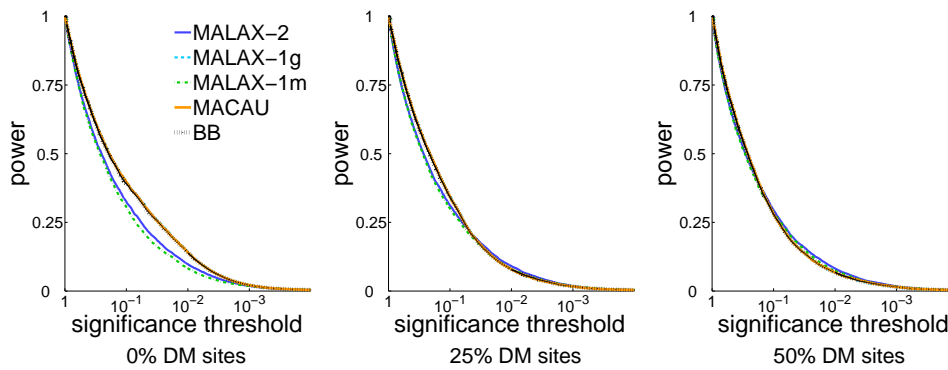
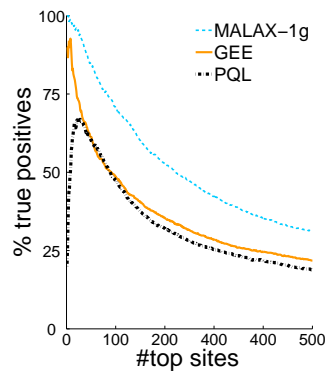


Association testing in bisulfite sequencing data via a Laplace approximation: Supplementary Information

1 Supplementary Figures



Supplementary Figure S1: Power measures for the evaluated methods under various proportions of differentially methylated sites. All results are averaged over 10 simulated data sets.



Supplementary Figure S2: The detection power of MALAX-1g and two additional evaluated methods under simulated data sets with no DM sites. The additional methods are GEE (based on CARAT) and PQL (based on GMMAT), as described in the supplementary note. All results are averaged over 10 simulated data sets.

2 Experiments with Additional Methods

We evaluated two additional methods for generalized linear mixed model (GLMM) approximation. Both methods were recently proposed for case-control association studies, but can be readily adapted to work with count-bases responses instead of binary responses.

The first method is based on the penalized quasi likelihood algorithm, as implemented in the GMMAT package [1]. This method is similar to the Laplace approximation but applies additional approximations, which results in faster computations at the cost of reduced accuracy.

The second method uses a generalized estimating equations (GEE) approach, as implemented in the CARAT package [2]. This method only models the first two moments of the model likelihood. We used a modified version of CARAT that was adapted to use a binomial instead of a binary response. We evaluated both a standard version and a version with an additional dispersion parameter, as described in [2]; both versions yielded effectively the same results.

GMMAT and CARAT were evaluated in a setting with no differentially methylated (DM) sites, which was the easiest setting in our experiments. Consequently, both methods used one variance component associated with a genetic similarity matrix, along with one variance component associated with the identity matrix, which can account for independent over-dispersion. The results demonstrate that MALAX-1g substantially outperformed both methods (Supplementary Figure S2). We therefore did not consider these methods in the remainder of the experiments. We note that we also evaluated versions of these approaches that use beta-binomial instead of binomial responses, but these versions were less accurate than the binomial ones.

3 Gradient Computation

The gradient of the approximate log likelihood described in the main text is required both for approximating the Hessian and for the maximum likelihood estimation procedure. Here we derive the gradient computation in detail. We first explicitly write the approximate log likelihood as follows:

$$\log P(\mathbf{y}^j | \mathbf{x}, \mathbf{W}, \mathbf{r}^j) \approx \underbrace{-\frac{1}{2}(\hat{\mathbf{l}} - \mathbf{m})^T \mathbf{G}^{-1} (\hat{\mathbf{l}} - \mathbf{m})}_{L_1} + \underbrace{\sum_{i=1}^n \log P(y_i^j | \hat{l}_i)}_{L_2} - \underbrace{\frac{1}{2} \log |\mathbf{B}|}_{L_3}, \quad (1)$$

where $\mathbf{B} \triangleq \mathbf{G}\mathbf{A}$. We denote the three terms on the right hand side of the above equation as L_1 , L_2 and L_3 , respectively. Note that the quantities \mathbf{G} and \mathbf{m} in L_1 depend explicitly on the model parameters, but $\hat{\mathbf{l}}$ and \mathbf{A} also implicitly depend on these parameters, where the dependence of \mathbf{A} is mediated entirely through its dependence on $\hat{\mathbf{l}}$. We therefore divide the partial derivative according to each parameter θ (which can represent variance components or fixed effects) into its explicit and implicit components, by using the chain rule as follows:

$$\frac{\partial \log P(\mathbf{y}^j | \mathbf{x}, \mathbf{W}, \mathbf{r}^j)}{\partial \theta} = \frac{\partial L_1}{\partial \theta} \Big|_{\text{explicit}} + \sum_{i=1}^n \frac{\partial \log P(\mathbf{y}^j | \mathbf{x}, \mathbf{W}, \mathbf{r}^j)}{\partial \hat{l}_i} \frac{\partial \hat{l}_i}{\partial \theta}. \quad (2)$$

We first derive the explicit components:

$$\begin{aligned}
\frac{\partial L_1}{\partial \sigma_v^2} \Big|_{\text{explicit}} &= -\frac{1}{2} (\hat{\mathbf{l}} - \mathbf{m})^T \mathbf{G}^{-1} \mathbf{K}_v \mathbf{G}^{-1} (\hat{\mathbf{l}} - \mathbf{m}) - \frac{1}{2} \text{tr} [\mathbf{G}^{-1} \mathbf{K}_v] \\
\frac{\partial L_1}{\partial \sigma_e^2} \Big|_{\text{explicit}} &= -\frac{1}{2} (\hat{\mathbf{l}} - \mathbf{m})^T \mathbf{G}^{-2} (\hat{\mathbf{l}} - \mathbf{m}) - \frac{1}{2} \text{tr} [\mathbf{G}^{-1}] \\
\frac{\partial L_1}{\partial \boldsymbol{\gamma}} \Big|_{\text{explicit}} &= -(\hat{\mathbf{l}} - \mathbf{m})^T \mathbf{G}^{-1} \mathbf{C}.
\end{aligned} \tag{3}$$

Here, \mathbf{C} is the matrix of covariates for the entire sample (the matrix \mathbf{W} with an appended column \mathbf{x}), and $\boldsymbol{\gamma} = [\boldsymbol{\alpha}^T \beta]^T$ is the vector of all fixed effects.

To derive the implicit components, we first note that $\frac{\partial(L_1+L_2)}{\partial \mathbf{l}} \Big|_{\mathbf{l}=\hat{\mathbf{l}}} = \mathbf{0}$ by definition. We therefore only need to compute $\frac{\partial \log |\mathbf{B}|}{\partial \hat{l}_i} \frac{\partial \hat{l}_i}{\partial \theta}$. To derive $\log |\mathbf{B}|$, we first explicitly write the negative Hessian \mathbf{A} as follows:

$$\mathbf{A} \triangleq -\nabla \nabla \log P(\mathbf{l} \mid \mathbf{x}, \mathbf{W}, \mathbf{y}^j, \mathbf{r}^j) \Big|_{\mathbf{l}=\hat{\mathbf{l}}} = -\nabla \nabla \log P(\mathbf{y}^j \mid \mathbf{l}) \Big|_{\mathbf{l}=\hat{\mathbf{l}}} + \mathbf{G}^{-1}. \tag{4}$$

Using this explicit notation, $\frac{\partial \log |\mathbf{B}|}{\partial \hat{l}_i}$ is given by:

$$\frac{\partial \log |\mathbf{B}|}{\partial \hat{l}_i} = -\frac{1}{2} \text{tr} \left[\mathbf{B}^{-1} \frac{\partial \mathbf{B}}{\partial \hat{l}_i} \Big|_{\mathbf{l}=\hat{\mathbf{l}}} \right] = \frac{1}{2} \text{tr} \left[\mathbf{B}^{-1} \mathbf{G} \text{diag} \left(\frac{\partial^3}{\partial \hat{l}_i^3} \log P(\mathbf{y}^j \mid \mathbf{l}) \Big|_{\mathbf{l}=\hat{\mathbf{l}}} \right) \right], \tag{5}$$

where $\frac{\partial \mathbf{B}}{\partial \hat{l}_i}$ is a matrix of element-wise partial derivatives.

Finally, $\frac{\partial \hat{l}_i}{\partial \theta}$ can be derived by using the fact that $\hat{\mathbf{l}} = \mathbf{G} \nabla (\log P(\mathbf{y}^j \mid \mathbf{l})) \Big|_{\mathbf{l}=\hat{\mathbf{l}}} + \mathbf{m}$, as explained in the main text. Therefore, for every parameter θ we have:

$$\frac{\partial \hat{\mathbf{l}}}{\partial \theta} = \frac{\partial \mathbf{G}}{\partial \theta} \nabla (\log P(\mathbf{y}^j \mid \mathbf{l})) \Big|_{\mathbf{l}=\hat{\mathbf{l}}} + \mathbf{G} \frac{\partial (\nabla (\log P(\mathbf{y}^j \mid \mathbf{l})) \Big|_{\mathbf{l}=\hat{\mathbf{l}}})}{\partial \hat{\mathbf{l}}} \frac{\partial \hat{\mathbf{l}}}{\partial \theta} + \frac{\partial \mathbf{m}}{\partial \theta}, \tag{6}$$

where the partial derivatives over matrices and vectors are computed element-wise. Writing these gradients explicitly for the model parameters, we have:

$$\begin{aligned}
\frac{\partial \hat{\mathbf{l}}}{\partial \sigma_v^2} &= \mathbf{K}_v \nabla (\log P(\mathbf{y}^j \mid \mathbf{l})) \Big|_{\mathbf{l}=\hat{\mathbf{l}}} + \mathbf{G} \nabla \nabla (\log P(\mathbf{y}^j \mid \mathbf{l})) \Big|_{\mathbf{l}=\hat{\mathbf{l}}} \frac{\partial \hat{\mathbf{l}}}{\partial \sigma_v^2} \\
\frac{\partial \hat{\mathbf{l}}}{\partial \sigma_e^2} &= \nabla (\log P(\mathbf{y}^j \mid \mathbf{l})) \Big|_{\mathbf{l}=\hat{\mathbf{l}}} + \mathbf{G} \nabla \nabla (\log P(\mathbf{y}^j \mid \mathbf{l})) \Big|_{\mathbf{l}=\hat{\mathbf{l}}} \frac{\partial \hat{\mathbf{l}}}{\partial \sigma_e^2} \\
\frac{\partial \hat{\mathbf{l}}}{\partial \alpha_k} &= \mathbf{G} \nabla \nabla (\log P(\mathbf{y}^j \mid \mathbf{l})) \Big|_{\mathbf{l}=\hat{\mathbf{l}}} \frac{\partial \hat{\mathbf{l}}}{\partial \alpha_k} + \mathbf{C}_k,
\end{aligned} \tag{7}$$

where \mathbf{C}_k is the k^{th} column of \mathbf{C} .

A careful analysis reveals that the only matrix that needs to be inverted is \mathbf{B} [3]. This is also the matrix whose determinant needs to be evaluated in the likelihood evaluation. Both quantities can be readily evaluated given the Cholesky decomposition of \mathbf{B} , whose computation scales cubically with the sample size. Hence, gradient computation incurs only a minor additional computational overhead over the likelihood approximation.

4 Data Simulations

The simulations procedure consisted of several stages: Genotype simulations, covariates simulations, cell-type simulations, phenotype simulation, methylation simulations and observed reads simulations. We describe each stage in turn.

4.1 Genome simulations

We simulated population structure via the Balding-Nichols model [4], wherein allele frequencies in the range [0.05, 0.5] were randomly drawn for an ancestral population, and frequencies for several subpopulations were drawn from a Beta distribution with parameters $f(1 - F_{ST})/F_{ST}$ and $(1 - f)(1 - F_{ST})/F_{ST}$, where f is the minor allele frequency in the ancestral population.

Next, we generated a mixture vector \mathbf{M}^i for every individual i from a symmetric Dirichlet distribution with a concentration parameter 0.25, such that $\sum_u \mathbf{M}_u^i = 1$, and u iterates over subpopulations. Finally, 60,000 SNPs were generated for every individual assuming Hardy-Weinberg equilibrium by sampling every SNP j of individual i from $\text{Bin}(2, p_j^i)$, where $p_j^i = \sum_u \mathbf{M}_u^i f_u^j$, and f_u^j is the MAF of SNP j in population u .

4.2 Covariate simulations

A vector of covariates for every individual i , \mathbf{C}^i , was generated by sampling each entry from a standard normal distribution, and adding an additional entry with the value 1.0 as an intercept.

4.3 Cell-type simulations

A vector of cell-types for every individual i , \mathbf{T}^i , was generated from a non-symmetric Dirichlet distribution with concentration parameters 1.33, 0.93, 4.299, 2.696, 16.344. These values were fitted using estimated cell type proportions for the GSE42861 dataset [5], a 450K array data set. Specifically, we obtained cell proportion estimates of five cell types (granulocytes, monocytes, B cells, NK cells, T cells) using the default implementation available in the minfi package [6], defined and assembled for the 450K array [7] based on the approach suggested by [8] and a 450K reference data set [9].

4.4 Phenotype simulation

To generate a phenotype, we first sampled effect sizes for the factors affecting the phenotype. Specifically, we sampled an effect size $\gamma_l^{\text{snp}} \sim \mathcal{N}(0, \sigma_{\text{snp}}^2)$ for every causal SNP l , an effect size $\gamma_u^{\text{pop}} \sim \mathcal{N}(0, \sigma_{\text{pop}}^2)$ for every subpopulation u (intended to model shared environmental factors within the subpopulations), and an effect size $\gamma_h^{\text{cell}} \sim \mathcal{N}(0, \sigma_{\text{cell}}^2)$ for every cell-type h .

Afterwards, the phenotype y^i for every individual i was generated as follows:

$$y^i = + \sum_{l \in \text{Causal}} s_l^i \gamma_l^{\text{SNP}} + \sum_j m_u^i \gamma_u^{\text{POP}} + \sum_h t_h^i \gamma_h^{\text{cell}} + \epsilon^i, \quad (8)$$

where s_l^i , m_u^i and t_h^i are the l^{th} causal SNP, u^{th} subpopulation ancestry fraction and h^{th} cell-type fraction of individual i , respectively, $\epsilon^i \sim \mathcal{N}(0, \sigma_\epsilon^2)$ is an environmental effect for individual i where σ_ϵ^2 guarantees that the phenotype variance is 1.0, and l iterates over indices of causal SNPs only. All effects were normalized to have a zero mean and unit variance.

4.5 Methylation simulation

Methylation levels π_j^i for every site j of every individual i were generated as follows: First, either 0%, 25% or 50% of the sites were randomly selected to be differentially methylated. Next, for every site j we generated effect sizes for the factors affecting the methylation. Specifically, we sampled an effect size $\delta_k^{j, \text{covar}} \sim \mathcal{N}(0, \nu_{\text{covar}}^{2, j})$ for every covariate k and site j , an effect size $\delta_l^{j, \text{SNP}} \sim \mathcal{N}(0, \nu_{\text{SNP}}^{2, j})$ for every causal SNP l and site j , an effect size $\delta_u^{j, \text{POP}} \sim \mathcal{N}(0, \nu_{\text{POP}}^{2, j})$ for every subpopulation u and site j (intended to model shared environmental factors within the subpopulations), an effect size $\delta_h^{j, \text{cell}} \sim \mathcal{N}(0, \nu_{\text{cell}}^{2, j})$ for every cell-type h and site j , and an effect size $\delta^{j, \text{pheno}} \sim \mathcal{N}(0, \nu_{\text{pheno}}^{2, j})$ for the phenotype. We also generated an environmental effect $e_j^i \sim \mathcal{N}(0, \nu_{\text{env}}^{2, j})$ for every site j of individual i .

Next, we generated $\text{logit}(\pi_j^{i, h})$ for every cell type h of site j of individual i as follows:

$$\text{logit}(\pi_j^{i, h}) = \sum_k c_k^i \delta_k^{j, \text{covar}} + \sum_{l \in \text{Causal}} s_l^i \delta_l^{j, \text{SNP}} + \sum_u m_u^i \delta_u^{j, \text{POP}} + \delta_h^{j, \text{cell}} + y^i \delta^{j, \text{pheno}} + e_j^i. \quad (9)$$

In the next step, the methylation level of every cell type h of site j of individual i , $\pi_j^{i, h}$, was computed as $\pi_j^{i, h} = 1 / (1 + \exp(-\text{logit}(\pi_j^{i, h})))$. Finally, a combined methylation level π_j^i for every individual i of every site j was computed via:

$$\pi_j^i = \sum_h t_h^i \pi_j^{i, h}. \quad (10)$$

4.6 Observed reads simulation

To simulate observed reads, we first sampled a total number of reads r_j^i for every site j of individual i from a negative Binomial distribution with parameters $n = 1.135515$, $p = 0.047623$. These values were fitted from the distribution of total number of reads of the Baboons data studied in the original MACAU paper [10].

Afterwards, an observed number of reads y_j^i for every site j of individual i was sampled from $\text{Bin}(r_j^i, \pi_j^i)$.

4.7 Default parameters

Unless otherwise stated, the default parameters used in the simulations are the following: Each data set consisted of 200 individuals, 3 covariates (including an intercept), 60,000 SNPs, 500 causal SNPs, 10,000 non-causal sites and 500 causal sites. The ancestry of every individual was a mixture of four different populations, sampled from a symmetric Dirichlet distribution with concentration parameter of 0.25. Additionally, every individual contained a mixture of 5 cell types, sampled from a Dirichlet distribution with parameters 1.33 , 0.93 , 4.299, 2.696, 16.344.

The default effect variances were the following: $\sigma_{\text{snp}}^2 = 0.25 / \#\text{SNPs}$, $\sigma_{\text{pop}}^2 = 0.025 / \#\text{populations}$, $\sigma_{\text{cell}}^2 = 0.25 / \#\text{cell-types}$, $\nu_{\text{covar}}^{2,j} = 0.01 / \#\text{covariates}$, $\nu_{\text{snp}}^{2,j} = 1 / \#\text{SNPs}$, $\nu_{\text{pop}}^{2,j} = 0.1 / \#\text{populations}$, $\nu_{\text{pheno}}^{2,j} = 0.025$, $\nu_{\text{env}}^{2,j} = 0.1$, $\nu_{\text{cell}}^{2,j} = 5.0$.

References

- [1] H. Chen et al. Control for Population Structure and Relatedness for Binary Traits in Genetic Association Studies via Logistic Mixed Models. *Am. J. Hum. Genet.* 98(4) (2016), 653–66.
- [2] D. Jiang, S. Zhong, and M.S. McPeck. Retrospective binary-trait association test elucidates genetic architecture of Crohn disease. *Am. J. Hum. Genet.* 98(2) (2016), 243–255.
- [3] C. Rasmussen and C. Williams. *Gaussian Processes for Machine Learning*. The MIT Press, 2006.
- [4] D.J. Balding and R.A. Nichols. A method for quantifying differentiation between populations at multi-allelic loci and its implications for investigating identity and paternity. *Genetica* 96(1-2) (1995), 3–12.
- [5] Y. Liu et al. Epigenome-wide association data implicate DNA methylation as an intermediary of genetic risk in rheumatoid arthritis. *Nat. Biotechnol.* 31(2) (2013), 142–147.
- [6] M.J. Aryee et al. Minfi: a flexible and comprehensive Bioconductor package for the analysis of Infinium DNA methylation microarrays. *Bioinformatics* 30(10) (2014), 1363–1369.
- [7] A.E. Jaffe and R.A. Irizarry. Accounting for cellular heterogeneity is critical in epigenome-wide association studies. *Genome Biol.* 15(2) (2014), 1.
- [8] E.A. Houseman et al. DNA methylation arrays as surrogate measures of cell mixture distribution. *BMC Bioinform.* 13(1) (2012), 1.
- [9] L.E. Reinius et al. Differential DNA methylation in purified human blood cells: implications for cell lineage and studies on disease susceptibility. *PLOS ONE* 7(7) (2012), e41361.
- [10] A.J. Lea, J. Tung, and X. Zhou. A Flexible, Efficient Binomial Mixed Model for Identifying Differential DNA Methylation in Bisulfite Sequencing Data. *PLoS Genet.* 11(11) (2015), e1005650.