

Outline

1. Write out and talk about model for the data (maybe discuss M vs. Beta value here)
2. Discuss methods without cell type data (i.e. Houseman and CATE).
3. When do these perform well?
 - (a) Standardized effect size for L , taking into account variability Ξ .
 - (b) Discuss $\frac{1}{p}\Gamma^T\Sigma^{-1}\Gamma$ result in Bai and Li, i.e. how informative methylation data are for cell type.
 - (c) Figures: simulation results. We will ALWAYS underestimate Ω when data are not informative. Effect on results is small when Ω is small (can we prove this?)
 - (d) Give Amish/Hutterite data example. Talk about size of $\frac{1}{p}\Gamma^T\Sigma^{-1}\Gamma$, MAJOR differences in the estimates for π_0 , the fact that the data suggest that there is confounding (p-value for α). Large difference in $\hat{\pi}_0$ along with the fact that $\frac{1}{p}\Gamma^T\Sigma^{-1}\Gamma = O\left(\frac{1}{n}\right)$ seem to indicate the data are not informative enough to estimate the correlation between C and X .
4. If data are not informative for cell type, we need another method to estimate the correlation between C and X . Obvious alternative: collect cell type or use training data to estimate the correlation between C and X .
5. In some circumstances, we can do well with only partially observed data (in the absence of other confounders). Plot simulation figures. When does this happen? What happens when we have additional confounders?
6. Recommendations: if you have strong prior assumptions that there are additional confounders that you can measure, it is always a good idea to measure them.
 - (a) Talk about how large correlation needs to be to start making serious errors. If standardized correlation is small, maybe you don't need to worry about it and can estimate it from the data.
 - (b) With strong prior knowledge, it is always a good idea to include covariates that may have an effect on response and are correlated with the variable of interest.

1 Notation

- All matrices (including vectors) are bold and all vectors are column vectors, unless otherwise stated.
- Throughout the paper, the matrix $\mathbf{X} = \mathbf{X}_{d \times n} \in \mathbb{R}^{d \times n}$ contains d covariates of interest (including an intercept) measured on n individuals. For example, in the real data example $d = 2$ where one covariate is the intercept and the other is the Amish vs. Hutterite contrast.
- Cell type proportions for K cells measured on n individuals is given as the matrix $\mathbf{C} = \mathbf{C}_{K \times n} \in [0 - 1]^{K \times n}$.
- $P_{\mathbf{A}}$ is the orthogonal projection onto the linear subspace generated by the columns of \mathbf{A} and $P_{\mathbf{A}}^\perp$ projects onto the subspace orthogonal to $\text{Image}(\mathbf{A})$.
- $\text{vec}(\mathbf{A})$ is the vectorized version of a matrix, obtain by column stacking.
- When the matrix $\mathbf{U} \in \mathbb{R}^{p \times n}$ is random and the covariance of the elements of \mathbf{U} are separable across rows and columns, we will write

$$\mathbf{U} \sim (\boldsymbol{\mu}_{p \times n}, [\mathbf{S}_{p \times p}, \mathbf{R}_{n \times n}])$$

to indicate that $\boldsymbol{\mu}$ is the mean of \mathbf{U} , \mathbf{S} is the covariance across rows and \mathbf{R} is the covariance across columns. That is, $\text{Var}(\text{vec}(\mathbf{U})) = \mathbf{R} \otimes \mathbf{S}$. For example, if we measured the methylation on p CpG sites across n independent individuals, $\mathbf{R} = \mathbf{I}_n$ and \mathbf{S} would be the covariance across the p CpG sites in a single sample. We characterize random variables by their first two moments, as it is often times more important to accurately estimate the mean and variance than it is to use the correct probability model (cite contiguity paper).

2 The Model

Let $Y_{p \times n}$ be the methylation response measured on p CpG sites across n independent individuals, $X_{d \times n} = [\vec{x}_1 \ \cdots \ \vec{x}_n]$ the matrix of d covariates of interest (i.e. disease status, sex, etc.) measured on the same n individuals and $C_{K \times n} = [\vec{c}_1 \ \cdots \ \vec{c}_n]$ be the cell composition matrix for K cell types across the n individuals. Assuming that methylation status is ONLY affect by the covariate(s) of interest X and cell type composition C , full data model for Y is then:

$$Y_{p \times n} = B_{p \times d} X_{d \times n} + L_{p \times K} C_{K \times n} + E_{p \times n} \quad (1)$$

Where $E_{p \times n} \sim (0_{p \times n}, [\Sigma_{p \times p} = \text{diag}(\sigma_1^2, \dots, \sigma_p^2), I_n])$, $L = \begin{bmatrix} \ell_1^T \\ \vdots \\ \ell_p^T \end{bmatrix}$ and $B = \begin{bmatrix} \beta_1^T \\ \vdots \\ \beta_p^T \end{bmatrix}$. That is,

$$Y | (X, C) \sim (BX + LC, [\Sigma, I_n]) \quad (2)$$

We are interested in estimating the effect β_i for each CpG site $i = 1, \dots, p$.

However, sometimes we do not have access to the cell type composition data. If cell type is uncorrelated with the covariate(s) of interest, we can account for it using factor analysis correction techniques (e.g. SVA). When cell type and the covariate(s) of interest are correlated, we assume that

$$C_{K \times n} = \Omega_{K \times d} X_{d \times n} + \Xi_{K \times n}, \quad \Xi = [\vec{\xi}_1 \ \cdots \ \vec{\xi}_n] \quad (3)$$

$$\Xi \sim (0, [\Lambda_{K \times K}, I_n]) \quad (4)$$

Note that since cell type composition is proportion data, the variance $E(\vec{\xi}_i \vec{\xi}_i^T | \vec{x}_i)$ is a function of the mean $E(\vec{c}_i | \vec{x}_i)$, which is not captured by our assumptions on Ξ . We show by simulation that we can ignore the mean-variance relationship and still perform accurate inference on \hat{B} .

When C is unobserved, we can write the full data model as

$$Y = BX + L\Omega X + L\Xi + E = BX + \tilde{L}\tilde{\Omega}X + \tilde{L}\tilde{\Xi} + E \quad (5)$$

where $\tilde{L} = \begin{bmatrix} \tilde{\ell}_1^T \\ \vdots \\ \tilde{\ell}_p^T \end{bmatrix} = L\Lambda^{1/2}$, $\tilde{\Omega} = \Lambda^{-1/2}\Omega$ and $\tilde{\Xi} \sim (0, [I_K, I_n])$. The first and second moments conditional on only observing the covariates X are

$$Y | X \sim (BX + \tilde{L}\tilde{\Omega}X, [\tilde{L}\tilde{L}^T + \Sigma, I_n]). \quad (6)$$

3 Estimation Methods With Unobserved Cell Type

We start by assuming X and cell type C are both observed. We see from equation 3 that if we were to estimate $\tilde{\Omega}$ using OLS (assuming Λ was known), we would get

$$\tilde{\Omega}^{(OLS)} = \Lambda^{-1/2} C X^T (X X^T)^{-1} = \tilde{\Omega} + \tilde{\Xi} X^T (X X^T)^{-1}. \quad (7)$$

Note that there are two sources of empirical correlation. The first is due to the true correlation Ω and the second is due to the residual $\tilde{\Xi} X^T (X X^T)^{-1}$, which we will refer to as spurious correlation. When C is observed (i.e. $\Omega X + \Xi$ is observed), the unbiased OLS estimate for B is

$$\hat{B}^{(OLS)} = Y X^T [X X^T]^{-1} - Y P_{X^\perp}^T \tilde{\Xi}^T [\tilde{\Xi} P_{X^\perp}^T \tilde{\Xi}^T]^{-1} \tilde{\Omega}^{(OLS)} = Y X^T [X X^T]^{-1} - \hat{\tilde{L}}^{(OLS)} \tilde{\Omega}^{(OLS)} \quad (8)$$

with variance

$$\text{Var}(\hat{\beta}_i^{(OLS)}) = \sigma_i^2 \left([X X^T]^{-1} + (\tilde{\Omega}^{(OLS)})^T [\tilde{\Xi} P_{X^\perp}^T \tilde{\Xi}^T]^{-1} \tilde{\Omega}^{(OLS)} \right) \quad (9)$$

which by Gauss-Markov is the smallest variance possible among all unbiased estimators for β_i . If cell type C is NOT observed, to correct for cell type heterogeneity it stands to reason that our goal should be to estimate the part of C

uncorrelated with \mathbf{X} and the part that is correlated with \mathbf{X} . That is, we wish to estimate $P_{\mathbf{X}^T}^\perp \tilde{\Xi}^T$ and $\tilde{\Omega}^{(\text{OLS})} = \text{true correlation} + \text{spurious correlation from the methylation data } \mathbf{Y}$.

Luckily, considerable attention has been given to developing methods to estimate \mathbf{B} in the presence of confounding that is correlated with the design matrix \mathbf{X} (LEAPP, Houseman et al., CATE). Using **Equ. 5** as a reference, all three methods assume that $\tilde{\Xi}$ is independent of \mathbf{X} and that \mathbf{B} is sparse. That is, the covariate of interest affects only a small

fraction of the CpG sites. Both of these assumptions tend to hold in real data. To get estimates $\hat{\mathbf{B}} = \begin{bmatrix} \hat{\beta}_1^T \\ \vdots \\ \hat{\beta}_p^T \end{bmatrix}$ and $\hat{\text{Var}}(\hat{\beta}_i)$

for CpG sites $i = 1, \dots, p$, all three methods use the following general procedure:

1. Compute $\hat{\tilde{\mathbf{L}}}$ and $\hat{\tilde{\Sigma}}$, estimates for $\tilde{\mathbf{L}}$ and $\tilde{\Sigma}$. This is done by regressing out \mathbf{X} from **Equ. 5** and then performing factor analysis on the residual matrix $\mathbf{Y}P_{\mathbf{X}^T}^\perp$.
2. Compute $\hat{\tilde{\Xi}}^\perp$, an estimate for $\tilde{\Xi}P_{\mathbf{X}^T}^\perp$. Two possibilities for $\hat{\tilde{\Xi}}^\perp$ are the BLUP for $\tilde{\Xi}P_{\mathbf{X}^T}^\perp$ and the GLS estimate $\left(\hat{\tilde{\mathbf{L}}}^T \hat{\tilde{\Sigma}}^{-1} \hat{\tilde{\mathbf{L}}}\right)^{-1} \hat{\tilde{\mathbf{L}}}^T \hat{\tilde{\Sigma}}^{-1} \mathbf{Y}P_{\mathbf{X}^T}^\perp$.
3. Compute $\hat{\tilde{\Omega}}$, an estimate for $\tilde{\Omega}^{(\text{OLS})}$ by regressing $\mathbf{Y}\mathbf{X}^T(\mathbf{X}\mathbf{X}^T)^{-1}$ onto $\hat{\tilde{\mathbf{L}}}$. The methods differ in the way they find $\hat{\tilde{\Omega}}$.
 - (a) Houseman uses the the GLS estimate for $\tilde{\Omega}^{(\text{OLS})}$, $\hat{\tilde{\Omega}} = \left(\hat{\tilde{\mathbf{L}}}^T \hat{\tilde{\Sigma}}^{-1} \hat{\tilde{\mathbf{L}}}\right)^{-1} \hat{\tilde{\mathbf{L}}}^T \hat{\tilde{\Sigma}}^{-1} \mathbf{Y}\mathbf{X}^T(\mathbf{X}\mathbf{X}^T)^{-1}$. If $\mathbf{B} \neq \mathbf{0}$, this is sub-optimal since $\mathbb{E}[\mathbf{Y}\mathbf{X}^T(\mathbf{X}\mathbf{X}^T)^{-1} | \mathbf{X}, \tilde{\Xi}] = \mathbf{B} + \tilde{\mathbf{L}}\tilde{\Omega}^{(\text{OLS})}$, meaning the GLS estimate will be contaminated by the non-zero entries of \mathbf{B} .
 - (b) CATE, LEAPP use a robust regression to find $\hat{\tilde{\Omega}}$. This is a more appropriate estimator when $\mathbf{B} \neq \mathbf{0}$.
4. Estimate \mathbf{B} as

$$\hat{\mathbf{B}} = \mathbf{Y}\mathbf{X}^T[\mathbf{X}\mathbf{X}^T]^{-1} - \mathbf{Y}\left(\hat{\tilde{\Xi}}^\perp\right)^T\left[\hat{\tilde{\Xi}}^\perp\left(\hat{\tilde{\Xi}}^\perp\right)^T\right]^{-1}\hat{\tilde{\Omega}} \approx \mathbf{Y}\mathbf{X}^T[\mathbf{X}\mathbf{X}^T]^{-1} - \hat{\tilde{\mathbf{L}}}\hat{\tilde{\Omega}} \quad (10)$$

We then set

$$\hat{\text{Var}}(\hat{\beta}_i) = \hat{\sigma}_i^2 \left([\mathbf{X}\mathbf{X}^T]^{-1} + \hat{\tilde{\Omega}}^T \left[\hat{\tilde{\Xi}}^\perp \left(\hat{\tilde{\Xi}}^\perp \right)^T \right]^{-1} \hat{\tilde{\Omega}} \right) \quad (11)$$

These methods (CATE and LEAPP in particular) are attractive because one can use all of the CpG sites to estimate the confounding in the data without actually having to measure the confounding (if that is at all possible). They tend to perform well when the observed methylation data are informative enough to estimate the cell type effect, $\tilde{\mathbf{L}}$ (simulation plots when data are informative enough to estimate cell type?). However, since we must use the noisy estimate $\hat{\tilde{\mathbf{L}}}$ to estimate $\tilde{\Omega}^{(\text{OLS})}$, a term that controls the accuracy of any inference we do on \mathbf{B} (see **equ. 10** and **equ. 11**), we should try to understand methylation data are not informative enough to infer differentially expressed CpGs in the presence of unmeasured cell type heterogeneity. We are particularly concerned with underestimating $\tilde{\Omega}^{(\text{OLS})}$, as this will tend to lead to non-conservative inference on \mathbf{B} .

4 When can we use these methods to correct for cell type and other confounding?

The fidelity of any (is this word too strong) cell type correction method depends on how informative the available methylation data are for estimating cell type heterogeneity. If $\frac{1}{n}\mathbf{X}\mathbf{X}^T \rightarrow \Sigma_X > 0$ and $\tilde{\mathbf{L}}$ is fixed, Hastie (cite paper) proves that under suitable conditions on \mathbf{B} , $\tilde{\mathbf{L}}$ and Σ ,

$$\hat{\beta}_i \xrightarrow{P} \beta_i \quad (12)$$

and

$$n\hat{\text{Var}}(\hat{\beta}_i) = \hat{\sigma}_i^2 \left(\left[\frac{1}{n} \mathbf{X} \mathbf{X}^T \right]^{-1} + \hat{\Omega}^T \left[\frac{1}{n} \hat{\Xi}^\perp \left(\hat{\Xi}^\perp \right)^T \right]^{-1} \hat{\Omega} \right) \xrightarrow{P} \sigma_i^2 (\Sigma_X^{-1} + \Omega^T \Lambda^{-1} \Omega) \quad (13)$$

which is exactly what we would expect had we measured cell type. We also show through simulation that for n large enough and $\tilde{\mathbf{L}}$ fixed, the above procedure with robust regression tends to control false discovery rate at a nominal level (show plots).

However, the methylation data may not be informative enough to estimate and correct for the effects of cell type heterogeneity for small to moderate sample sizes. Recall that when we correct for cell type we must first compute $\hat{\tilde{\mathbf{L}}}$, and use that noisy estimate of $\tilde{\mathbf{L}}$ to regress $\mathbf{Y} \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1}$ onto $\hat{\tilde{\mathbf{L}}}$. If the residual $\hat{\tilde{\mathbf{L}}} - \tilde{\mathbf{L}}$ is large, then our estimate $\hat{\tilde{\Omega}}$ will be shrunk closer to $\mathbf{0}$. The intuition here is that if $\tilde{\mathbf{L}}$ was just a random matrix, the regression $\mathbf{Y} \mathbf{X}^T (\mathbf{X} \mathbf{X}^T)^{-1} \sim \hat{\tilde{\mathbf{L}}}$ will give regression coefficients that are very close to 0. If our estimate for $\hat{\tilde{\Omega}}$ is too small, we will underestimate the variance of $\hat{\beta}_i$ and any inference we do will be anti-conservative.

To study this phenomenon, we will assume $\mathbf{B} = \mathbf{0}$ and that there is only one covariate of interest (we assume that the covariate and response have been mean centered, so we may ignore the intercept and let $d = 1$). That is, variability in the data is caused only by cell type heterogeneity which may be correlated with the variable of interest:

$$\mathbf{Y} = \tilde{\mathbf{L}} \tilde{\Omega} \mathbf{X} + \tilde{\mathbf{L}} \tilde{\Xi} + \mathbf{E}. \quad (14)$$

If we let $\mathbf{E}^* = (\tilde{\mathbf{L}} - \hat{\tilde{\mathbf{L}}})(\tilde{\Omega} \mathbf{X} + \tilde{\Xi}) + \mathbf{E}$, we can rewrite 11 as

$$\mathbf{Y} = \hat{\tilde{\mathbf{L}}}(\tilde{\Omega} \mathbf{X} + \tilde{\Xi}) + \mathbf{E}^*. \quad (15)$$

If the residual $\hat{\tilde{\mathbf{L}}} - \tilde{\mathbf{L}}$ is large in comparison to $\hat{\tilde{\mathbf{L}}}$, we have no way of separating the effect $\hat{\tilde{\mathbf{L}}}(\tilde{\Omega} \mathbf{X} + \tilde{\Xi})$ from the error \mathbf{E}^* , which causes our estimate $\hat{\tilde{\Omega}}$ to shrink toward $\mathbf{0}$.

If cell type were known, (i.e. if $\tilde{\Omega} \mathbf{X} + \tilde{\Xi}$ were known), we can estimate $\tilde{\mathbf{L}}$ using OLS where the estimates $\hat{\tilde{\ell}}_i^{(\text{OLS})} \approx \tilde{\ell}_i \pm 1.96 \frac{\sigma_i}{\sqrt{n}}$. Surprisingly, even if we do not know cell type, we can use factor analysis to get an estimate $\hat{\tilde{\ell}}_i^{(\text{FA})} \approx \tilde{\ell}_i \pm 1.96 \frac{\sigma_i}{\sqrt{n}}$! Therefore, the effect $\hat{\tilde{\mathbf{L}}}(\tilde{\Omega} \mathbf{X} + \tilde{\Xi})$ is roughly equal to or smaller than the noise \mathbf{E}^* if, on average, $\tilde{\ell}_i \leq O\left(\frac{\sigma_i}{\sqrt{n}}\right)$ for CpG sites $i = 1, \dots, p$. Another way of saying this is $\frac{1}{p} \tilde{\mathbf{L}}^T \Sigma^{-1} \tilde{\mathbf{L}} \leq O\left(\frac{1}{n}\right)$. From here on out, we will define

$$\mathcal{I}_{\text{confounding}} = \frac{1}{p} \tilde{\mathbf{L}}^T \Sigma^{-1} \tilde{\mathbf{L}} \quad (16)$$

as the confounding information matrix. This is not to be confused with the Fisher Information matrix, although it has a similar role. $\mathcal{I}_{\text{confounding}}$ determines how much information about the confounding (in our case cell type heterogeneity) can be gleaned from the observed methylation data. When this information is on the order of the statistical error (i.e. when $\mathcal{I}_{\text{confounding}} \leq O\left(\frac{1}{n}\right)$), we have difficulty estimating and correcting for cell type heterogeneity.

Note that since we are assuming $\tilde{\Xi} \sim (\mathbf{0}, [I_K, I_n])$ in **equ. 14**, we can rotate $\tilde{\Xi}$ without changing the first and second moments. Therefore, the first and second moments of $\tilde{\mathbf{L}} \tilde{\Omega} \mathbf{X} + \tilde{\mathbf{L}} \tilde{\Xi} + \mathbf{E}$ and $(\tilde{\mathbf{L}} \mathbf{U})(\mathbf{U}^T \tilde{\Omega}) \mathbf{X} + (\tilde{\mathbf{L}} \mathbf{U})(\mathbf{U}^T \tilde{\Xi}) + \mathbf{E}$ are the same for any rotation matrix $\mathbf{U} \in \mathbb{R}^{K \times K}$. This means that it suffices to assume that

$$\mathcal{I}_{\text{confounding}} = \text{diag}(\delta_1, \dots, \delta_K) \text{ with } \delta_1 \geq \delta_2 \geq \dots \geq \delta_K > 0. \quad (17)$$

Lastly, since we have assumed $d = 1$, we can re-write **equ. 14** as

$$\mathbf{Y} = \tilde{\mathbf{L}} \tilde{\Omega}^{\text{OLS}} \mathbf{X} + \tilde{\mathbf{L}} \tilde{\Xi} P_{\mathbf{X}^\perp}^\perp + \mathbf{E} = \|\tilde{\Omega}^{\text{OLS}}\|_2 (\tilde{\mathbf{L}} \mathbf{Q} \mathbf{e}_1) \mathbf{X} + \tilde{\mathbf{L}} \tilde{\Xi} P_{\mathbf{X}^\perp}^\perp + \mathbf{E} \quad (18)$$

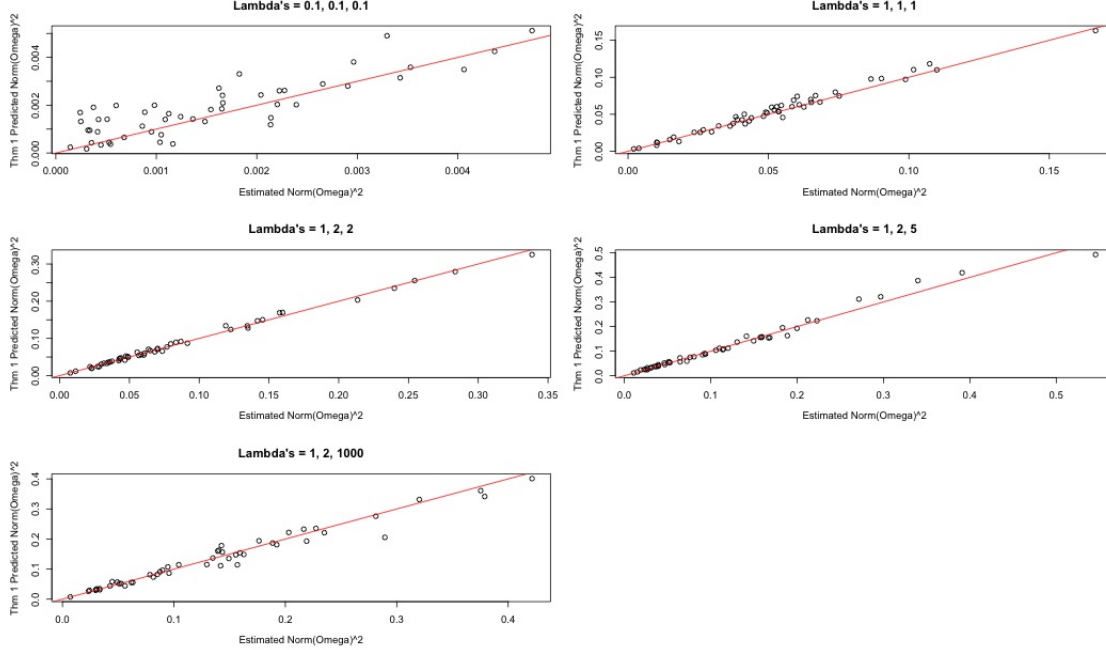
where $\tilde{\Omega}^{\text{OLS}} = \|\tilde{\Omega}^{\text{OLS}}\|_2 \mathbf{Q} \mathbf{e}_1$ is the QR decomposition of $\tilde{\Omega}^{\text{OLS}}$. From **equ. 11** and **equ. 18**, the accuracy of our inference is dependent on how well we estimate $\|\tilde{\Omega}^{\text{OLS}}\|_2$ (it also depend on $\mathbf{Q} \mathbf{e}_1$, but we will ignore this for now). When we underestimate $\|\tilde{\Omega}^{\text{OLS}}\|_2$, we tend to perform anti-conservative inference on \mathbf{B} .

We can consolidate these ideas into a single statement about the behavior of $\|\hat{\tilde{\Omega}}\|_2$ when the methylation data are not informative enough to accurately estimate the cell type effect $\tilde{\mathbf{L}}$:

Theorem 1 Suppose that $\tilde{\Omega} = (\omega_1, \dots, \omega_K)^T \neq 0$ and $I_{\text{confounding}} = \text{diag}(\delta_1, \dots, \delta_K)$ such that $\delta_i n \rightarrow \lambda_i$ for $i \geq k$ and $\delta_i n \rightarrow \infty$ for $i < k$. Then as $n, p \rightarrow \infty$,

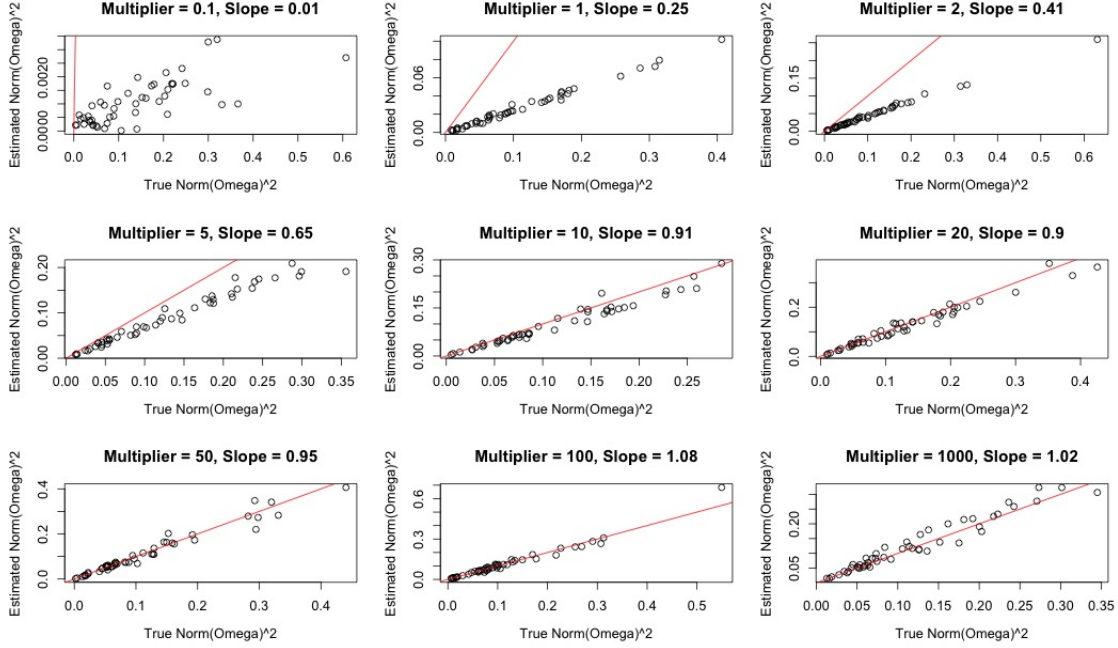
$$\|\hat{\tilde{\Omega}}\|_2^2 \rightarrow \sum_{j=1}^{k-1} \omega_j^2 + \sum_{i=k}^K \left(\frac{\lambda_i}{\lambda_i + 1} \right)^2 \omega_i^2 < \|\tilde{\Omega}\|_2^2 \quad (19)$$

I have a back-of-the-envelope proof of this. Rigorously proving this statement may take some work. Below is a proof of principle that leads me to believe the statement is correct:



(The above data were simulated assuming $\tilde{\Xi} \sim MN_{3 \times 100}(\mathbf{0}, I_3, I_{100})$, $E \sim MN_{p \times 100}(\mathbf{0}, \text{diag}(0.15, \dots, 0.15), I_{100})$. The elements of \tilde{L} are independent normals with variance chosen to get the correct $I_{\text{confounding}} = \text{diag}(\frac{\lambda_1}{100}, \frac{\lambda_2}{100}, \frac{\lambda_3}{100})$)

We can illustrate the effect of the size of $I_{\text{confounding}}$ on the extent of the shrinkage in $\hat{\tilde{\Omega}}$ with a simple simulation. We first fix the effect due to cell type, \tilde{L} , and the residual variability, Σ , s.t. $\frac{1}{p} \tilde{L}^T \Sigma^{-1} \tilde{L} = \frac{1}{n} I_K$. We also fix a $2 \times n$ design matrix $X = \begin{bmatrix} \text{Intercept} \\ \text{Disease Status (0's and 1's)} \end{bmatrix}$. We then simulate data according to **Equ. 11** by varying the observed correlation between cell type and covariate of interest ($\tilde{\Omega}^{(\text{OLS})}$) and see how significantly $\hat{\tilde{\Omega}}$ underestimates the observed correlation. Below are our results:



It is clear that the less informative the methylation data are for cell type (i.e. the smaller $\frac{1}{p} \tilde{\mathbf{L}}^T \Sigma^{-1} \tilde{\mathbf{L}}$ is), the more difficult it is to estimate the observed correlation between cell type and the covariate(s) interest.

Based on the above results, we can identify realistic scenarios when it will be difficult to estimate the correlation between cell type and the covariate(s) of interest. The first is when the non-standardized cell type effect \mathbf{L} is small, causing only subtle differences in methylation between cell types. This can occur when cells are taken from the same tissue with little observable cellular heterogeneity. Another related possibility is that the majority of CpG sites have the same methylation status across cell types and only a few CpG's have different methylation patterns across cells. If the cell type effects for only a few of the CpG sites are large, then underestimating the variance given by **equ. 11** can have serious repercussions on the accuracy of any false discovery procedure (need to add an FDR plot). A third may occur when the covariate(s) of interest explains the majority of the variability in cellular heterogeneity. For example, if our covariate of interest is asthma status and everyone with asthma has one cell type composition and everyone without asthma has a different cell type composition, then there is no way to distinguish the cell type effect from the asthma effect and it is impossible to estimate $\tilde{\mathbf{L}}$, regardless of the sample size. Mathematically speaking, this occurs when the residual variance $\mathbf{\Lambda}$ from **3** and **4** is small, meaning $\tilde{\mathbf{L}} = \mathbf{L}\mathbf{\Lambda}^{1/2}$ is small.

5 Real Data Example

Our motivating example came from whole-blood methylation data collected from individuals in the Amish and Hutterite communities, two founder populations with little genetic diversity. The goal of this experiment was to study any potential differences in the two population's methylome. We quantified the methylation from $p = 327,273$ CpG sites from 30 Amish and 30 Hutterite adolescents (ages 7-14) using 450K Illumina chips and measured the whole-blood cell composition for 29 of the Amish and all 30 of the Hutterite adolescents using flow cytometry. The design matrix for this experiment was

$$\mathbf{X} = \begin{bmatrix} \text{Intercept} \\ \text{Amish (1) or Hutterite (0)} \end{bmatrix} \quad (20)$$

and the cell type composition matrix was

$$\mathbf{C} = \begin{bmatrix} \text{Tcells} \\ \text{Bcells} \\ \text{Eosinophils} \\ \text{Neutrophils} \\ \text{Monocytes} \end{bmatrix}. \quad (21)$$

Since we have access to whole-blood cell composition for 59 of the 60 study subjects, we have the luxury of comparing inference on \mathbf{B} when cell type is and is not observed (we only use the $n = 59$ subjects with observed cell type to compare inference). When we assumed cell type was observed, we modeled the data as

$$\mathbf{Y}_{p \times n} = \mathbf{B}_{p \times d} \mathbf{X}_{d \times n} + \mathbf{L}_{p \times K} \mathbf{C}_{K \times n} + \mathbf{\Gamma}_{p \times r} \mathbf{V}_{r \times n} + \mathbf{E}_{p \times n} \quad (22)$$

$$\mathbf{B} = \begin{bmatrix} \mu_1 & b_{(A-H)_1} \\ \vdots & \vdots \\ \mu_p & b_{(A-H)_p} \end{bmatrix} \quad (23)$$

where the term $\mathbf{\Gamma}_{p \times r} \mathbf{V}_{r \times n}$ was included to take into account the effect of unmeasured covariates \mathbf{V} and μ_i and $b_{(A-H)_i}$ are the global mean and Amish-Hutterite contrast at CpG i . We further assumed that $\mathbf{V} \sim MN_{r \times n}(\boldsymbol{\alpha}_{r \times d} \mathbf{X}_{d \times n}, \mathbf{I}_r, \mathbf{I}_n)$. For simplicity we set $\boldsymbol{\alpha} = \mathbf{0}$, although the inference we did on \mathbf{B} would not change if we assumed $\boldsymbol{\alpha} \neq \mathbf{0}$ because we needed to account for any spurious correlation between the unmeasured covariates \mathbf{V} and \mathbf{X} . In order to correct for cell type heterogeneity in this scenario, we simply regressed out \mathbf{C} , meaning the cell type-corrected model was

$$\mathbf{Y}^{(c)} = \mathbf{B} \mathbf{X}^{(c)} + \mathbf{\Gamma} \mathbf{V}^{(c)} + \mathbf{E}^{(c)} \quad (24)$$

$$\mathbf{Y}^{(c)} = \mathbf{Y} P_{\mathbf{C}^\top}^\perp, \quad \mathbf{V}^{(c)} = \mathbf{V} P_{\mathbf{C}^\top}^\perp, \quad \mathbf{E}^{(c)} = \mathbf{E} P_{\mathbf{C}^\top}^\perp \quad (25)$$

We estimated $r = 4$ additional unobserved factors using BCV (cite Wang, Hastie paper) and used CATE to compute

any spurious correlation between $\mathbf{V}^{(c)}$ and $\mathbf{X}^{(c)}$, as well as $\hat{\mathbf{B}} = \begin{bmatrix} \hat{\mu}_1^{(\text{cell})} & \hat{b}_{(A-H)_1}^{(\text{cell})} \\ \vdots & \vdots \\ \hat{\mu}_p^{(\text{cell})} & \hat{b}_{(A-H)_p}^{(\text{cell})} \end{bmatrix}$. Our test statistic used to infer the difference in methylation at CpG site i between Amish and Hutterite adolescents was then

$$\hat{t}_i^{(\text{cell})} = \frac{\hat{b}_{(A-H)_i}^{(\text{cell})}}{\sqrt{\hat{\text{Var}}(\hat{b}_{(A-H)_i}^{(\text{cell})})}} \quad (26)$$

where $\hat{\text{Var}}(\hat{b}_{(A-H)_i}^{(\text{cell})})$ was computed according to **equ. 11**. We used a t -distribution with $n - K - d - r = 48$ degrees of freedom to compute p-values at each CpG site $i = 1, \dots, p$ to test for differential methylation.

When cell type is UNOBSERVED, we combined unmeasured covariates \mathbf{C} and \mathbf{V} in **equ. 22** and modeled the data as

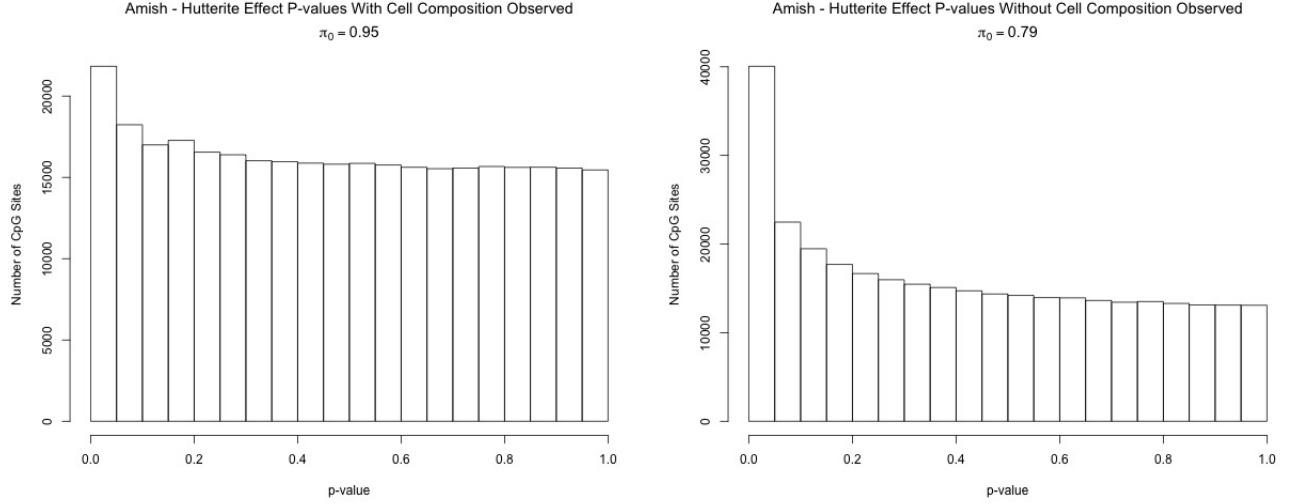
$$\mathbf{Y}_{p \times n} = \mathbf{B}_{p \times d} \mathbf{X}_{d \times n} + \tilde{\mathbf{L}}_{p \times K'} \tilde{\boldsymbol{\Omega}}_{K' \times d} \mathbf{X}_{d \times n} + \tilde{\mathbf{L}}_{p \times K'} \tilde{\boldsymbol{\Xi}}_{K' \times n} + \mathbf{E}_{p \times n} \quad (27)$$

where $\tilde{\boldsymbol{\Omega}} = \begin{bmatrix} \Lambda^{-1/2} \boldsymbol{\Omega} \\ \boldsymbol{\alpha} \end{bmatrix}$, $\tilde{\mathbf{L}} = \begin{bmatrix} \mathbf{L} \Lambda^{1/2} & \mathbf{\Gamma} \end{bmatrix}$ and $\tilde{\boldsymbol{\Xi}} = \begin{bmatrix} \Lambda^{-1/2} \boldsymbol{\Xi} \\ \mathbf{V} \end{bmatrix}$. We used the exact same procedure as above to estimate

$K' = 6$, $\hat{\tilde{\boldsymbol{\Omega}}}$ and $\hat{\tilde{\mathbf{B}}} = \begin{bmatrix} \hat{\mu}_1^{(\text{no cell})} & \hat{b}_{(A-H)_1}^{(\text{no cell})} \\ \vdots & \vdots \\ \hat{\mu}_p^{(\text{no cell})} & \hat{b}_{(A-H)_p}^{(\text{no cell})} \end{bmatrix}$. And just as we did when we observed each individual's whole-blood cell composition, we let

$$\hat{t}_i^{(\text{no cell})} = \frac{\hat{b}_{(A-H)_i}^{(\text{no cell})}}{\sqrt{\hat{\text{Var}}(\hat{b}_{(A-H)_i}^{(\text{no cell})})}} \quad (28)$$

and computed p-values to test for differential methylation between both groups using a t -distribution with $n - d - K' = 51$ degrees of freedom. Below are p-value histograms that compare our results when we assume cell composition are and are not observed:



The stark difference between these results is because of the correlation between whole-blood cell composition and group status. The table below gives the standardized regression coefficients after regressing cell type C onto Amish vs. Hutterite group status X :

Cell Type	Amish - Hutterite Standardized Effect
T cell	-0.62
B cell	0.65
Eosinophil	-0.40
Neutrophil	1.41
Monocyte	0.65

Since cell heterogeneity and group status are so correlated, there is little variability left in X to estimate B after we correct for cell type (i.e. we are in the scenario when $\Lambda^{1/2}$ is small), which in turn cause our test statistics to move closer to 0. We also estimated \tilde{L} and found that $\left(\frac{\frac{1}{p}\tilde{L}^T\hat{\Sigma}^{-1}\tilde{L}}{\left|\frac{1}{59}I_5\right|}\right)^{1/5} \approx 1.3$. This implies that if \tilde{L} is non-zero and cellular heterogeneity were not observed, the observed CpG methylation data would not be informative enough to estimate cell type and we would perform considerable anti-conservative inference. However, just because $\left|\frac{1}{p}\tilde{L}^T\hat{\Sigma}^{-1}\tilde{L}\right|$ is slightly larger than $\left|\frac{1}{59}I_5\right|$ does not mean the inference we do without measured cell type is anti-conservative. In fact, if $\tilde{L} = 0$, we would expect $\frac{1}{p}\tilde{L}^T\hat{\Sigma}^{-1}\tilde{L} \approx \frac{1}{59}I_5$. Therefore, we cannot say with 100% confidence that the results without cell type measurements are anti-conservative. We can only use our prior understanding of methylation patterns across different cell types and assume that $\tilde{L} \neq 0$. If this were in fact the case, our sample size is not large enough to perform accurate inference without having observed cell composition.

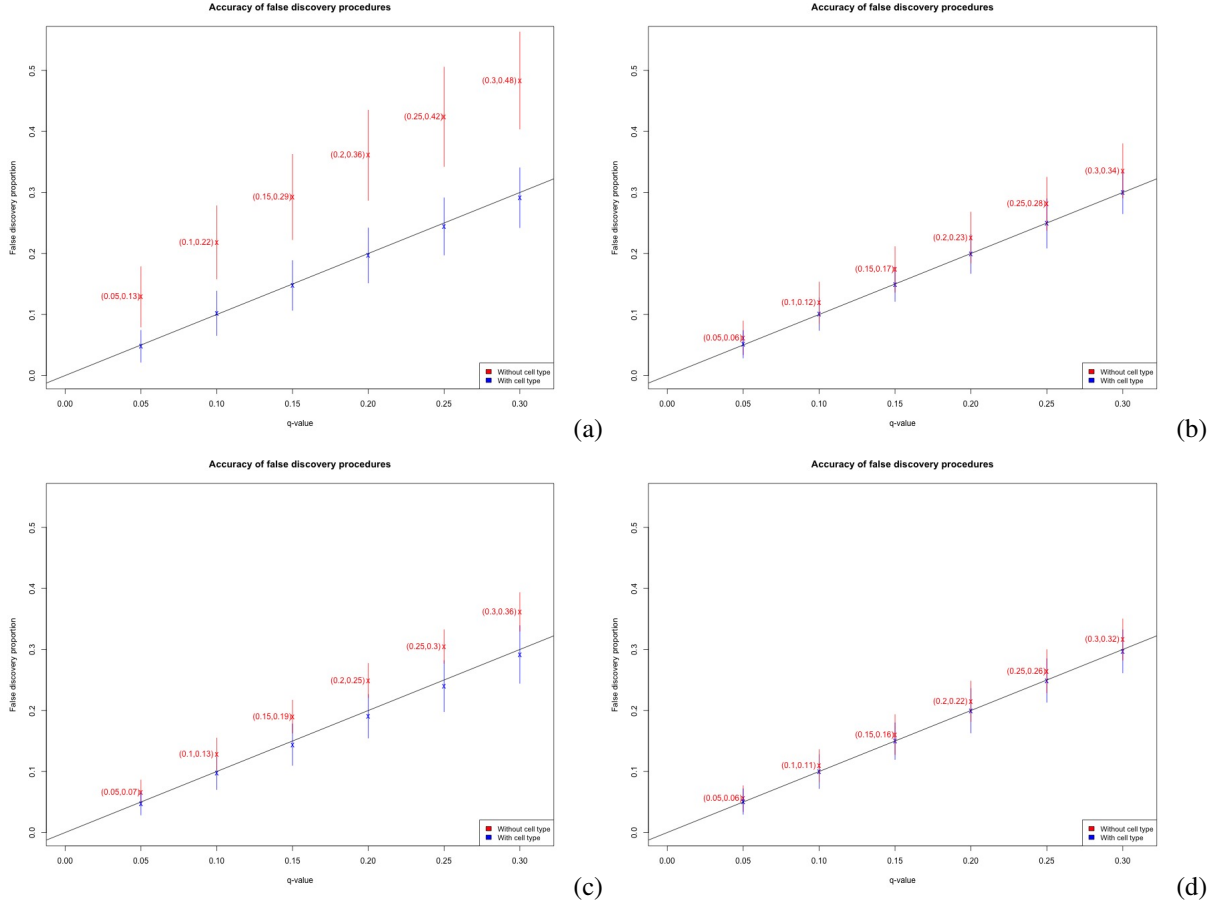
6 Simulation Study

Incorrectly estimating the correlation between cell type heterogeneity and the covariate(s) of interest can have deleterious effects on any procedure that attempts to control for false discovery. To explore this, we simulated data according to **equ. 5** with parameters based on estimates in the Amish-Hutterite study (see below table). We partitioned $n = 100$ individuals into two groups, 50 of them Amish and 50 of them Hutterites. In all of our simulations, methylation differed in only one cell type, i.e. $K = 1$.

Simulation ID	Amish-Hutterite Stand. Effect, B	Cell Type Stand. Effect, L	$\mathcal{I}_{\text{confounding}} = \frac{1}{p} L^T \Sigma^{-1} L$	Observed Correlation btwn. X and C , $\tilde{\Omega}^{\text{OLS}}$
(a)	$0.95\delta_0 + 0.05N(0, 0.4)$	$0.7\delta_0 + 0.3N(0, \tau_1)$	$\frac{1.3}{p} I_1$	1.41
(b)	$0.95\delta_0 + 0.05N(0, 0.4)$	$0.7\delta_0 + 0.3N(0, \tau_3)$	$\frac{1.3 \times 10}{p} I_1$	1.41
(c)	$0.95\delta_0 + 0.05N(0, 0.4)$	$0.7\delta_0 + 0.3N(0, \tau_2)$	$\frac{1.3}{10p} I_1$	1.41
(d)	$0.95\delta_0 + 0.05N(0, 0.4)$	$0.7\delta_0 + 0.3N(0, \tau_4)$	$\frac{1.3}{n} I_1$	1.41/3

The above table gives the parameters for the four different simulations, where the variances τ_i were chosen to get the correct $\mathcal{I}_{\text{confounding}}$. The below figures plot Storey's q-value against the true false discovery proportion

$$FDP(q) = \frac{\#\{\text{False positives with q-values} \leq q\}}{\#\{\text{CpG sites with q-values} \leq q\}} \quad (29)$$



It is not hard to see that when the correlation $\tilde{\Omega}^{\text{OLS}}$ is large and we have little power to estimate it (plot (a)), we incur significantly more false positives than we would expect using a reasonable false discovery control procedure. The reason for this is two fold. First, when we underestimate $\tilde{\Omega}^{\text{OLS}}$ our estimator for \hat{B} is significantly biased and second, we will underestimate the variance of the components of \hat{B} (see **equ. 9**). However, if we have enough statistical power to accurately estimate $\tilde{\Omega}^{\text{OLS}}$ (plot (b)), our false discovery control procedure is accurate. We also note that if the cell type effect is very small in comparison to the effect of interest (simulation (c)), we significantly underestimate $\tilde{\Omega}^{\text{OLS}}$ but still can accurately control false discovery. This is because the effect $\tilde{L}\tilde{\Omega}^{\text{OLS}}X$ is much smaller in comparison to BX , so we tend to find true non-zero Amish-Hutterite effects before Amish-Hutterite effects that are just due to differences in cell composition. Similar reasoning explains why we can control for false discovery when $\tilde{\Omega}^{\text{OLS}}$ is small (as in simulation (d)). Even though we underestimate $\tilde{\Omega}^{\text{OLS}}$, $\tilde{\Omega}^{\text{OLS}} - \hat{\Omega}$ is innocuously small.

7 Partially Observed Cell Type

8 What I need to add

- My model with partially observed cell type. I only have simulation results here. I should include a brief discussion about the caveats with real data, since additional batch effects may be difficult to estimate.
- A discussion when the residual E is not normally distributed. The confidence we have in factor analysis will depend on the higher order cumulants (notably the fourth kumulant) when the residuals are not normally distributed. This is important when we attempt to diagnose datasets with confounding we cannot estimate (i.e. how much does $\frac{1}{p}\hat{\Gamma}^T\hat{\Sigma}^{-1}\hat{\Gamma}$ differ from $\frac{1}{p}\Gamma^T\Sigma^{-1}\Gamma$ as a function of the distribution of E).

9 Junk

$X^T = [Q_1 \quad Q_2] \begin{bmatrix} R_{d \times d} \\ 0_{n-d \times d} \end{bmatrix}$ be the QR decomposition of X^T . $W_{r \times n} \sim MN_{r \times n}(\alpha_{r \times d} X_{d \times n}, I_r, I_n)$ be a random effects matrix that represents additional confounding in the experiment that may be correlated with the covariates of interest. In most data sets, α is small or 0, i.e. the correlation between the confounders and the covariates of interest is small. For simplicity, we assume $\alpha = 0$, although the model and methods can be easily extended for $\alpha \neq 0$. $\tilde{\Gamma} = \begin{bmatrix} \Gamma & L\Lambda^{1/2} \end{bmatrix}$, $\tilde{W} = \begin{bmatrix} W \\ \Lambda^{-1/2}\Xi \end{bmatrix}$ and $\tilde{\Omega} = \begin{bmatrix} 0_{r \times d} \\ \Lambda^{-1/2}\Omega \end{bmatrix}$. Note that $\tilde{W} \sim MN_{(r+K) \times n}(\mathbf{0}, I_{r+K}, I_n)$.