

Evaluation and application of summary statistic imputation to discover new height-associated loci

Sina Rüeger, Aaron McDaid and Zoltán Kutalik

March 13, 2018

S3 Appendix: Accounting for varying sample size and missingness

All previously published methods assume that all effect estimates are based on the same set of N individuals. This assumption does not always hold, for example when meta-analysing studies use different genotyping chips or different imputation panels. As a result, the covariance between effect estimates will change. In the extreme case when effect estimates are computed in two non-overlapping samples, the correlation will be zero even if there is very high LD between the two SNVs.

\mathbf{C} , as defined above, is an estimate of $\Gamma_{\mathcal{M}\mathcal{M}}$, which is the correlation matrix due to LD among the tag SNVs in the current region. We define \mathbf{N} as a vector recording the sample size of each tag SNV, N_{max} as the maximum in \mathbf{N} , and assume that every tag SNV k the sample of individuals is a subset of a complete sample of N_{max} individuals. Unlike the situation where the sample size is the same for all tag SNVs, varying sample size requires to first impute the Z-statistic z before computing the standardised effect size a . For each tag SNV k , we have observed a ‘partial’ Z-statistic, z_k° , computed over the N_k individuals. For a SNV u , our goal is to impute a ‘complete’ Z-statistic, z_u , estimated from a complete sample of N_{max} people. In other words, for any SNV u we wish to impute $z_u|z_{\mathcal{M}}^\circ$.

To perform imputation, we require the correlation between any target complete Z-statistic, z_u , and any observed partial Z-statistic, z_k° , (with $k \in \mathcal{M}$),

$$\mathbf{d}_k := \text{Cor}[z_u, z_k^\circ] = c_{uk} \sqrt{\frac{N_k}{N_{max}}}$$

where $c_{uk} = \widehat{\Gamma}_{uk}$ is our estimate of the LD-correlation between the two SNVs. We also require the correlations among the partial Z-statistics, $z_{\mathcal{M}}^\circ$. For any two SNVs $k, l \in \mathcal{M}$, the correlation of their observed partial Z-statistics can be calculated as:

$$\text{Cor}[z_k^\circ, z_l^\circ] = c_{kl} \frac{N_{k \cap l}}{\sqrt{N_k N_l}}$$

29 where N_k and N_l are the number of individuals for which SNV k and l are available,
 30 respectively, and $N_{k \cap l}$ is the number of individuals that contributed to the calculation
 31 of the effect estimates for both SNVs k and l . The estimation of $N_{k \cap l}$ is discussed in the
 32 next section.

33 We can use this to adjust the correlation matrices \mathbf{C} and \mathbf{c} , respectively \mathbf{D} and \mathbf{d} , the
 34 elements of which are

$$\mathbf{D}_{k,l} = c_{kl} \delta_{kl},$$

35 By defining $\delta_{kl} := \frac{N_{k \cap l}}{\sqrt{N_k N_l}}$, we can calculate the adjusted (estimated) correlation matrix
 36 \mathbf{D} , where each element is calculated as follows:

$$\mathbf{D}_{kl} = c_{kl} \delta_{kl}.$$

37 \mathbf{D} and \mathbf{d} are therefore defined as adjusted versions of \mathbf{C} and \mathbf{c} respectively. \mathbf{C} and \mathbf{D}
 38 are $q \times q$ matrices, where q is the number of tag SNVs, and they are recomputed in each
 39 region. \mathbf{c} and \mathbf{d} are vectors of length q , and are recomputed for each target SNV u .

40 The conditional distribution is

$$(z_u - \mathbb{E}[z_u]) | \mathbf{z}_{\mathcal{M}}^\circ \sim \mathcal{N}(\mathbf{d}' \mathbf{D}^{-1} (\mathbf{z}_{\mathcal{M}}^\circ - \mathbb{E}[\mathbf{z}_{\mathcal{M}}^\circ]), 1 - \mathbf{d}' \mathbf{D}^{-1} \mathbf{d}).$$

41 Applying the simplifying assumption that $\mathbb{E}[z_u] \approx \mathbf{d}' \mathbf{D}^{-1} \mathbb{E}[\mathbf{z}_{\mathcal{M}}^\circ]$, similar to the assump-
 42 tion that took us from 1 to 2,

$$z_u | \mathbf{z}_{\mathcal{M}}^\circ \sim \mathcal{N}(\mathbf{d}' \mathbf{D}^{-1} \mathbf{z}_{\mathcal{M}}^\circ, 1 - \mathbf{d}' \mathbf{D}^{-1} \mathbf{d})$$

43 and therefore we impute $z_u | \mathbf{z}_{\mathcal{M}}^\circ$ as

$$\hat{z}_u = \mathbb{E}[z_u | \mathbf{z}_{\mathcal{M}}^\circ] = \mathbf{d}' \mathbf{D}^{-1} \mathbf{z}_{\mathcal{M}}^\circ. \quad (\text{S1})$$

44 In order to convert \hat{z}_u into the corresponding estimate of the standardised effect, we
 45 consider the (hypothetical) process of imputing each individual genotype.

46 If we had the individual-level genetic data, with j to index individuals, each element
 47 of the N_{max} -element vector \mathbf{g}^u for SNV u could be imputed using genotypes from the
 48 tag SNVs $G^{\mathcal{M}}$ via $\hat{g}_j^u = \mathbf{c}'_{\mathcal{M}(j),u} \mathbf{C}_{\mathcal{M}(j),\mathcal{M}(j)}^{-1} G_j^{\mathcal{M}}$, where the set $\mathcal{M}(j)$ can be different for
 49 each individual as each individual has a different set of tagged SNVs. The corresponding
 50 standardised effect estimate, based on linear regression, would be

$$\hat{a}_u = \mathbb{E}[a_u | \mathbf{a}_{\mathcal{M}}] = \frac{(\hat{\mathbf{g}}^u)' \mathbf{y}}{(\hat{\mathbf{g}}^u)' \hat{\mathbf{g}}^u}$$

51 The denominator of this is $(\hat{\mathbf{g}}^u)' \hat{\mathbf{g}}^u = N_{max} \mathbf{d}' \mathbf{D}^{-1} \mathbf{d}$, as opposed to $(\mathbf{g}^u)' \mathbf{g}^u = N_{max}$, and

52 we define the *effective sample size* as $N_{max}\mathbf{d}'\mathbf{D}^{-1}\mathbf{d}$. Therefore, even though we do not
 53 have the per-individual genetic data, we can impute the standardised effect a_u via

$$\hat{a}_u = \mathbb{E}[a_u | \mathbf{z}_{\mathcal{M}}^{\circ}] = \frac{\hat{z}_u}{\sqrt{N_{max}\mathbf{d}'\mathbf{D}^{-1}\mathbf{d}}}. \quad (\text{S2})$$

54 **Estimating overlap $N_{k \cap l}$ and δ**

55 Typically, we do not know the details of the exact sample overlap for every pair of SNVs,
 56 $n_{k \cap l}$, and instead simply know N_{max} and the vector \mathbf{N} . Therefore, we must derive the
 57 sample overlap based on assumptions about the dependence structure of missingness.

If each SNV has a corresponding binary missingness vector, the correlation between these missingness vectors will be maximised when the sample overlap is at its maximum, $N_{k \cap l} = \min(N_k, N_l)$. To enable the *dependent* approach, we construct a \mathbf{D} matrix by replacing $N_{k \cap l}$ with $\min(N_k, N_l)$,

$$\mathbf{D}_{kl}^{(dep)} = \mathbf{C}_{kl}\hat{\delta}_{kl}^{(dep)} = \mathbf{C}_{kl} \min\left(\frac{\sqrt{N_k}}{\sqrt{N_l}}, \frac{\sqrt{N_l}}{\sqrt{N_k}}\right). \quad (\text{S3})$$

58 and plug $\mathbf{D}^{(dep)}$ into Eqs. (S1) and (S2).

If the missingness vectors are *independent* of each other, the expected overlap can be estimated as

$$\mathbf{D}_{kl}^{(ind)} = \mathbf{C}_{kl}\hat{\delta}_{kl}^{(ind)} = \mathbf{C}_{kl} \frac{\sqrt{N_k N_l}}{N_{max}}. \quad (\text{S4})$$