1 **Detection and accurate False Discovery Rate control of differentially methylated regions**
2 **from Whole Genome Bisulfite Sequencing**
3

4 Keegan D. Korthauer[1,2], Sutirtha Chakraborty[3], Yuval Benjamini[4], and Rafael A. Irizarry[1,2]

5 [1]Department of Biostatistics & Computational Biology, Dana-Farber Cancer Institute
6 [2]Department of Biostatistics, Harvard T.H. Chan School of Public Health
7 [3]Novartis
8 [4]Department of Statistics, Hebrew University
9

10 **Summary**

11 With recent advances in sequencing technology, it is now feasible to measure DNA methylation

12 at tens of millions of sites across the entire genome. In most applications, biologists are

13 interested in detecting differentially methylated regions, composed of multiple sites with

14 differing methylation levels among populations. However, current computational approaches for

15 detecting such regions do not provide accurate statistical inference. A major challenge in

16 reporting uncertainty is that a genome-wide scan is involved in detecting these regions, which

17 needs to be accounted for. A further challenge is that sample sizes are limited due to the costs

18 associated with the technology. We have developed a new approach that overcomes these

19 challenges and assesses uncertainty for differentially methylated regions in a rigorous manner.

20 Region-level statistics are obtained by fitting a generalized least squares (GLS) regression model

21 with a nested autoregressive correlated error structure for the effect of interest on transformed

22 methylation proportions. We develop an inferential approach, based on a pooled null

23 distribution, that can be implemented even when as few as two samples per population are

24 available. Here we demonstrate the advantages of our method using both experimental data and

25 Monte Carlo simulation. We find that the new method improves the specificity and sensitivity of

26 list of regions and accurately controls the False Discovery Rate (FDR).

27

1   **Keywords:** bisulfite sequencing, differential methylation, false discovery rate, generalized least

2   squares, inference

3

4   **1. Introduction**

5   DNA methylation is an important epigenetic modification that plays a role in a wide variety of

6   biological processes. Numerous studies have been carried out to locate CpG loci where DNA

7   methylation may be involved in gene regulation, differentiation, and cancer. With recent

8   advances in sequencing technology such as Whole Genome Bisulfite Sequencing (WGBS), it is

9   now possible to measure DNA methylation at single base resolution across all CpGs in the

10  genome. Even though the most common application of the technology is to detect differentially

11  methylated regions (DMRs) between populations, most methods for analysis of WGBS

12  experiments focus on statistical differences for CpG loci one at a time (Akalin *et al.*, 2012,

13  Dolzhenko and Smith, 2014, Lee and Morris, 2016, Park *et al.*, 2014, Park and Wu, 2016). While

14  useful, approaches for identification of differentially methylated loci (DML) have many practical

15  limitations in both implementation and interpretation. Here, we discuss these limitations as well

16  as outline the challenges of performing inference at the region level. Finally, we introduce a

17  rigorous statistical approach that overcomes these challenges to construct de novo DMRs with

18  accurate FDR control.

19         Methods to identify DMLs in WGBS experiments are greatly hindered by the high-

20  dimensionality and low sample size setting that is common in high-throughput genomics studies.

21  The number of tests performed is equal to the number of loci analyzed, which is very large in

22  typical WGBS studies. In the human genome, for example, there are close to thirty million CpG

23  loci (Smith and Meissner, 2013). Further, DML methods generally do not account for the well-

24  known fact that measurements are spatially correlated across the genome (Leek *et al.*, 2010) and

1    instead treat measurements from all loci as independent. Correcting for multiple comparisons

2    without taking into account these correlations can result in a loss of power.

3        Additionally, methods for assessing the significance of DMLs typically require large

4    sample sizes due to reliance on large sample approximations (Dolzhenko and Smith, 2014,

5    Hansen, Langmead and Irizarry, 2012, Hebestreit, Dugas and Klein, 2013, Lee and Morris,

6    2016). Although WGBS is the current gold standard for estimating whole genome methylation

7    profiles (Marx, 2016), cost limitations are still a barrier to acquiring more than a few individuals

8    per biological condition in many studies (Ziller *et al.*, 2015). This is reflected in the study design

9    of major consortiums that aim to characterize the epigenome. For example, WGBS experiments

10   in murine embryos carried out as part of the ENCODE project are limited to two biological

11   replicates per tissue type and developmental time point combination (He *et al.*, 2017). In

12   addition, the number of biological replicates measured with WGBS in the UCSD Human

13   Reference Epigenome Mapping Project (Schultz *et al.*, 2015) is also limited to 2-3 per tissue

14   type. As such, we aim to maximize power while controlling the false discovery rate even with

15   sample sizes as small as two samples per condition.

16       Methods for identifying DMLs also need to properly model count data that does not

17   conform to standard Gaussian models. This is in contrast to methylation array analysis, where

18   Gaussian models performed well (Jaffe *et al.*, 2012). One option is to assume that methylation

19   proportions, defined as the number of methylated reads divided by the number of total reads

20   covering a given CpG locus, follow a normal distribution (Hansen, Langmead and Irizarry,

21   2012). However this assumption clearly does not hold when the total reads covering the CpG,

22   referred to as the coverage, is small, a common occurrence in these datasets. The approach also

23   ignores that variance of this proportion depends on the coverage. To overcome these limitations,

1    DML approaches have also modeled WGBS count data using Binomial models (Saito, Tsuji and

2    Mituyama, 2014). However, Binomial models on their own cannot account for biological

3    variability within sample groups. In order to account for biological variability in count data,

4    Beta-Binomial models (Park *et al.*, 2014, Sun *et al.*, 2014) are a natural extension. However they

5    come at the cost of increased computational burden when testing millions of loci.

6          Beyond implementation challenges, DML approaches also suffer from limited

7    interpretability. In general, identifying DMRs is more biologically relevant than reporting DMLs.

8    Apart from the so-called 'CpG traffic lights' (Khamis *et al.*, 2017), most individual CpG loci

9    likely do not have a large impact on epigenetic function on their own, but rather through a

10   biochemical modification that involves several loci. Most notably, regional DNA methylation

11   levels are correlated with the expression levels of nearby genes. Specifically, methylation gain is

12   associated with stable transcriptional silencing of nearby genes (Bird, 2002). In the context of

13   differential methylation analysis, Aryee *et al.* (2014) found that differentially expressed genes

14   were consistently more likely to be located near DMRs than DMLs.

15          While DML approaches may construct DMRs by chaining together neighboring

16   significant loci, this type of approach will not yield a proper assessment of the statistical

17   significance of the constructed regions, nor will the False Discovery Rate (FDR) be properly

18   controlled (Robinson *et al.*, 2014). This is because controlling the FDR at the level of individual

19   loci is not the same as controlling FDR of regions, as has been noted in the context of peak

20   calling in ChIP-seq experiments (Lun and Smyth, 2014, Siegmund, Zhang and Yakir, 2011).

21   FDR correction at the level of individual loci means that the proportion of expected false positive

22   loci is controlled, not the proportion of false positive regions. Statistically, this is a critical point

23   since FDR control of DMR detection is not guaranteed under the DML setting. In fact, many

1    discoveries at the loci level may constitute only a single discovery. This means that a large

2    number of correct rejections at the loci level can inflate the denominator in the FDR calculation,

3    which will artificially lower the false discovery rate of loci as compared to regions (Figure 1).

4    We were motivated to develop a procedure to control FDR at the region level and provide an

5    accurate measure of statistical significance for each region.

6        Many recent computational approaches have been developed with the goal of identifying

7    DMRs, but most do not provide formal inference for regions (Hansen, Langmead and Irizarry,

8    2012, Saito, Tsuji and Mituyama, 2014, Wu *et al.*, 2015, Yu and Sun, 2016) and instead join

9    together significant DMLs. This type of procedure will suffer from the problems outlined above.

10   Other approaches can perform inference at the region level, but only for predefined regions of

11   interest or fixed sliding windows (Hebestreit, Dugas and Klein, 2013, Sun *et al.*, 2014). Though

12   useful in targeted settings such as Reduced Representation Bisulfite Sequencing (RRBS), or

13   when we have prior knowledge of the DMR size, they are not applicable to identifying DMRs of

14   arbitrary size from WGBS. Those methods that scan the genome for DMRs and provide

15   inference at the region level do not properly control FDR (Juhling *et al.*, 2016, Wen *et al.*, 2016).

16   This is evidenced, for example, by the FDRs reported in the simulation studies of Wen *et al.*

17   (2016), which were as high as 0.85 and widely varied across scenarios. Juhling *et al.* (2016) also

18   do not achieve accurate FDR control in simulation studies (see Section 4.1).

19       The challenge of performing inference at the region level is complicated by several

20   factors in addition to the challenges already discussed in the context of DML analysis. The first

21   challenge is in defining the region boundaries themselves. Without prior knowledge or

22   predefined regions, we need to construct data-driven regions. Calculating a test statistic for these

23   data-driven regions of varying sizes with a known null distribution is not straightforward. In

1   addition, challenges are presented by the complex statistical dependencies observed in

2   measurements from nearby loci (Benjamini, Taylor and Irizarry, 2016), as well as different

3   within group variability across loci (Hansen, Langmead and Irizarry, 2012). Some methods

4   ignore correlation across loci (Wen *et al.*, 2016) or biological variability from sample to sample

5   (Saito, Tsuji and Mituyama, 2014, Wu *et al.*, 2015). Not properly accounting for both of these

6   sources of variability in DNA methylation data, however, results in misleading conclusions or

7   loss of power. For a full review of DML and DMR methods, see Shafi *et al.* (2017).

8       To overcome the limitations and challenges detailed above, we propose a two-stage

9   approach that first detects candidate regions and then explicitly evaluates statistical significance

10   at the region level while accounting for known sources of variability. Candidate DMRs are

11   defined by segmenting the genome into groups of CpGs that show consistent evidence of

12   differential methylation. Because the methylation levels of neighboring CpGs are highly

13   correlated, we first smooth the signal to combat loss of power due to low coverage as done by

14   Hansen, Langmead and Irizarry (2012). In the second stage, we compute a statistic for each

15   candidate DMR that takes into account variability between biological replicates and spatial

16   correlation among neighboring loci. Significance of each region is assessed via a permutation

17   procedure which uses a pooled null distribution that can be generated from as few as two

18   biological replicates, and false discovery rate is controlled using the procedure of Benjamini and

19   Hochberg (1995). Code to reproduce the analyses presented in this paper is provided in

20   Supplementary material and the open-source R package dmrseq that implements the approach is

21   available on GitHub.

22       In Section 2, we provide a detailed description of the datasets used. We describe the

23   methodological details of the approach and detail the data processing and analysis procedure in

1    Section 3. In Section 4, we present our findings using both experimental data and simulations.

2    We demonstrate that the proposed approach assigns greater statistical significance to regions that

3    have greater biological significance in terms of potential functional roles in the regulation of

4    gene expression. We also evaluate sensitivity and specificity of the approach by analyzing null

5    comparisons of samples from the same biological condition, with and without adding simulated

6    DMRs. We demonstrate that dmrseq has higher sensitivity than existing approaches and

7    accurately assesses statistical significance of regions through False Discovery Rate estimation. A

8    discussion of the advantages and limitations of the method are given in Section 5.

9

10   **2. Data Description**

11   dmrseq is generally applicable to WGBS data which contains the counts for both methylated and

12   unmethylated reads mapping to each CpG loci. This information can be obtained from raw

13   sequencing reads using the mapping software Bismark (Krueger and Andrews, 2011), as

14   described in the Supplementary materials. Specifically, CpG loci that are covered by at least one

15   read in every sample should be used in the analysis. Other methods for analysis of WGBS data

16   recommend removing CpG sites that have only a few reads in each sample, and while processed

17   data of this form may be analyzed by our approach, it is important to note that this may result in

18   a loss of power to detect regions in low-coverage areas of the genome.

19          In this study, we use our approach to identify DMRs using publically available WGBS

20   data from two different case studies, as described below. We also evaluate sensitivity and

21   specificity of DMR methods by applying them to simulated data. Summary of coverage and

22   methylation values for all datasets used can be found in Table 1 and Supplementary Figure S2.

23   For more details on data processing, see Section 1 of the Supplementary materials.

1   *2.1 Simulated data*

2   Two sets of simulated data were constructed: one representing a null comparison (with no

3   DMRs) and another containing simulated DMRs. To ensure that the simulated datasets closely

4   match the characteristics of the observed experimental data, they were generated based on

5   WGBS data from a study of human dendritic cells (Pacis *et al.*, 2015). This study estimated

6   methylation profiles of human dendritic cells from six donors before and after infection with a

7   pathogen. The null comparison was constructed by randomly partitioning the six control samples

8   (before infection) into two groups of three samples each, denoted Simulation N3. The same is

9   done for a subset of four of the samples to evaluate performance when there are only two

10  samples in each population, denoted Simulation N2.

11       Starting with the null comparisons, 3,000 simulated DMRs were added to each dataset in

12  order to evaluate specificity and sensitivity. These are denoted Simulations D2 and D3 for two

13  and three samples per population, respectively. Briefly, a DMR is constructed by sampling a

14  cluster of neighboring CpGs and simulating the number of methylated reads, conditional on

15  observed coverage, for the samples from one population from a binomial distribution. The

16  binomial probabilities are equal to the observed methylation proportions plus or minus a

17  randomly sampled difference, which varies smoothly over the region according to a function

18  similar to the tricube kernel (Cleveland, 1979) (see Section 2.4 of the Supplementary materials).

19

20  *2.2 UCSD Human Reference Epigenome Mapping Project*

21  Data from several human tissue samples from the UCSD Human Reference Epigenome Mapping

22  Project (Schultz *et al.*, 2015) was used to identify DMRs related to tissue type.  Specifically, four

1    tissues were selected for performing pairwise comparisons: (1) Heart, left ventricle, (2) Heart,

2    right ventricle, (3) Sigmoid colon, and (4) Small intestine.

3

4    *2.3 Murine models of leukemia*

5    In this study, marrow or thymus cells from two biological replicates form each of three different

6    murine lines were extracted and genome-wide methylation levels measured with WGBS. One

7    condition consisted of a wild-type control mouse. The other two had alterations in one or both of

8    the DNMT3a or FLT3 loci, both of which have previously demonstrated implications in the

9    development of leukemia (Pacis *et al.*, 2015). The mouse model with a wild-type DNMT3a locus

10   and a duplication of the FLT3 locus has been shown to induce ALL. The mouse model with the

11   same duplication of the FLT3 locus as well as a knock out of DNMT3a has been shown to

12   induce the more lethal and aggressive AML. The DNMT3a also plays a role in promoting DNA

13   methylation, so it is of interest to characterize the resulting differences in methylated regions

14   among the control and two different leukemia models.

15

16   **3. Analysis Framework**

17   A two-step procedure is carried out to (1) construct de novo candidate regions, and (2) score

18   candidate regions to quantify the effect of the covariate of interest on methylation level, and

19   evaluate statistical significance by comparing them to null regions. Here we detail each stage of

20   the approach.

21

22   *3.1 Construction of candidate regions*

1 In step 1, we detect candidate regions that contain multiple loci showing evidence of a difference

2 in the smoothed pooled methylation proportion between biological conditions. For simplicity of

3 presentation, we assume there are two biological conditions $s \in \{1,2\}$, with sample indices

4 $j \in C_s$ (see Supplementary materials Section 2.7 for the case of more than two conditions). Let

5 $M_{ij}$ be the number of methylated reads and $U_{ij}$ the number of unmethylated reads for locus $i$ of

6 sample $j$ from condition $s$. The coverage is denoted $N_{ij}$, where $N_{ij} = M_{ij} + U_{ij}$. The estimate of

7 the mean methylation proportion $\hat{\pi}_{is}$ for loci $i$ in condition $s$ is taken to be the sum of methylated

8 reads from all samples in that condition divided by the sum of all reads (i.e. the coverage) from

9 all samples in condition $s$:

$$\hat{\pi}_{is} = \sum_{j \in C_s} M_{ij} \Big/ \sum_{j \in C_s} N_{ij}$$

10 This leads to the following estimate of methylation proportion difference $\hat{\beta}_i$ between condition $s$

11 and $s'$ at loci $i$:

$$\hat{\beta}_i = \hat{\pi}_{is} - \hat{\pi}_{is'}$$

12 In order to give more weight to measurements with higher coverage, this estimate pools together

13 samples within the same condition. To account for biological variability between samples and

14 further reduce influence of observations with low coverage, smoothed individual loci estimates

15 $\hat{\beta}_i^{Smooth}$ are obtained using a local-likelihood smoother (Loader, 1999) with smoothing weights

16 $w_i$ equal to the median coverage at loci $i$ scaled by the average Median Absolute Deviation

17 (MAD) within the sample groups $\bar{\delta}_i$:

18 $$w_i = \frac{median_j(N_{ij})}{\bar{\delta}_i}, \text{ where } \bar{\delta}_i = \frac{1}{2}\sum_s Median_{j \in C_s} \left| \frac{M_{ij}}{N_{ij}} - Median_{k \in C_s}\left(\frac{M_{ik}}{N_{ik}}\right) \right|$$

19

20 This places more emphasis on observations with high coverage and low variability within sample

21 group (see Section 2.1 of the Supplementary material for more details).

1          Candidate regions are defined by segmenting the genome into groups of loci with a

2      smoothed and scaled pooled proportion difference $\hat{\beta}_i^{Smooth} / \hat{\sigma}_i(\hat{\beta}_i)$ in the same direction that is

3      greater than some threshold in absolute value (refer to Supplementary materials Section 2.2 for

4      more details). Maximum spacing between loci within a candidate region is controlled by a

5      predetermined value, and loci at the start and end of the region with low difference values are

6      trimmed (refer to Supplementary materials Section 2.3 for more details). The threshold value

7      should be chosen liberally so that it will more or less capture all of the true differences without

8      regard to false positives, as significance of the candidate regions is assessed in the next step.

9

10     *3.2 Assessing significance of regions*

11     In the second step, we assess the significance of candidate regions. This task is complicated by

12     the fact that the null statistics are calculated on an enriched set of regions. In general, the null

13     distribution generated by the type of selection procedure described in the previous section is not

14     known. A natural approach would be to carry out a permutation test to control FWER (family-

15     wise error rate), which is done by Jaffe *et al.* (2012) to infer DMRs from array data. However,

16     this is not feasible when we have only a few samples per population as is most often the case

17     with WGBS. Thus, we set out to construct a statistic that can be comparable across the genome

18     so that the signal can be compared among regions. Such an exchangeable statistic allows us to

19     generate an approximate null distribution by pooling genomewide candidate regions detected

20     from permutations.

21          To generate an approximately exchangeable region statistic that measures the strength of

22     methylation difference, we need to account for sources of variation that are known to vary across

23     the genome, including biological variability from sample to sample (Hansen, Langmead and

1    Irizarry, 2012), as well as covariance of nearby loci (Benjamini, Taylor and Irizarry, 2016).

2    Failing to do so may result in large test statistics just by chance for regions with high variability,

3    leading to increased FDR or decreased power. For example, if we use an area-based statistic

4    (Hansen, Langmead and Irizarry, 2012) or a mean difference statistic averaged across loci, power

5    to detect DMRs is greatly reduced in simulation studies (Supplementary Figure S5 and

6    Supplementary materials Section 4.1).

7         Since we need to compute the statistic over potentially hundreds of thousands of

8    candidate regions, we also favor an approach that provides efficient and stable estimation

9    procedures. For these reasons, we make use of generalized least squares (GLS) regression model

10   with a nested autoregressive correlated error structure for the effect of interest on transformed

11   methylation proportions, the advantages of which are described in detail in the next subsections.

12

13   *3.2.1 Estimation of region statistics with Generalized Least Squares models*

14   To account for sampling variability, we assume that methylation counts for region $r$ are

15   Binomially-distributed with probability $p_{ijr}$, where

16   $$M_{ijr} \mid N_{ijr}, p_{ijr} \sim Bin(N_{ijr}, p_{ijr}).$$

17   To model biological variability, we allow the binomial proportion for samples in condition

18   $s \in \{1,2\}$ to vary according to a beta distribution with shape parameters $\alpha_{irs}$ and $\beta_{irs}$, where

19   $$p_{ijr} \sim Beta(\alpha_{irs}, \beta_{irs}).$$

20   Let $\pi_{irs} = \frac{\alpha_{irs}}{\alpha_{irs} + \beta_{irs}}$ denote the mean of this Beta distribution. We are interested in estimating and

21   assessing the significance of the difference in mean methylation levels across a region $r$ for two

22   biological conditions.

1      Our approach models transformed methylation proportions using GLS to obtain an

2      approximation of the effect of interest. While directly modeling counts with either a Beta-

3      Binomial Generalized Linear Model (GLM) or a Generalized Linear Mixed Model (GLMM)

4      would allow us to accommodate complex covariance structures across samples and loci, it also

5      results in complex likelihoods that require iterative maximization for each candidate region.

6      Further, these procedures are subject to instability of estimation for methylation levels near the

7      boundaries (zero and one) or non-identifiability in the case of separation as they occur in GLM

8      (Gelman *et al.*, 2008) and GLMM (Abrahantes and Aerts, 2012) estimation. GLS models, in

9      contrast, are efficient and stable to estimate due to the availability of approximate closed-form

10      parameter estimates. Though GLS does not model counts directly, we incorporate information

11      lost after transformation of methylation proportions through specification of a variance estimate

12      that depends on coverage.

13      We choose the arcsine link function $Z_{ijr} = arcsin\left( 2M_{ijr}/N_{ijr} - 1\right)$ to obtain

14      transformed methylation proportions, as proposed by (Park and Wu, 2016) for DML analysis, for

15      its desirable ability to stabilize the dependence of the variance on the mean methylation level.

16      While the variance of methylation proportions $M_{ijr}/N_{ijr}$ depends on the mean parameter $\pi_{ijr}$,

17      the variance of $Z_{ijr}$ only depends on coverage $N_{ijr}$ and the dispersion of the Beta-Binomial

18      distribution (refer to Supplementary materials Section 2.6 for more details). This helps us to form

19      a statistic involving the transformed proportions that is exchangeable across regions that have

20      different mean methylation values.

21      We assume a linear effect on the arcsine link-transformed methylation proportion

22      parameters:

$$arcsin(2\pi_{ijr} - 1) = \sum_{l=1}^{L_r} \beta_{0lr} 1_{[i=l]} + \beta_{1r}X_j = X\boldsymbol{\beta}_r$$

1    Here $\beta_{0lr}$ are loci-specific intercept terms that account for variation on overall methylation levels

2    across the region, where $l = 1, \dots L_r$ and $L_r$ denotes the number of loci in region $r$. The

3    coefficient for the effect of interest (e.g. biological group) is $\beta_{1r}$. We denote the design matrix as

4    $X$ and the $(L_r + 1)$-length vector of all coefficients $(\beta_{01r}, \beta_{02r}, \dots, \beta_{0L_r r}, \beta_{1r})$ as $\boldsymbol{\beta}_r$. This leads

5    to the following model for the transformed response $Z_r = (Z_{11r}, \dots, Z_{L_r Jr})$ in region $r$

$$Z_r = X\boldsymbol{\beta}_r + \boldsymbol{\varepsilon}_r$$

6    where we assume that $E[\boldsymbol{\varepsilon}_r] = 0$ and $Var[\boldsymbol{\varepsilon}_r] = V_r$, which can be fit by GLS given an estimate

7    of the covariance matrix $V_r$. Since GLS allows arbitrary covariance structures, we use an

8    autoregressive correlation structure to account for the correlation of methylation levels among

9    nearby loci. To account for the dependence of the variance on coverage as mentioned above, we

10    use variance weights. More details on the specific structure and estimation of $V_r$ are given in the

11    next section.

12            With the above model, we assess the strength of the effect of the covariate of interest on

13    methylation level within region $r$ using the t-statistic $t_r$ from the Wald test of the null hypothesis

14    that $\beta_{1r} = 0$. Parameter estimates and their standard errors are obtained with the `gls` function in

15    the `nlme` package (Pinheiro *et al.*, 2017). Significance is evaluated by permutation using a

16    pooled null distribution as described in detail in Section 3.2.3.

17

18    *3.2.2 Covariance of methylation levels within regions*

19    In the estimation of the covariance matrix $V_r$, we take into account biological variability through

20    variance weighting, and correlation of nearby loci through an autocorrelation structure. The

21    variance weighting is done to account for the dependence of the variance of transformed values

1   $Z_{ijr}$ on coverage. This variance depends non-linearly on $N_{ijr}$ (Supplement Section 2.6), but in

2   order to enable efficient closed-form estimation with GLS, we further approximate it by

$$Var(Z_{ijr}) \approx \frac{\sigma_r^2}{N_{ijr}}$$

3   In addition, in order to construct a valid permutation test where the variance conditional on the

4   effect of interest is invariant to permutation, we assume this variance identical for all samples at

5   a given loci by approximating $N_{ijr}$ by $median_j(N_{ijr}) = N_{i.r}$.

6           To model correlation of nearby loci, we use the flexible continuous autoregressive

7   correlation structure of order 1, abbreviated CAR(1). Under CAR(1), the correlation parameter

8   depends on the length of the interval between the two observations considered in the following

9   manner

$$\rho_r(\tau) = e^{-\phi_r|\tau|}$$

10  where $\tau$ is the length of the interval between two observations and $\phi_r$ is the positive continuous-

11  time autoregressive coefficient (following the notation of Jones and Boadi-Boateng (1991)) for

12  region $r$. Thus, for subject $j$, the predicted methylation value for loci $i$ at location $t_{ijr}$ in region $r$

13  given the methylation value at loci $i - 1$ is

$$\hat{Z}_{ijr} = \hat{Z}_{i-1,jr} \, e^{-\phi_r|t_{ijr}-t_{i-1,jr}|}$$

14  If the error variance of the CAR1 process is $\sigma_{ir}^2 = \frac{\sigma_r^2}{N_{i.r}}$, and we let the correlation structure be

15  nested within subject (i.e. such that observations from two subjects are independent), it follows

16  that the covariance matrix for a given sample can be written

$$\boldsymbol{V}_{jr} = \sigma_r^2 \begin{pmatrix} \dfrac{1}{N_{1.r}} & \dfrac{e^{-\phi_r|t_{1jr}-t_{2jr}|}}{\sqrt{N_{1.r}N_{2.r}}} & \cdots & \dfrac{e^{-\phi_r|t_{1jr}-t_{L_rjr}|}}{\sqrt{N_{1.r}N_{L_r.r}}} \\[2ex] \dfrac{e^{-\phi_r|t_{2jr}-t_{1jr}|}}{\sqrt{N_{1.r}N_{2.r}}} & \dfrac{1}{N_{2.r}} & \cdots & \dfrac{e^{-\phi_r|t_{2jr}-t_{L_rjr}|}}{\sqrt{N_{2.r}N_{L_r.r}}} \\[2ex] \vdots & \vdots & \ddots & \vdots \\[2ex] \dfrac{e^{-\phi_r|t_{L_rjr}-t_{1jr}|}}{\sqrt{N_{1.r}N_{L_r.r}}} & \dfrac{e^{-\phi_r|t_{L_rjr}-t_{2jr}|}}{\sqrt{N_{2.r}N_{L_r.r}}} & \cdots & \dfrac{1}{N_{L_r.r}} \end{pmatrix}$$

1    and for two subjects $j$ and $j'$, $Cov(Z_{ijr}, Z_{ij'r}) = 0$.

2        The estimation of $\phi_r$ is computationally efficient to carry out on small to moderately sized

3    regions. However, for larger regions with more than 40 loci we use the slightly simpler AR(1)

4    correlation structure since it is many times faster to compute. This discrete formulation assumes

5    that observations are equally spaced, and that observations that are separated by lag 1 are

6    correlated with region-specific correlation parameter $\rho_r$. In addition, observations that are

7    separated by $m$ positions are correlated by $\rho_r^m$. This results in a covariance matrix for $\boldsymbol{Z}_r$ from

8    region $r$, subject $j$ of

9    $$\boldsymbol{V}_{jr} = \sigma_r^2 \begin{pmatrix} \dfrac{1}{N_{1.r}} & \dfrac{\rho_r}{\sqrt{N_{1.r}N_{2.r}}} & \cdots & \dfrac{\rho_r^{L_r-1}}{\sqrt{N_{1.r}N_{L_r.r}}} \\[2ex] \dfrac{\rho_r}{\sqrt{N_{1.r}N_{2.r}}} & \dfrac{1}{N_{2.r}} & \cdots & \dfrac{\rho_r^{L_r-2}}{\sqrt{N_{2.r}N_{L_r.r}}} \\[2ex] \vdots & \vdots & \ddots & \vdots \\[2ex] \dfrac{\rho_r^{L_r-1}}{\sqrt{N_{1.r}N_{L_r.r}}} & \dfrac{\rho_r^{L_r-2}}{\sqrt{N_{2.r}N_{L_r.r}}} & \cdots & \dfrac{1}{N_{L_r.r}} \end{pmatrix}$$

10    and again we assume that for two subjects $j$ and $j'$, $Cov(Z_{ijr}, Z_{ij'r}) = 0$.

11        The CAR(1) structure simplifies to the AR(1) process under certain conditions when

12    observations are equally spaced (Jones and Boadi-Boateng, 1991). Thus the discrete AR(1) can

13    be viewed as an approximation of the CAR(1) when correlations are positive and the two provide

14    increasingly more similar estimates as observations approach constant spacing. Indeed, when

15    comparing model fits under both correlation structures in simulated data, the t-statistics for the

1    coefficient of interest under CAR(1) generally converge to the estimates under AR(1) as the

2    number of loci increases (Supplementary Figure S1 and Section 2.5).

3

4    *3.2.3 Permutation to generate a null set of regions*

5    The values of the covariate of interest (e.g. biological group) are permuted and the previous steps

6    repeated in order to generate a set of statistics under the null hypothesis. Since the statistics

7    account for known sources of variation that would otherwise prevent to comparison of regions

8    across the genome, we can pool them together to form an approximate null distribution with as

9    few as two samples per population. The empirical p-value is calculated by comparing the

10   observed test statistics to the entire null set of statistics from all permutations. Control of FDR is

11   carried out by adjusting the p-values using the procedure of Benjamini and Hochberg (1995).

12

13   **4. Results**

14   For each of the datasets described in Section 2, we applied dmrseq, as well as three widely used

15   methods for DMR detection: BSmooth (Hansen, Langmead and Irizarry, 2012), DSS (Park and

16   Wu, 2016), and metilene (Juhling *et al.*, 2016). Each approach was evaluated based on the

17   criteria detailed in the next subsections. For specific details on software implementation, refer to

18   the Supplementary materials (Section 3).

19

20   **4.1 Simulation using dendritic cell data**

21   Specificity was evaluated by identifying DMRs in null comparisons of two (N2) and three (N3)

22   samples per group. Sensitivity was evaluated by identifying simulated DMRs in comparisons of

23   two (D2) and three (D3) samples per group. Performance of each method is assessed by its

1    ability to identify as many of the simulated DMRs as possible, while identifying as few DMRs as

2    possible in the null comparison.

3        dmrseq did not identify any DMRs at the 0.05 level for the null comparisons N2 or N3

4    (Table 2). This remains true even when increasing the FDR threshold to 0.5 in both settings. In

5    contrast, metilene identified a small number of DMRs, DSS identified many hundreds, and

6    BSmooth tens of thousands using default settings (specific parameter specifications provided in

7    Supplementary materials Section 2.6). When applied to the datasets with simulated DMRs (D2

8    and D3), dmrseq is able to accurately control the False Discovery Rate, whereas metilene cannot

9    (Figure 2, Supplementary Figure S3). Note that analogous results cannot be obtained from DSS

10   or BSmooth, as there is no way to specify FDR level.

11       BSmooth and DSS identify similar numbers of False Positive regions in D2 and D3

12   compared to the null setting of N2 and N3, and far more than dmrseq and metilene (Table 3).

13   Although both BSmooth and DSS have favorable numbers of TPs, it is clear that this comes at

14   the expense of lack of control of FDR (Figure 3). Similarly, metilene has favorable numbers of

15   FPs, but this comes at the expense of low power. Further, even at similar observed FDR levels,

16   dmrseq achieves higher power levels than the alternative methods.

17       Although FDR thresholds are not available for BSmooth or DSS, we also investigated the

18   sensitivity and specificity of other settings beyond defaults of the thresholds at the single-loci

19   level (the loci t-statistic cutoff for BSmooth, and the loci p-value for DSS). Making these

20   thresholds more conservative generally reduced the numbers of False Positives, but once again

21   dmrseq was consistently able to identify more True Positives at similar numbers of False

22   Positives (See Supplementary Results and Figure S4).

1    We also stress that although lower False Positive rates could be achieved in this

2    simulation study for BSmooth and DSS, individual loci thresholds do not correspond directly to

3    specific FDRs at the region level.  As a result, in practice, one must choose a threshold either by

4    default settings, or by trial and error.

5

6    **4.2 Human tissue and murine leukemia experimental data**

7    The human tissue and murine leukemia studies were evaluated empirically based on the observed

8    association of DMRs with differential expression by RNA-seq. Differentially expressed (DE)

9    genes were identified using DESeq2 version 1.14.1 (Love, Huber and Anders, 2014). To assess

10   functional relevance of the results, detected DMRs that overlap promoter regions of DE genes

11   were assessed for signal in the expected direction. Specifically, a DMR - DE gene pair is

12   expected to have higher methylation values in the sample group with lower expression. The odds

13   that the DMR and DE statistics are in opposing directions are calculated at various FDR cutoffs

14   for dmrseq and metilene to assess whether top-ranked DMRs are more likely to be biologically

15   relevant. The same is done for various cutoffs for the numbers of top-ranking regions by effect

16   size. Additionally, for each cutoff we calculate the number of CpGs covered and the proportion

17   of detected DMRs that are within 2kb (from the center of the region) of a promoter region of a

18   DE gene.

19   To qualitatively assess the ability of the dmrseq region-level summary statistic to rank

20   DMRs as compared to other methods, we display example regions from the human tissue and

21   murine leukemia studies. These examples illustrate the increased variability of regions that are

22   highly ranked by naïve statistics but not dmrseq (Figure 5). We include a DMR with concordant

23   rankings that exhibits clear differences between two human tissue types (Figure 5A). In contrast,

1    the regions with discordant rankings between dmrseq q-value and mean difference (Figure 5B)

2    and area statistics (Figure 5C) exhibit considerable variability between samples or loci (See

3    Supplementary materials Section 2.8 for more details).

4

5    *4.2.1 Tissue specificity in human samples*

6    For DSS, metiline, and dmrseq, the number of DMRs found (Table 4) parallels the numbers of

7    DE genes found by DESeq2 (Supplementary Table S2), but DSS generally found far more

8    DMRs and metline far fewer.  For BSmooth, however, the number of DMRs identified was

9    similar for all comparisons. This happens because the cutoff for the individual loci statistics is set

10    by default at a quantile of the observed statistics, resulting in a similar number of loci being

11    deemed significant.

12          The tissue-specific DMRs found by dmrseq are enriched for inverse associations with DE

13    genes, and this enrichment is stronger for DMRs with lower FDRs (Figure 4). Additionally,

14    enrichment of dmrseq DMRs is generally stronger than that of alternative methods. While

15    metiline also provides an FDR estimate, there is no consistent association between the FDR

16    ranking and strength of association with expression. DMRs identified by BSmooth and DSS

17    cannot be ranked by FDR and the default settings may not be ideal, so we also rank DMRs by

18    effect size (raw methylation difference) with optimized parameter settings (see Supplementary

19    materials Section 3.2). The BSmooth and DSS DMRs with highest effect sizes exhibit

20    comparable enrichment to dmrseq, with metilene considerably lower (Supplementary Figures S6

21    and S7). However, arbitrary cutoffs of effect size do not directly correspond to significance level,

22    and the enrichment when including all DMRs is highest for dmrseq (Figure 4).

23

1    *4.2.2 DNMT3a loss in murine leukemia models*

2    In the murine leukemia models, dmrseq finds the most DMRs in the comparison of AML and the

3    control (Table 5), which is also the comparison for which the most DE genes were identified (see

4    Supplementary Table S4). In contrast, DSS and metline both find the most DMRs in the

5    comparison with the fewest DE genes identified, and BSmooth identified similar numbers of

6    DMRs in each comparison, each with far more DMRs than the other methods.

7        The murine leukemia DMRs found by dmrseq are enriched for inverse associations with

8    DE genes, and this enrichment is stronger for DMRs with lower FDRs (Figure 4). Additionally,

9    enrichment is generally stronger than that of BSmooth, DSS, and metline. While metiline also

10    provides an FDR estimate, there is no consistent association between the FDR ranking and

11    strength of association with expression. Similar to the tissue specificity analysis, BSmooth and

12    DSS DMRs with highest effect sizes exhibit comparable enrichment to dmrseq, with metilene

13    considerably lower, and the enrichment when including all DMRs often drops lower for

14    BSmooth, DSS, or metline than for dmrseq (Supplementary Figures S8 and S9).

15

16    **5. Discussion**

17    We have described dmrseq, a method useful for discovering and prioritizing DMRs from WGBS

18    data. The approach is based on rigorous statistical reasoning and is the first method that permits

19    accurate inference on DMRs that are found by scanning the genome. By developing a

20    transformation that results in summary statistics from candidate regions being exchangeable, we

21    are able to borrow strength across the genome to build a null distribution that permits inference

22    with a sample size as small as 2. We have demonstrated how the method clearly outperforms

1   currently used tools with several experimental data examples and Monte Carlo simulation. The

2   method is implemented as open source software in the form of an R package.

3

4   **Supplementary Material**

5   The reader is referred to the online Supplementary Materials for further details of data

6   acquisition and processing, additional methodological details, software implementation details,

7   and supplementary results. In addition, annotated R scripts for the simulation and case study

8   analyses are available in the GitHub repository https://github.com/kdkorthauer/dmrseqPaper, and

9   the R package dmrseq is available on GitHub at https://github.com/kdkorthauer/dmrseq.

10

11  **Acknowledgement**

12  *Conflict of Interest:* None declared

13

14  **Funding**

16

**References**

17  ABRAHANTES, J. C. and AERTS, M. (2012). A solution to separation for clustered binary data.
18          *Statistical Modelling* **12**, 3-27.

19  AKALIN, A., KORMAKSSON, M., LI, S., GARRETT-BAKELMAN, F. E., FIGUEROA, M. E., MELNICK,
20          A. and MASON, C. E. (2012). methylKit: a comprehensive R package for the analysis of
21          genome-wide DNA methylation profiles. *Genome Biol* **13**, R87.

22  ARYEE, M. J., JAFFE, A. E., CORRADA-BRAVO, H., LADD-ACOSTA, C., FEINBERG, A. P., HANSEN,
23          K. D. and IRIZARRY, R. A. (2014). Minfi: a flexible and comprehensive Bioconductor
24          package for the analysis of Infinium DNA methylation microarrays. *Bioinformatics* **30**,
25          1363-1369.

BENJAMINI, Y. and HOCHBERG, Y. (1995). Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society Series B* **57**, 289-300.

BENJAMINI, Y. and HOCHBERG, Y. (1995). Controlling the False Discovery Rate: a Practical and Powerful Approach to Multiple Testing. *Journal of the Royal Statistical Society B* **57**, 289-300.

BENJAMINI, Y., TAYLOR, J. and IRIZARRY, R. A. (2016). Selection Corrected Statistical Inference for Region Detection with High-througput Assays. *bioRxiv*.

BIRD, A. (2002). DNA methylation patterns and epigenetic memory. *Genes & Development* **16**, 6-21.

CLEVELAND, W. S. (1979). Robust locally weighted regression and smoothing scatterplots. *Journal of the American Statistical Association* **74**, 829-836.

DOLZHENKO, E. and SMITH, A. D. (2014). Using beta-binomial regression for high-precision differential methylation analysis in multifactor whole genome bisulfite sequencing experiments. *BMC Bioinformatics* **15**, 215.

DOLZHENKO, E. and SMITH, A. D. (2014). Using beta-binomial regression for high-precision differential methylation analysis in multifactor whole-genome bisulfite sequencing experiments. *BMC Bioinformatics* **15**, 215.

GELMAN, A., JAKULIN, A., PITTAU, M. G. and SU, Y.-S. (2008). A weakly informative default prior distribution for logistic and other regression models. *The Annals of Applied Statistics* **2**, 1360-1383.

HANSEN, K. D., LANGMEAD, B. and IRIZARRY, R. A. (2012). BSmooth: from whole genome bisulfite sequencing reads to differentially methylated regions. *Genome Biol* **13**, R83.

HE, Y., HARIHARAN, M., GORKIN, D. U., DICKEL, D. E., LUO, C., CASTANON, R. G., NERY, J. R., LEE, A. Y., WILLIAMS, B. A., TROUT, D., AMRHEIN, H., FANG, R., CHEN, H., LI, B., VISEL, A., PENNACCHIO, L. A., REN, B. and ECKER, J. R. (2017). Spatiotemporal DNA Methylome Dynamics of the Developing Mammalian Fetus. *bioRxiv*.

HEBESTREIT, K., DUGAS, M. and KLEIN, H. U. (2013). Detection of significantly differentially methylated regions in targeted bisulfite sequencing data. *Bioinformatics* **29**, 1647-1653.

JAFFE, A. E., MURAKAMI, P., LEE, H., LEEK, J. T., FALLIN, M. D., FEINBERG, A. P. and IRIZARRY, R. A. (2012). Bump hunting to identify differentially methylated regions in epigenetic epidemiology studies. *Int J Epidemiol* **41**, 200-209.

JONES, R. H. and BOADI-BOATENG, F. (1991). Unequally Spaced Longitudinal Data with AR(1) Serial Correlation. *Biometrics* **47**, 161-175.

1    JUHLING, F., KRETZMER, H., BERNHART, S. H., OTTO, C., STADLER, P. F. and HOFFMANN, S.
2         (2016). metilene: fast and sensitive calling of differentially methylated regions from
3         bisulfite sequencing data. *Genome Res* **26**, 256-262.

4    KHAMIS, A. M., LIOZNOVA, A. V., ARTEMOV, A. V., RAMENSKY, V., BAJIC, V. B. and
5         MEDVEDEVA, Y. A. (2017). CpG traffic lights are markers of regulatory regions in
6         humans. *bioRxiv*.

7    KRUEGER, F. and ANDREWS, S. R. (2011). Bismark: a flexible aligner and methylation caller for
8         Bisulfite-Seq applications. *Bioinformatics* **27**, 1571-1572.

9    LEE, W. and MORRIS, J. S. (2016). Identification of Differentially Methylated Loci Using
10        Wavelet-Based Functional Mixed Models. *Bioinformatics* **32**, 664-672.

11   LEEK, J. T., SCHARPF, R. B., BRAVO, H. C., SIMCHA, D., LANGMEAD, B., JOHNSON, W. E.,
12        GEMAN, D., BAGGERLY, K. and IRIZARRY, R. A. (2010). Tackling the widespread and
13        critical impact of batch effects in high-throughput data. *Nat Rev Genet* **11**, 733-739.

14   LOADER, C. (1999). *Local Regression and Likelihood*. New York: Springer.

15   LOVE, M. I., HUBER, W. and ANDERS, S. (2014). Moderated estimation of fold change and
16        dispersion for RNA-seq data with DESeq2. *Genome Biol* **15**, 550.

17   LUN, A. T. and SMYTH, G. K. (2014). De novo detection of differentially bound regions for ChIP-
18        seq data using peaks and windows: controlling error rates correctly. *Nucleic Acids Res*
19        **42**, e95.

20   MARX, V. (2016). Genetics: profiling DNA methylation and beyond. *Nat Methods* **13**, 119-122.

21   PACIS, A., TAILLEUX, L., MORIN, A. M., LAMBOURNE, J., MACISAAC, J. L., YOTOVA, V.,
22        DUMAINE, A., DANCKAERT, A., LUCA, F., GRENIER, J. C., HANSEN, K. D., GICQUEL, B.,
23        YU, M., PAI, A., HE, C., TUNG, J., PASTINEN, T., KOBOR, M. S., PIQUE-REGI, R., GILAD,
24        Y. and BARREIRO, L. B. (2015). Bacterial infection remodels the DNA methylation
25        landscape of human dendritic cells. *Genome Res* **25**, 1801-1811.

26   PARK, Y., FIGUEROA, M. E., ROZEK, L. S. and SARTOR, M. A. (2014). MethylSig: a whole
27        genome DNA methylation analysis pipeline. *Bioinformatics* **30**, 2414-2422.

28   PARK, Y. and WU, H. (2016). Differential methylation analysis for BS-seq data under general
29        experimental design. *Bioinformatics* **32**, 1446-1453.

30   PINHEIRO, J., BATES, D., SARKAR, S. D. D. and R CORE TEAM (2017). nlme: Linear and
31        Nonlinear Mixed Effects Models. pp. https://CRAN.R-project.org/package=nlme.

32   ROBINSON, M. D., KAHRAMAN, A., LAW, C. W., LINDSAY, H., NOWICKA, M., WEBER, L. M. and
33        ZHOU, X. (2014). Statistical methods for detecting differentially methylated loci and
34        regions. *Front Genet* **5**, 324.

1 SAITO, Y., TSUJI, J. and MITUYAMA, T. (2014). Bisulfighter: accurate detection of methylated
2    cytosines and differentially methylated regions. *Nucleic Acids Res* **42**, e45.

3 SCHULTZ, M. D., HE, Y., WHITAKER, J. W., HARIHARAN, M., MUKAMEL, E. A., LEUNG, D.,
4    RAJAGOPAL, N., NERY, J. R., URICH, M. A., CHEN, H., LIN, S., LIN, Y., JUNG, I., SCHMITT,
5    A. D., SELVARAJ, S., REN, B., SEJNOWSKI, T. J., WANG, W. and ECKER, J. R. (2015).
6    Human body epigenome maps reveal noncanonical DNA methylation variation. *Nature*
7    **523**, 212-216.

8 SHAFI, A., MITREA, C., NGUYEN, T. and DRAGHICI, S. (2017). A survey of the approaches for
9    identifying differential methylation using bisulfite sequencing data. *Brief Bioinform*.

10 SIEGMUND, D. O., ZHANG, N. R. and YAKIR, B. (2011). False discovery rate for scanning
11    statistics. *Biometrika* **98**, 979-985.

12 SMITH, Z. D. and MEISSNER, A. (2013). DNA methylation: roles in mammalian development. *Nat*
13    *Rev Genet* **14**, 204-220.

14 SUN, D., XI, Y., RODRIGUEZ, B., PARK, H. J., TONG, P., MEONG, M., GOODELL, M. A. and LI, W.
15    (2014). MOABS: model based analysis of bisulfite sequencing data. *Genome Biol* **15**,
16    R38.

17 WEN, Y., CHEN, F., ZHANG, Q., ZHUANG, Y. and LI, Z. (2016). Detection of differentially
18    methylated regions in whole genome bisulfite sequencing data using local Getis-Ord
19    statistics. *Bioinformatics* **32**, 3396-3404.

20 WU, H., XU, T., FENG, H., CHEN, L., LI, B., YAO, B., QIN, Z., JIN, P. and CONNEELY, K. N.
21    (2015). Detection of differentially methylated regions from whole-genome bisulfite
22    sequencing data without replicates. *Nucleic Acids Res* **43**, e141.

23 YU, X. and SUN, S. (2016). HMM-DM: identifying differentially methylated regions using a
24    hidden Markov model. *Stat Appl Genet Mol Biol* **15**, 69-81.

25 ZILLER, M. J., HANSEN, K. D., MEISSNER, A. and ARYEE, M. J. (2015). Coverage
26    recommendations for methylation analysis by whole-genome bisulfite sequencing. *Nat*
27    *Methods* **12**, 230-232, 231 p following 232.
28

## Tables

**Table 1: Summary of datasets used.** Summary measures include the number of samples per population ('Samples'), the number of CpGs with at least one read in all samples in the population ('CpGs Covered'), median number of reads mapping to each covered CpG ('Median Coverage'), minimum and maximum number of reads mapping to each covered CpG ('Coverage Range'). Since the number of CpGs and their coverage are identical in the null comparisons and DMR simulations, the entries for N2 and D2 are combined. Likewise for N3 and D3.

| Dataset | Populations | Samples | CpGs Covered | Median Coverage Range | Maximum Coverage Range |
|---|---|---|---|---|---|
| Human Tissues | Heart, Left Ventricle | 2 | 27458696 | 59-71 | 453000-1473499 |
| | Heart, Right Ventricle | 2 | 27340755 | 27-59 | 554455-779621 |
| | Sigmoid Colon | 2 | 27477877 | 70-76 | 564656-671429 |
| | Small Intestine | 2 | 27344594 | 22-71 | 269326-758025 |
| Murine Leukemia | ALL | 2 | 17666741 | 5 | 2848074-3274608 |
| | AML | 2 | 18306783 | 6-8 | 2279583-2491520 |
| | Control | 2 | 18661620 | 7-9 | 3207310-4909532 |
| Simulated | Simulations N2 & D2 | 2 | 22015096 | 9-10 | 200-236 |
| | Simulations N3 & D3 | 3 | 21795211 | 9-10 | 200-236 |

**Table 2: Null comparison results for sample size 2 (N2) and sample size 3 (N3).** Numbers of DMRs identified by dmrseq and metilene are shown at the 0.05 FDR level. Default settings were used for BSmooth and DSS.

| Null Comparison | Method | DMRs (FPs) |
|---|---|---|
| N2 | dmrseq | 0 |
| | BSmooth | 76,563 |
| | DSS | 661 |
| | metilene | 31 |
| N3 | dmrseq | 0 |
| | BSmooth | 76,319 |
| | DSS | 770 |
| | metilene | 27 |

**Table 3: Simulated DMR results for sample size 2 (D2) and sample size 3 (D3).** Numbers of DMRs identified by dmrseq and metilene are shown at the 0.05 FDR level. Default settings were used for BSmooth and DSS. True Positives (TPs) is the number of simulated DMRs that are overlapped by at least one identified DMR. False Positives (FPs) are DMRs that do not overlap any of the simulated DMRs.

| Simulation | Method | DMRs | TPs (unique) | FPs |
|---|---|---|---|---|
| | dmrseq | 914 | 816 | 42 |
| D2 | BSmooth | 73,252 | 2,466 | 70,688 |
| | DSS | 2,086 | 762 | 655 |
| | metilene | 329 | 210 | 30 |
| | dmrseq | 1,620 | 1,455 | 78 |
| D3 | BSmooth | 72,764 | 2,646 | 69,999 |
| | DSS | 2,858 | 1,257 | 763 |
| | metilene | 652 | 441 | 27 |

**Table 4: Tissue-specific DMR results.** Number of DMRs found by dmrseq and metilene at FDR level 0.05, and BSmooth and DSS under default settings.

| Tissue Comparison | dmrseq | BSmooth | DSS | metilene |
|---|---|---|---|---|
| Left Ventricle vs Right Ventricle | 0 | 88,443 | 6,312 | 24 |
| Sigmoid Colon vs Small Intestine | 14,695 | 75,968 | 51,744 | 949 |
| Left Ventricle vs Small Intestine | 33,740 | 76,078 | 153,217 | 6,344 |
| Left Ventricle vs Sigmoid Colon | 106,461 | 76,307 | 229,729 | 8,133 |
| Right Ventricle vs Small Intestine | 32,143 | 76,334 | 129,106 | 5,756 |
| Right Ventricle vs Sigmoid Colon | 73,431 | 76,643 | 196,998 | 7,692 |

**Table 5: Murine Leukemia model DMR results.** Number of DMRs found by dmrseq and metilene at FDR level 0.10, and BSmooth and DSS under default settings.

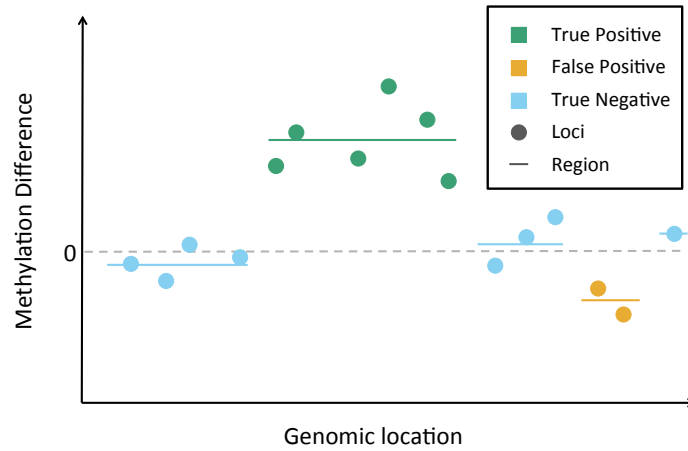| Condition Comparison | dmrseq | BSmooth | DSS | metilene |
|---|---|---|---|---|
| AML vs Control | 16,465 | 43,818 | 21,256 | 3,004 |
| ALL vs Control | 8,855 | 51,723 | 16,478 | 3,182 |
| AML vs ALL | 9,253 | 50,004 | 23,582 | 3,360 |

# Figures



**Figure 1: Illustration of why FDR at the loci level is not the same as FDR at the region level.** This schematic shows a plot of genomic location versus methylation difference estimates at several neighboring loci. The individual CpGs (points) are shaded by whether they are a true or false positive. Regions are denoted by lines. The loci FDR is $FDR_{loci} = \frac{\# \text{ False Positive Loci}}{\text{Total } \# \text{ of Significant Loci}}$, which is equal to 0.25 in this example. The region FDR is $FDR_{region} = \frac{\# \text{ False Positive Regions}}{\text{Total } \# \text{ of Significant Regions}}$, which is equal to 0.50 in this example.



**Figure 2: dmrseq provides accurate FDR control of regions.** Specified versus observed region-level FDR level is plotted for two different sample size settings from simulated data for dmrseq. Note that region-level FDR cannot be specified for BSmooth or DSS, and results for metilene are shown in Supplementary Figure S3.
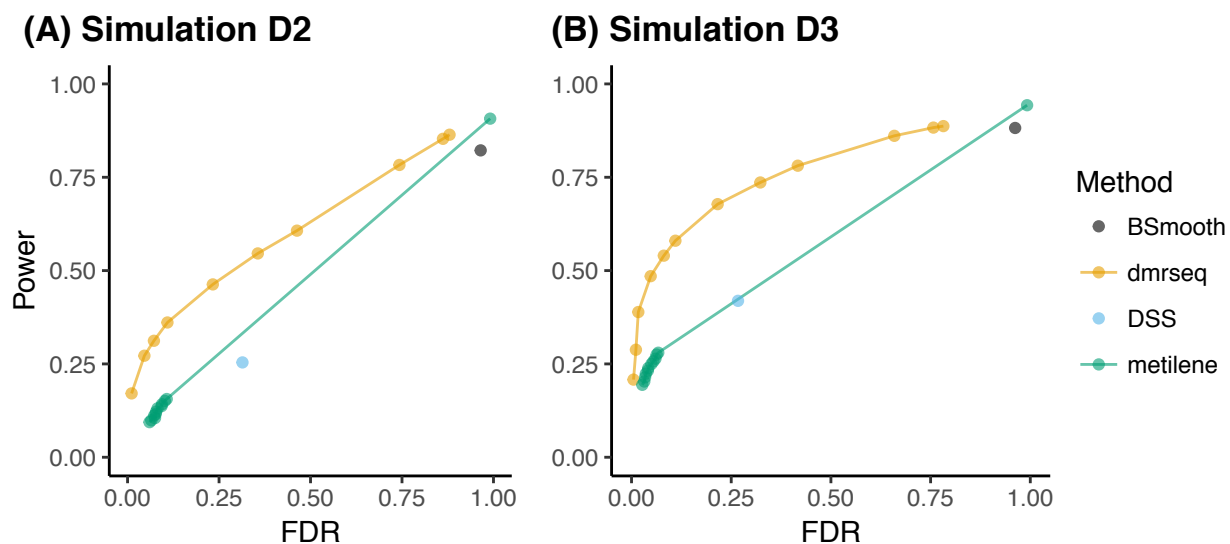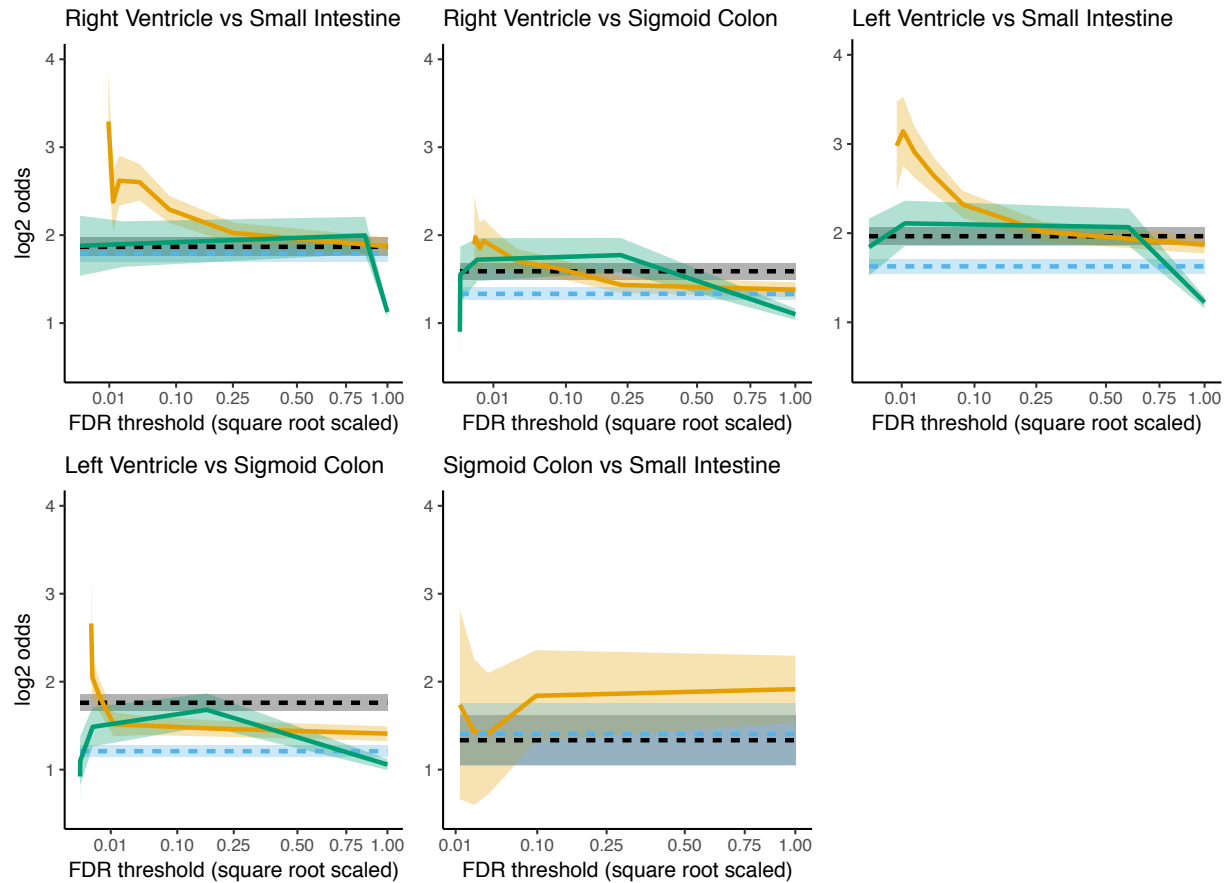
**Figure 3: dmrseq is more powerful than other methods.** FDR and power results for (A) Simulation D2 and (B) Simulation D3, with method denoted by color. dmrseq and metilene results are displayed for several different FDR cutoffs. Since region level FDR control is not possible for BSmooth and DSS, results using default settings are displayed. Power is calculated as the proportion of simulated DMRs overlapped by at least one identified DMR. FDR is calculated as the proportion of DMRs identified that do not overlap with any of the simulated DMRs.
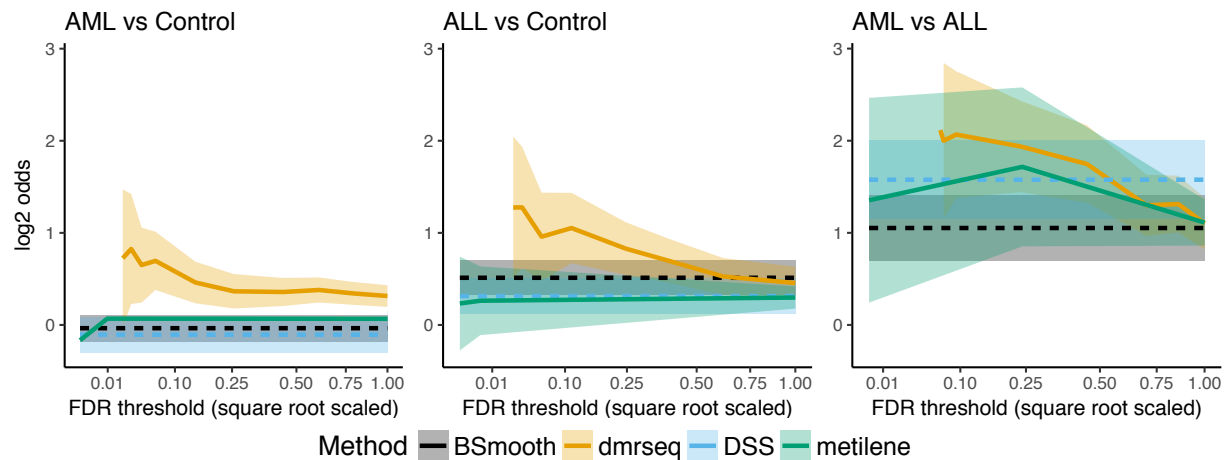
**Figure 4: dmrseq achieves stronger inverse association of methylation and differential expression at lower FDR thresholds.** Odds of inverse association between methylation difference of (A) tissue-specific DMRs and (B) murine leukemia DMRs with differential expression of nearby DE genes (log2 transformed) is displayed on the y-axis. For dmrseq and metilene, the x-axis represents the FDR threshold (square-root scaled) for which the odds calculation (cumulative) is performed. Since FDR cannot be specified for BSmooth or DSS, the odds are calculated over all DMRs identified and displayed as a horizontal line. Note that the comparison between Left and Right Ventricles is not shown, since no DE genes were identified.
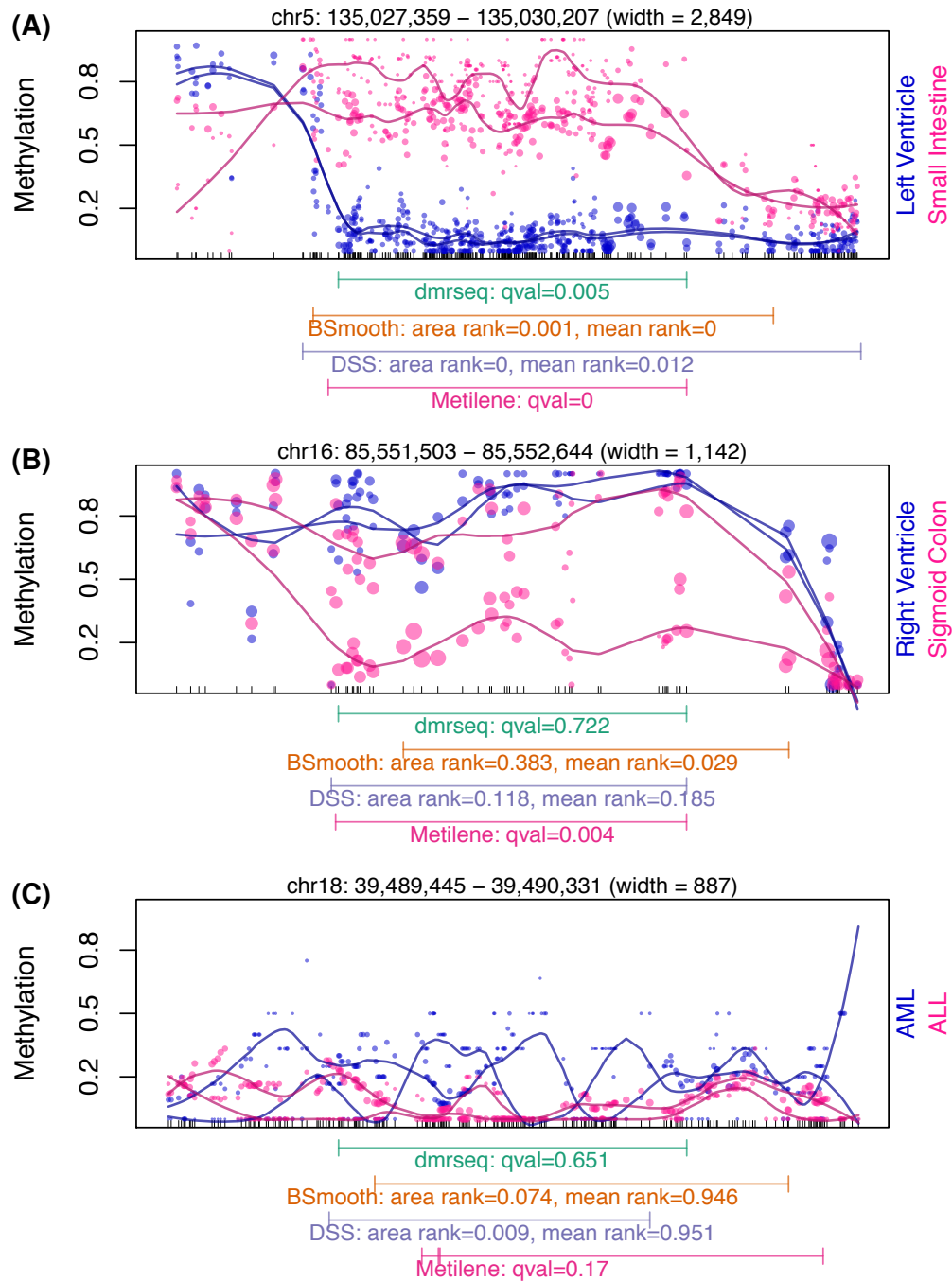
**Figure 5: dmrseq ranks regions by statistical significance.** Example regions from the human tissue and murine leukemia studies are displayed for three cases that illustrate the increased variability of regions that are highly ranked by area or mean difference statistics of BSmooth and DSS but not dmrseq. For each case, the q-value is shown for dmrseq and metiline, and the rank percentile by the area statistic and mean difference statistics are both shown for BSmooth and DSS (see Supplement Section 2.8 for details). (A) All methods assign a consistently high rank. (B) dmrseq assigns a low rank, but the mean difference statistic of BSmooth and DSS assign a high rank. (C) dmrseq assigns a low rank, but the area statistic of BSmooth and DSS assign a high rank. The condition comparison is indicated by the labels to the right of each plot.