

Suppose we have methylation data for p CpG sites on the logit scale (i.e. $\text{logit}(\text{beta value})$) for two sets of individuals, Y_1, Y_2 , where $Y_i \in \mathbb{R}^{p \times n_i}$ and n_i the number of individuals in each group. Let $X_i \in \mathbb{R}^{d \times n_i}$ be the covariate matrix for d covariates of interest (i.e. sex, asthma status, etc.) and $C_i \in [0, 1]^{K \times n_i}$ be the cell type matrix for $K + 1$ cell types ($K + 1$ since X_i has an intercept). If $c_s^{(i)}$ is the s^{th} column of C_i , then $\mathbf{1}^T c_s^{(i)} < 1$. The data are then modeled linearly as

$$Y_i = B_{p \times d} X_i + L_{p \times K} C_i + E_i, \quad E_i \sim MN_{p \times n_i}(0, \Sigma_{p \times p}, I_{n_i})$$

In our set up, we observed X_1, X_2 and C_1 but NOT C_2 . For now, we assume that the columns of C_i follow a dirichlet distribution with

$$C_i | X_i, \Omega \sim \underbrace{\left(\text{Dir}(\alpha \Omega x_{1_1})^T, \dots, \text{Dir}(\alpha \Omega x_{1_{n_i}})^T \right)}_{n_i \text{ independent Dirichlet distributions}}$$

where $\Omega \in \mathbb{R}^{K \times d}$ and $\alpha > 0$ is a concentration parameter. Note that

$$E(C_i | X_i, \Omega) = \Omega X_i^T$$

$$\text{Var}(c_s^{(i)} | X_i, \Omega) = \frac{1}{\alpha + 1} \left(\text{diag}(\Omega x_{i_s}) - \Omega x_{i_s} (\Omega x_{i_s})^T \right)$$

and

$$E(Y_1 | X_1, X_2, C_1) = B X_1 + L C_1$$

$$E(Y_2 | X_1, X_2, C_1) = B X_2 + L E(\Omega | X_1, X_2, C_1) X_2$$

0.1 Estimating Ω

Let $l(\Omega, \alpha | C_1) = \log p(C_1 | \Omega, \alpha)$ and let $\hat{\Omega}, \hat{\alpha} = \hat{\Omega}(C_1), \hat{\alpha}(C_1)$ be the MLE. In order to estimate B and L I condition on C_1 , which meaning there is no randomness in $\hat{\Omega}, \hat{\alpha}$ in the usual frequentist setting. In order to introduce uncertainty back into the estimates $\hat{\Omega}, \hat{\alpha}$, I assume there is a prior $p(\Omega, \alpha)$ with compact support on Ω, α . If we define $h = (\text{vec}(\Omega), \alpha)$, we have

$$\log p(\Omega, \alpha | C_1) = \underbrace{\text{Const} + l(\Omega, \alpha | C_1)}_{O_p(n_1^{1/2})} + \underbrace{\log p(h)}_{O_p(1)} = \text{Const} + l(\hat{h} | C_1) - \underbrace{\frac{1}{2} (h - \hat{h})^T H(\hat{h}) (h - \hat{h})}_{\text{Kernel of a } N(\hat{h}, H(\hat{h})^{-1})} + o_p(\|h - \hat{h}\|^2) + \log p(h)$$

where $H = -\nabla_{hh}^2 l|_{\hat{h}} \approx I(\hat{h})$, the Fisher information at \hat{h} . Since $\frac{1}{2} (h - \hat{h})^T H(\hat{h}) (h - \hat{h})$ dominates the above expression for large n_1 , we may approximate $p(h | C_1)$ with a normal distribution:

$$p(h | C_1) \approx N(\hat{h}, I(\hat{h})^{-1}).$$

This is exactly the Bayesian central limit theorem.

For our asymptotic results, we only need the asymptotic first and second moments of $p(h | C_1)$ (and possibly the assumption that higher moments exist, which is why I assume the prior for h has compact support), along with $n_1 = O(n_2)$. The intuition for this is that the variance matrix for the second set of individuals will have n_2^2 entries, where the off-diagonal elements all are on the order $\frac{1}{n_1}$ (from the variance of $\Omega, \alpha | C_1$). In order to ensure that the maximum eigenvalue of this matrix doesn't explode as $n_2 \rightarrow \infty$, we require $n_1 = O(n_2)$. When this is the case, I am pretty sure I can prove consistency of the the quasi-likelihood estimator.

0.2 Quasi-Likelihood Estimator

Since the Normal + Dirichlet is a difficult likelihood to work with, I use a quasi-likelihood estimator. For each CpG site $g = 1, \dots, p$, let y_{i_g} be the g^{th} row of Y_1 , ℓ_g the g^{th} column of L and β_g the g^{th} column of B . Let $\tilde{\beta}_g = \begin{pmatrix} \beta_g \\ \ell_g \end{pmatrix}$. The model for these individual sites is

$$y_{1_g} = X_1^T \beta_g + C_1^T \ell_g + \epsilon_{1_g}, \quad \epsilon_{1_g} \sim N(0, \sigma_g^2)$$

$$y_{2_g} = X_2^T \beta_g + C_2^T \ell_g + \epsilon_{2_g}, \quad \epsilon_{2_g} \sim N(0, \sigma_g^2)$$

where ϵ_{i_g} is independent of C_2 . The **means** of the random variables are

$$\begin{aligned} \mu_{1_g} &= \mu_{1_g}(\beta_g, \ell_g) = E(y_{1_g} | X_1, X_2, C_1) = X_1^T \beta_g + C_1^T \ell_g \\ \mu_{2_g} &= \mu_{2_g}(\beta_g, \ell_g) = E(y_{2_g} | X_1, X_2, C_1) = X_2^T \beta_g + \underbrace{E(C_2 | X_2, C_1)^T}_{=X_2^T(\hat{\Omega}^T + O(\frac{1}{\sqrt{n_1}}))} \ell_g. \end{aligned}$$

The **variances** are

$$\begin{aligned} \text{Var}(y_{1_g} | X_1, X_2, C_1) &= \sigma_g^2 I_{n_1} \\ \text{Var}(y_{2_g} | X_1, X_2, C_1) &= \text{Var}(C_2^T \ell_g | X_2, X_1, C_1) + \sigma_g^2 I_{n_2} = G(\ell_g). \end{aligned}$$

Let c_s and c_t be the s and t^{th} columns of C_2 , respectively. Define

$$R_s = \frac{1}{\alpha + 1} \left(\text{diag}(\Omega x_{2_s}) - \Omega x_{2_s} (\Omega x_{2_s})^T \right) \text{ and } V_s = E_{\Omega, \alpha | C_1, X_1, X_2} R_s = O(1).$$

$$S_{ss} = \underbrace{\text{Var}_{h|C_1, X_1, X_2}(E(c_s | C_1, X_1, X_2, \Omega))}_{\text{Variance taken over } p(h | C_1, X_1, X_2)} = O\left(\frac{1}{n_1}\right)$$

$$S_{st} = \text{Var}_{h|C_1, X_1, X_2}(E(c_s | C_1, X_1, X_2, \Omega), E(c_t | C_1, X_1, X_2, \Omega)) = O\left(\frac{1}{n_1}\right)$$

Finally, we can get an analytic expression for $G(\ell_g)$ by noting that

$$\text{Var}(C_2^T \ell_g | X_2, X_1, C_1) = \begin{pmatrix} \ell_g^T (V_1 + S_{11}) \ell_g & \cdots & \ell_g^T S_{n_2 1} \ell_g \\ \vdots & \ddots & \vdots \\ \ell_g^T S_{1 n_2} \ell_g & \cdots & \ell_g^T (V_{n_2} + S_{n_2 n_2}) \ell_g \end{pmatrix} \in \mathbb{R}^{n_2 \times n_2}$$

Next, we define

$$U = \begin{bmatrix} X_1 & X_2 \\ C_1 & E(C_2 | X_1, X_2, C_1) \end{bmatrix} \begin{pmatrix} \frac{1}{\sigma_g^2} (y_{1_g} - \mu_{1_g}) \\ G(\ell_g)^{-1} (y_{2_g} - \mu_{2_g}) \end{pmatrix}$$

and

$$\begin{aligned} T &= -E\left(\frac{dU}{d\beta} | X_1, X_2, C_1\right) = \text{Var}(U | X_1, X_2, C_1) = \begin{bmatrix} X_1 & X_2 \\ C_1 & E(C_2 | X_1, X_2, C_1) \end{bmatrix} \begin{pmatrix} \frac{1}{\sigma_g^2} I_{n_1} & 0 \\ 0 & G(\ell_g)^{-1} \end{pmatrix} \begin{bmatrix} X_1 & X_2 \\ C_1 & E(C_2 | X_1, X_2, C_1) \end{bmatrix}^T = \\ &= \underbrace{\begin{bmatrix} X_1 & X_2 \\ C_1 & \hat{\Omega} X_2 \end{bmatrix} \begin{pmatrix} \frac{1}{\sigma_g^2} I_{n_1} & 0 \\ 0 & G(\ell_g)^{-1} \end{pmatrix} \begin{bmatrix} X_1 & X_2 \\ C_1 & \hat{\Omega} X_2 \end{bmatrix}^T}_{\lambda_{\min} > \delta, \lambda_{\max} < Const} \left(1 + O_P\left(\frac{1}{\sqrt{n_1}}\right) + O_P\left(\frac{1}{n_1}\right)\right) \approx \underbrace{\begin{bmatrix} X_1 & X_2 \\ C_1 & \hat{\Omega} X_2 \end{bmatrix} \begin{pmatrix} \frac{1}{\sigma_g^2} I_{n_1} & 0 \\ 0 & G(\ell_g)^{-1} \end{pmatrix} \begin{bmatrix} X_1 & X_2 \\ C_1 & \hat{\Omega} X_2 \end{bmatrix}^T}_{\lambda_{\min} > \delta, \lambda_{\max} < Const}. \end{aligned}$$

The **Newton Updates** are then

$$\hat{\beta}_{k+1} = \hat{\beta}_k + T_k^{-1} U_k.$$

0.3 What I need to Prove

I still need prove that the estimator $\hat{\beta} \xrightarrow{P} \tilde{\beta}$ when $n_1 = O(n_2)$ and under what conditions we have asymptotic normality. Once I prove these, I already have a proof showing that if consistency and asymptotic normality hold when σ_g^2 is known, it also holds when it uses the OLS plugin estimator for σ_g^2 based on the n_1 individuals.

0.4 Other Possible Models

In the above model I assumed that the methylation response Y are M-values (i.e. $\text{logit}\left(\frac{\# \text{methylated} + \alpha}{\# \text{methylated} + \# \text{unmethylated} + \alpha}\right)$, α recommended by Illumina). The data we collect are actually count data, so a more appropriate model for the number of methylation events is binomial, i.e. if y_{ig} is the number of observed methylated residues for site g and n_{ig} the total number of times we observe site g for individual i ,

$$y_{ig} \sim \text{Bin}(n_{ig}, \pi_{ig}), \pi_{ig} = f(\text{covariates for person } i, \tilde{\beta}_g)$$

We can use the quasi-likelihood method above with a modification to the variance, assuming a logit link function. That is, if $V = V(\tilde{\beta}_g) = \begin{pmatrix} \sigma_g^2 I_{n_1} & 0 \\ 0 & G(\ell_g) \end{pmatrix}$, the new variance we would use would be

$$V_{\text{Bin}} = \Gamma^{1/2} V \Gamma^{1/2}$$

where $\Gamma = \text{diag}(\text{var}(\mu_1), \dots, \text{var}(\mu_{n_1+n_2}))$. Note that for $n_2 = 0$, this is just the dispersed binomial variance. This type of variance modification was used in McPeck, 2016 to correct for population specific effects in binary response data.

The only benefit to using M-values that we can better correct batch effects present in the data.

0.5 Review of Current Methods

To my knowledge there are three current methods to correct for cell type in methylation experiments, 2 by Houseman and 1 by Zhou.

1. **Houseman, 2012:** Houseman uses a training set consisting of methylation measured in cell types sorted by flow cytometry. The user can then feed in their methylation data matrix (on the beta scale) and get back the predicted cell types for each individuals. The problems with this method are:
 - (a) Predicted cell type depends on population. Michelle compared her measured cell proportions with Houseman's predicted values and found that the flow cytometry data was very different that Houseman's predictions.
 - (b) Houseman does not provide any sort of uncertainty in the prediction, only a point estimate. This is important in assessing the uncertainty in our parameter estimates.
 - (c) Their model does not use covariates as information in deconvolving cell type. If two people have similar covariates, then they should also have similar cell types.
2. **Houseman, 2014:** Houseman uses an unsupervised method to estimating regression coefficients in the presence of unmeasured cell type confounders. The shortcomings of this method are:
 - (a) One is forced to make restrictive assumptions about the sparsity of B and $L\Omega$. In an independent paper, Hastie 2016 shows that in order to get consistent estimators of B in the unsupervised regime, one needs to make the assumption that the sparsity of B INCREASES as the sparsity of L INCREASES. That is, if the unobserved cell type does not explain a large portion of the variability in methylation (which we have seen happen, e.g. Michelles Hutterite data), then one has no hope of estimating B . In fact, whenever cell type DOES NOT depend on your covariates you drastically reduce your power to do inference on B when you do not observe cell type.
 - (b) Even when $L\Omega$ is dense and B is sparse, Houseman uses a poor estimator for \hat{B} . The asymptotic properties of his estimator are not known and it is unrealistic to get p-values for individual sites. The best we can do is get site-specific confidence intervals. We can actually improve upon this method by using a different estimator.
3. **Zhou, 2014:** Zhou uses a linear regression treating methylation as the covariate and disease status as the response. He first forms that individual relatedness matrix Φ using the centered methylation matrix as a proxy for individual relatedness. He then iteratively adds principle components of Φ until only a small fraction of sites are significantly correlated with the phenotype. The problems here are:
 - (a) The method is very adhoc and the statistical properties are not very well understood.

- (b) You can only look at one covariate at a time. You cannot test the hypothesis that methylation is uncorrelated with ≥ 2 covariates at once.
- (c) You lose power when performing OLS with binary response (in comparison to a GLM).

The criticism of all of the above methods is one has no idea if the assumptions you are making on the contribution of the effect due to cell type are correct.

I think the method I describe above is important for those performing large DNA methylation experiments when experimenters are unsure how methylation changes across cell type and how cell type varies across covariates. Michelle has Hutterite/Amish data at both extremes:

1. 30 Hutterite children ages 7-14 where cell type appears to be INDEPENDENT of the primary covariates (age, asthma, sex). If Houseman or Zhou were used here, we would incorrectly remove the majority of the effect.
2. Hutterite and Amish individuals whose cell type differs due to community, but no other covariates. No unsupervised model would be able to account for this type of data structure.