j = \frac{1}{n} \sum_{i=1}^n z_{ij}y_i,$$

and the naive estimate of σ_A^2 is

$$\hat{\sigma}_A^2 = \sum_{j=1}^d (\tilde{\beta}_j^2 + \tilde{\gamma}_j^2).$$

Proposition 3.

$$\mathbb{E}(\hat{\sigma}_A^2) = \frac{n-1}{n} \sigma_A^2 + \frac{1}{n} \mathbb{E}((\|\mathbf{x}\|^2 + \|\mathbf{z}\|^2)y^2).$$

Proof. The same argument used in Proposition 3 applies here. \square

Letting

$$\hat{B}_A = \frac{1}{n^2} \sum_{i=1}^n (\|\mathbf{x}_i\|^2 + \|\mathbf{z}_i\|^2) y_i^2 \quad \text{and} \quad \tilde{B}_A = \frac{2d}{n^2} \sum_{i=1}^n y_i^2,$$

we obtain estimates

$$\begin{aligned} \hat{\sigma}_{A,BC}^2 &= \frac{n}{n-1} (\hat{\sigma}_A^2 - \hat{B}_A) \quad \text{and} \\ \hat{\sigma}_{A,BC,CM}^2 &= \frac{n}{n-1} (\hat{\sigma}_A^2 - \tilde{B}_A). \end{aligned}$$

As in the linear case we find that $\mathbb{E}(\hat{\sigma}_{A,BC}^2) = \sigma_A^2$. Also from “A = S + L” we get

$$\mathbb{E}((\tilde{B}_A - \hat{B}_A)^2) \leq c^4 (\bar{\mu}_4(x) + \bar{\mu}_4(z) + 2\sqrt{\bar{\mu}_4(x)\bar{\mu}_4(z)}) \frac{d}{n^2}$$

where $\bar{\mu}_4(x)$ and $\bar{\mu}_4(z)$ are average fourth moments of x_j and z_j respectively.

Theorem 2. Let $\hat{\sigma}_A^2$ be the naive quasi-regression variance estimate of σ_A^2 and assume that the phenotype satisfies the bound $|y| \leq c$. Then

$$\text{Var}(\hat{\sigma}_{A,BC}^2) = O\left(\frac{1}{n} + \frac{d^2}{n^3}\right).$$

Proof. The additive variance explained is $\sigma_A^2 = \sigma_L^2 + \sigma_S^2$, and the estimate partitions accordingly, as $\hat{\sigma}_{A,BC}^2 = \hat{\sigma}_{L,BC}^2 + \hat{\sigma}_{S,BC}^2$. We can apply Theorem 1 to both parts, letting z_j take the role of x_j in the second application. \square

For the quadratic model, we incorporate estimates

$$\tilde{\beta}_{jk} = \frac{1}{n} \sum_{i=1}^n x_{ij} x_{ik} y_i$$

for $1 \leq j < k \leq d$. The naive estimate of σ_Q^2 is

$$\hat{\sigma}_Q^2 = \sum_{j=1}^d (\tilde{\beta}_j^2 + \tilde{\gamma}_j^2) + \sum_{j < k} \tilde{\beta}_{jk}^2.$$

Proposition 4.

$$\mathbb{E}(\hat{\sigma}_Q^2) = \frac{n-1}{n}\sigma_Q^2 + \frac{1}{n}\mathbb{E}((\|\mathbf{x}\|^2 + \|\mathbf{z}\|^2 + \|\mathbf{x}\|^4)y^2).$$

Proof. Recall that $\sigma_Q^2 = \sigma_L^2 + \sigma_S^2 + \sigma_I^2$. To take account of the σ_I^2 component, write for $j < k$,

$$\begin{aligned} n\mathbb{E}(\tilde{\beta}_{jk}^2) &= \frac{1}{n} \sum_i \sum_{i'} \mathbb{E}(x_{ij}x_{ik}x_{i'j}x_{i'k}y_i y_{i'}) \\ &= \mathbb{E}(x_j^2 x_k^2 y^2) + (n-1)\beta_{jk}^2. \end{aligned}$$

Summing over $1 \leq j < k \leq d$, yields

$$\sum_{j < k} \mathbb{E}(\tilde{\beta}_{jk}^2) = \frac{1}{n}\mathbb{E}(\|\mathbf{x}\|^4 y^2) + \frac{n-1}{n}\sigma_I^2. \quad \square$$

Once again, we get a bias correction and a concentration of measure approximation:

$$\begin{aligned} \hat{B}_Q &= \frac{1}{n^2} \sum_{i=1}^n (\|\mathbf{x}_i\|^4 + \|\mathbf{x}_i\|^2 + \|\mathbf{z}_i\|^2) y_i^2 \quad \text{and} \\ \tilde{B}_Q &= \frac{d^2 + (2 + \bar{\mu}_4)d}{n^2} \sum_{i=1}^n y_i^2. \end{aligned}$$

The bias is now of much larger order d^2 , consistent with there being more than $d^2/2$ parameters to estimate.

Using these expressions, we obtain estimates

$$\begin{aligned} \hat{\sigma}_{Q,BC}^2 &= \frac{n}{n-1} (\hat{\sigma}_Q^2 - \hat{B}_Q) \quad \text{and} \\ \hat{\sigma}_{Q,BC,CM}^2 &= \frac{n}{n-1} (\hat{\sigma}_Q^2 - \tilde{B}_Q). \end{aligned}$$

Here $\mathbb{E}(\hat{\sigma}_{Q,BC}^2) = \sigma_Q^2$.

In the model equation “Q = L + S + I” we have already seen how the linear and square parts have an accurate concentration of measure approximation. It remains to consider $\tilde{B}_I - \hat{B}_I$.

Proposition 5. *For the interaction model, the bias adjustments satisfy*

$$\mathbb{E}((\tilde{B}_I - \hat{B}_I)^2) = O\left(\frac{d^3}{n^2}\right).$$

Proof.

$$\mathbb{E}((\tilde{B}_I - \hat{B}_I)^2) = \frac{1}{n^4} \mathbb{E}\left(\left(\sum_{i=1}^n (d^2 + d\bar{\mu}_4 - \|\mathbf{x}_i\|^4) y_i^2\right)^2\right)$$

$$\begin{aligned}
&= \frac{1}{n^3} \mathbb{E} \left((d^2 + d\bar{\mu}_4 - \|\mathbf{x}\|^4)^2 y^4 \right) + \frac{n-1}{n^3} \mathbb{E} \left((d^2 + d\bar{\mu}_4 - \|\mathbf{x}\|^4) y^2 \right)^2 \\
&= \frac{c^4}{n^2} \text{Var}(\|\mathbf{x}\|^4) \\
&= O\left(\frac{d^3}{n^2}\right). \quad \square
\end{aligned}$$

Theorem 3. Let $\hat{\sigma}_{\mathbb{Q}}^2$ be the naive quasi-regression variance estimate of $\sigma_{\mathbb{Q}}^2$ and assume that the phenotype satisfies the bound $|y| \leq c$. Then

$$\text{Var}(\hat{\sigma}_{\mathbb{Q},\text{BC}}^2) = O\left(\frac{1}{n} + \frac{d^4}{n^3}\right).$$

Proof. The error from the additive parts is already seen to be $O(1/n + d^2/n^3)$. Therefore it suffices to show that $\text{Var}(\hat{\sigma}_{\text{I,BC}}^2) = O(1/n + d^4/n^3)$. This is exactly what we would expect from a linear model with $d^* = d(d-1)/2$ independent predictors. But the predictors in the interaction model, while uncorrelated, are not independent.

We can follow the argument for the linear case using d^* predictors of the form $x_j x_k$ for $j < k$ in Lemma 3, replacing each x_k by $x_r x_s$ for $r < s$ and then each x_j by $x_j x_k$ for $j < k$ and each \mathbf{x} by $\mathbf{x}^* = (x_1 x_2, x_1 x_3, \dots, x_{d-1} x_d) \in \mathbb{R}^{d^*}$. We then use $\mathbb{E}(x_j x_k y) = \beta_{jk}$ and replace ε_{L} by $\varepsilon_{\text{I}} = y - \sum_{j < k} x_j x_k \beta_{jk}$.

The analogy holds in a straightforward way except for the term which becomes

$$\frac{n-1}{n^3} \sum_{j < k} \sum_{r < s} \mathbb{E}(x_j x_k x_r x_s y^2)^2$$

which now involves four way sums over components of \mathbf{x} . We need this to be $O(d^*/n^2) = O(d^2/n^2)$. If we use the artificial phenotype $y^2 - \mathbb{E}(y^2)$ and consider a pure quartic model with $d(d-1)(d-2)(d-3)/24$ predictors $x_j x_k x_r x_s$ with j, k, r, s all distinct, then the terms with no ties among j, k, r, s sum to at most $\mathbb{E}(y^2)^2$. What remains are the terms where $\{j, k\} \cap \{r, s\} \neq \emptyset$. There are $O(d^2)$ such terms and they are uniformly bounded. \square

Thus while it is possible to estimate the interaction inheritance and hence the quadratic inheritance with $n \ll d^2$, we still need n of larger order than d . Practically, this implies that we may be able to investigate interactions among strategically chosen subsets of SNPs but perhaps not the entire interaction at practical sample sizes.

The concentration of measure approximation for the interaction inheritance is $\|\mathbf{x}\|^4 \approx d^2 + d\bar{\mu}_4$. The root mean square change from this shortcut is $d^{3/2}/n$, while the standard deviation of the estimate is $O(n^{-1/2} + d^2 n^{-3/2})$. Suppose that n is just barely large enough to allow estimation of σ_{I}^2 , perhaps $n \propto d^{4/3+\epsilon}$. Then the RMS difference between $\hat{\sigma}_{\text{I,BC}}^2$ and $\hat{\sigma}_{\text{I,BC,CM}}^2$ is of order $d^{3/2} n^{-1} = d^{1/6-\epsilon}$ which diverges. As a result, the concentration of measure shortcut is not recommended for the interaction or quadratic model.

Acknowledgments

I thank Or Zuk for helpful discussions. This work was supported by the NSF under grant DMS-0906056.

References

- Candès, E. and Davenport, M. A. (2011). How well can we estimate a sparse vector? Technical report, Stanford University.
- Manolio, T. A., Collins, F. S., Cox, N. J., Goldstein, D. B., Hindorf, L. A., Hunter, D. J., McCarthy, M. I., Ramos, E. M., Cardon, L. R., Chakravarti, A., Cho, J. H., Guttmacher, A. E., Kong, A., Kruglyak, L., Mardis, E., Rotimi, C. N., Slatkin, M., Valle, D., Whittemore, A. S. Boehnke, M., Clark, A. G., Eichler, E. E., Gibson, G., Haines, J. L., Mackay, T. F. C., McCarroll, S. A., and Visscher, P. M. (2009). Finding the missing heritability of complex diseases. *Nature*, 461:747–753.
- Owen, A. B. (2000). Assessing linearity in high dimensions. *The Annals of Statistics*, 28(1):1–19.
- Zhao, P. and Yu, B. (2006). On model selection consistency of the lasso. *Journal of Machine Learning Research*, 7:2541–2563.

Appendix

Here we prove some basic Lemmas used in the main Theorem. The following basic properties of $\|\mathbf{x}\|^2$ are used below:

$$\begin{aligned} \mathbb{E}(\|\mathbf{x}\|^2) &= d, & \mathbb{E}(\|\mathbf{x}\|^4) &= d^2 + d\bar{\mu}_4 \\ \mathbb{E}((\|\mathbf{x}\|^2 - d)^2) &= d\bar{\mu}_4, & \text{and } \mathbb{E}((\|\mathbf{x}\|^2 - d)^r) &= O(d^{r-1}), \quad \text{integer } r \geq 3. \end{aligned}$$

Lemma 1 uses the artificial phenotype $y^2 - \mathbb{E}(y^2)$ to bound a quantity which appears in $\text{Var}(\hat{\sigma}_L^2)$.

Lemma 1.

$$\sum_j \sum_k \mathbb{E}(x_j x_k y^2)^2 \leq (d\bar{\mu}_4 + 2)\mathbb{E}(y^4).$$

Proof. We write

$$\sum_j \sum_k \mathbb{E}(x_j x_k y^2)^2 = \sum_j \sum_{k \neq j} \mathbb{E}(x_j x_k y^2)^2 + \sum_j \mathbb{E}(x_j^2 y^2)^2.$$

For the first term, consider an artificial phenotype $y^2 - \mathbb{E}(y^2)$. Then

$$\sum_j \sum_{k \neq j} \mathbb{E}(x_j x_k y^2)^2 = 2 \sum_{j < k} \mathbb{E}(x_j x_k (y^2 - \mathbb{E}(y^2)))^2$$

$$\begin{aligned} &\leq 2\mathbb{E}((y^2 - \mathbb{E}(y^2))^2) \\ &\leq 2\mathbb{E}(y^4), \end{aligned}$$

because of inequality (7) applied to the interaction model for this artificial phenotype. Next

$$\sum_j \mathbb{E}(x_j^2 y^2)^2 \leq \sum_j \mathbb{E}(x_j^4) \mathbb{E}(y^4) \leq \mathbb{E}(y^4) d \bar{\mu}_4. \quad \square$$

Lemma 2.

$$\mathbb{E}((y - \varepsilon) \|\mathbf{x}\|^2 y^3) \leq \sigma_L c^3 (d^2 + d \bar{\mu}_4)^{1/2}.$$

Proof. Using Cauchy-Schwarz, boundedness of y , and the above results on $\|\mathbf{x}\|$:

$$\begin{aligned} \mathbb{E}((y - \varepsilon) \|\mathbf{x}\|^2 y^3) &\leq \mathbb{E}((y - \varepsilon)^2)^{1/2} \mathbb{E}(\|\mathbf{x}\|^4 y^6)^{1/2} \\ &\leq \sigma_L c^3 (d^2 + d \bar{\mu}_4)^{1/2}. \quad \square \end{aligned}$$

Lemma 3 is the main lemma.

Lemma 3. *Under the conditions of Theorem 1*

$$\text{Var}(\hat{\sigma}_L^2) \leq \frac{4}{n} c^2 \sigma_L^2 + \frac{1}{n^2} \left(4\sigma_L c^3 (d^2 + d \bar{\mu}_4)^{1/2} + 2(d \bar{\mu}_4 + 2) \mathbb{E}(y^4) \right) + \frac{1}{n^3} c^4 (d^2 + d \bar{\mu}_4).$$

Proof. First,

$$\text{Var}(\hat{\sigma}_L^2) = \sum_j \sum_k \mathbb{E}(\tilde{\beta}_j^2 \tilde{\beta}_k^2) - \left(\mathbb{E} \left(\sum_j \tilde{\beta}_j^2 \right) \right)^2$$

We simplify some expressions below using $\mathbb{E}(x_j y) = \beta_j$ and $\varepsilon_L \equiv y - \sum_j x_j \beta_j$.

For each pair $j, k \in \{1, 2, \dots, d\}$,

$$\begin{aligned} \mathbb{E}(\tilde{\beta}_j^2 \tilde{\beta}_k^2) &= \mathbb{E} \left(\left(\frac{1}{n} \sum_i x_{ij} y_i \right)^2 \left(\frac{1}{n} \sum_i x_{ik} y_i \right)^2 \right) \\ &= \frac{1}{n^4} \sum_i \sum_{i'} \sum_{i''} \sum_{i'''} \mathbb{E}(x_{ij} x_{i'j} x_{i''k} x_{i'''k} y_i y_{i'} y_{i''} y_{i'''}) \\ &= \frac{(n-1)(n-2)(n-3)}{n^3} \mathbb{E}(x_j y)^2 \mathbb{E}(x_k y)^2 \\ &\quad + \frac{(n-1)(n-2)}{n^3} \left(\mathbb{E}(x_j y)^2 \mathbb{E}(x_k^2 y^2) + \mathbb{E}(x_j^2 y^2) \mathbb{E}(x_k y)^2 + 4\mathbb{E}(x_j x_k y^2) \mathbb{E}(x_j y) \mathbb{E}(x_k y) \right) \\ &\quad + \frac{n-1}{n^3} \left(2\mathbb{E}(x_j y) \mathbb{E}(x_j x_k^2 y^3) + 2\mathbb{E}(x_k y) \mathbb{E}(x_k x_j^2 y^3) \right. \\ &\quad \quad \left. + \mathbb{E}(x_j^2 y^2) \mathbb{E}(x_k^2 y^2) + 2\mathbb{E}(x_j x_k y^2)^2 \right) \\ &\quad + \frac{1}{n^3} \mathbb{E}(x_j^2 x_k^2 y^4). \end{aligned}$$

Summing over j and k ,

$$\begin{aligned}
\sum_j \sum_k \mathbb{E}(\tilde{\beta}_j^2 \tilde{\beta}_k^2) &= \frac{(n-1)(n-2)(n-3)}{n^3} \sigma_L^4 \\
&+ \frac{(n-1)(n-2)}{n^3} \left(2\sigma_L^2 \mathbb{E}(\|\mathbf{x}\|^2 y^2) + 4\mathbb{E}((y - \varepsilon_L)^2 y^2) \right) \\
&+ \frac{n-1}{n^3} \left(4\mathbb{E}((y - \varepsilon_L) \|\mathbf{x}\|^2 y^3) + \mathbb{E}(\|\mathbf{x}\|^2 y^2)^2 + 2 \sum_j \sum_k \mathbb{E}(x_j x_k y^2)^2 \right) \\
&+ \frac{1}{n^3} \mathbb{E}(\|\mathbf{x}\|^4 y^4).
\end{aligned}$$

Next $\mathbb{E}(\tilde{\beta}_j^2) = \mathbb{E}(x_j^2 y^2)/n + (n-1)\beta_j^2/n$. Therefore

$$\begin{aligned}
\left(\mathbb{E} \left(\sum_j \tilde{\beta}_j^2 \right) \right)^2 &= \frac{1}{n^2} \sum_j \sum_k \mathbb{E}(x_j^2 y^2) \mathbb{E}(x_k^2 y^2) \\
&+ \frac{(n-1)^2}{n^2} \sum_j \sum_k \beta_j^2 \beta_k^2 + 2 \frac{n-1}{n^2} \sum_j \sum_k \mathbb{E}(x_j^2 y^2) \beta_k^2 \\
&= \frac{1}{n^2} \mathbb{E}(\|\mathbf{x}\|^2 y^2)^2 + \frac{(n-1)^2}{n^2} \sigma_L^4 + 2 \frac{n-1}{n^2} \mathbb{E}(\|\mathbf{x}\|^2 y^2) \sigma_L^2.
\end{aligned}$$

Subtracting,

$$\begin{aligned}
\text{Var}(\hat{\sigma}_L^2) &= \left(\frac{(n-1)(n-2)(n-3)}{n^3} - \frac{(n-1)^2}{n^2} \right) \sigma_L^4 \\
&+ 2 \left(\frac{(n-1)(n-2)}{n^3} - \frac{n-1}{n^2} \right) \mathbb{E}(\|\mathbf{x}\|^2 y^2) \sigma_L^2 \\
&+ \left(\frac{n-1}{n^3} - \frac{1}{n^2} \right) \mathbb{E}(\|\mathbf{x}\|^2 y^2)^2 \\
&+ 4 \frac{(n-1)(n-2)}{n^3} \mathbb{E}((y - \varepsilon_L)^2 y^2) \\
&+ \frac{n-1}{n^3} \left(4\mathbb{E}((y - \varepsilon_L) \|\mathbf{x}\|^2 y^3) + 2 \sum_j \sum_k \mathbb{E}(x_j x_k y^2)^2 \right) \\
&+ \frac{1}{n^3} \mathbb{E}(\|\mathbf{x}\|^4 y^4).
\end{aligned}$$

The first three terms above are less than or equal to zero. Therefore

$$\begin{aligned}
\text{Var}(\hat{\sigma}_L^2) &\leq 4 \frac{(n-1)(n-2)}{n^3} \mathbb{E}((y - \varepsilon_L)^2 y^2) \\
&+ \frac{n-1}{n^3} \left(4\mathbb{E}((y - \varepsilon_L) \|\mathbf{x}\|^2 y^3) + 2 \sum_j \sum_k \mathbb{E}(x_j x_k y^2)^2 \right) \\
&+ \frac{1}{n^3} \mathbb{E}(\|\mathbf{x}\|^4 y^4).
\end{aligned}$$

We can simplify this bound using $\mathbb{E}((y - \varepsilon_L)^2 y^2) \leq c^2 \sigma_L^2$, Lemmas 1 and 2, and $\mathbb{E}(\|\mathbf{x}\|^4 y^4) \leq c^4 (d^2 + d\bar{\mu}_4)$, yielding

$$\begin{aligned} \text{Var}(\hat{\sigma}_L^2) &\leq 4 \frac{(n-1)(n-2)}{n^3} c^2 \sigma_L^2 \\ &\quad + \frac{n-1}{n^3} \left(4\sigma_L c^3 (d^2 + d\bar{\mu}_4)^{1/2} + 2(d\bar{\mu}_4 + 2)\mathbb{E}(y^4) \right) + \frac{1}{n^3} c^4 (d^2 + d\bar{\mu}_4) \\ &\leq \frac{4}{n} c^2 \sigma_L^2 + \frac{1}{n^2} \left(4\sigma_L c^3 (d^2 + d\bar{\mu}_4)^{1/2} + 2(d\bar{\mu}_4 + 2)\mathbb{E}(y^4) \right) + \frac{1}{n^3} c^4 (d^2 + d\bar{\mu}_4). \end{aligned}$$

□