

Consider P traits (Y_1, \dots, Y_P) , where number of observations of each trait $Y_i = N$. When stacked, Y is a $NP \times 1$ vector. We won't consider covariates other than sites since we assume they're orthogonal and can just project them away. We consider the site covariates $X_{n \times \#sites}$ and site effects for each phenotype γ_j which are each $P \times 1$ vectors. Lastly the error ($\epsilon \sim N(0, \sigma^2)$) is scaled for each site-phenotype combination, where each δ_j is an $\#nsites$ vector that scales the error variance

1 The Model

$$\begin{aligned} Y &= (Y_1, \dots, Y_P) \\ Y &= X_g \beta + X_s \gamma + \epsilon \\ \text{vec}(Y) &= (\beta'_g \otimes I_n) \text{vec}(X_G) + (\gamma'_s \otimes I_n) \text{vec}(X_s) + \text{vec}(\epsilon) \\ \text{var}(\text{vec}(\epsilon)) &= \Sigma_e \otimes I_n \end{aligned}$$

Both methods will be GRM based so

$$\text{var}(\text{vec}(X_g \beta)) = \Sigma_g \otimes A = \begin{bmatrix} \sigma_{g1}^2 A & \dots & \rho_{g,1P} A \\ \vdots & \ddots & \vdots \\ \rho_{g,P1} A & \dots & \sigma_{gP}^2 A \end{bmatrix}$$

$$\text{Var}(Y) = \begin{bmatrix} \text{Var}(Y_1) & \dots & \Sigma_{1P} \\ \vdots & \ddots & \vdots \\ \Sigma_{P1} & \dots & \text{Var}(Y_P) \end{bmatrix}$$

Considering just two phenotypes we have

$$\begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} \sim N\left(\begin{bmatrix} X_s \gamma_1 \\ X_s \gamma_2 \end{bmatrix}, \begin{bmatrix} \sigma_{g1}^2 A + \sigma_{e1}^2 I & \rho_g A + \rho_e I \\ \rho_g A + \rho_e I & \sigma_{g2}^2 A + \sigma_{e2}^2 I \end{bmatrix}\right)$$

2 AdjHE RE method

Treating site effects as random, and then no covariance between site, genetic, or error we get

$$\begin{bmatrix} Y_1 \\ Y_2 \end{bmatrix} \sim N(0, \begin{bmatrix} \sigma_{g1}^2 A + \sigma_{s1}^2 S + \sigma_{e1}^2 I & \rho_g A + \rho_s S + \rho_e I \\ \rho_g A + \rho_s S + \rho_e I & \sigma_{g2}^2 A + \sigma_{s2}^2 S + \sigma_{e2}^2 I \end{bmatrix})$$

This means

$$\text{Cov}(Y_1, Y_2) = \rho_g A + \rho_s S + \rho_e I$$

So the covariance between different traits has parts based in the similarities in the genetics, site, and experimental error.

3 ComBat

Long story short, CovBat technique only addresses the covariance due to the sites.

First they use Combat to residualize by subtracting the empirical Bayes estimators of the site mean (γ_S^*) and variance (δ_S^*). And we know that they empirical bayes estimators are consistent ($\gamma_S^* \xrightarrow{P} \gamma_S, \delta_S^* \xrightarrow{P} \delta_S$). Looking at the combat adjustment

$$Y^{combat} = (Y - X_s \gamma^*) \delta^{*-1}$$

We have

$$X_s \beta_s^* \xrightarrow{P} X_s \beta_s$$

Therefore a projection defined by the site means we have

$$Q_s^* \xrightarrow{P} Q_s = X_s (X_s' X_s)^{-1} X_s'$$

However, since $X \not\perp PC$ we run into the problem that

$$Q_s Y = Q_s X_g \beta + \epsilon$$

Assuming that we have a perfectly balanced study s.t. $PC \perp X$ Then we'd have

$$Q_s Y = X_g \beta + \epsilon \sim N(0, \sigma_g^2 A + \sigma_e^2)$$

$$\delta^* \xrightarrow{P} \delta \quad \therefore \quad \delta^{*-1} \stackrel{CMT}{=} \text{diag}(1/\delta_j^*) \xrightarrow{P} \delta^{-1}$$

By Slutsky's theorem and CMT and assuming no differences ($\delta = I$)

$$Y^{combat} \xrightarrow{d} \epsilon \sim N(0, \sigma_g^2 A + \sigma_e^2 I)$$

4 Covbat

Taking the residuals from Combat Y^{combat} (which have mean 0) we denote the covariance matrices of each residualized phenotype as $\text{var}(Y_j^{combat}) = \Sigma_j$. They

take a PCA decomposition of the residualized phenotypes to get the covariance matrix

$$\Sigma = \sum_{k=1}^p \lambda_k \phi_k \phi_k^T$$

the residualized phenotypes are expressed using the coordinates (η) along each of the first p eigenvectors. However, because part of the covariance contains the GRM, this would affect the heritability estimate.