

Genome analysis

D³M: detection of differential distributions of methylation levels

Yusuke Matsui^{1,*}, Masahiro Mizuta², Satoshi Ito³, Satoru Miyano³ and Teppei Shimamura^{1,*}

¹Nagoya University Graduate School of Medicine, Nagoya 466-8550, Japan, ²Information Initiative Center, Hokkaido University, Sapporo 060-0811, Japan and ³Institute of Medical Science, The University of Tokyo, Tokyo 108-8639, Japan

*To whom correspondence should be addressed.

Associate Editor: John Hancock

Received on June 22, 2015; revised on February 19, 2016; accepted on March 8, 2016

Abstract

Motivation: DNA methylation is an important epigenetic modification related to a variety of diseases including cancers. We focus on the methylation data from Illumina's Infinium HumanMethylation450 BeadChip. One of the key issues of methylation analysis is to detect the differential methylation sites between case and control groups. Previous approaches describe data with simple summary statistics or kernel function, and then use statistical tests to determine the difference. However, a summary statistics-based approach cannot capture complicated underlying structure, and a kernel function-based approach lacks interpretability of results.

Results: We propose a novel method D³M, for detection of differential distribution of methylation, based on distribution-valued data. Our method can detect the differences in high-order moments, such as shapes of underlying distributions in methylation profiles, based on the Wasserstein metric. We test the significance of the difference between case and control groups and provide an interpretable summary of the results. The simulation results show that the proposed method achieves promising accuracy and shows favorable results compared with previous methods. Glioblastoma multiforme and lower grade glioma data from The Cancer Genome Atlas show that our method supports recent biological advances and suggests new insights.

Availability and Implementation: R implemented code is freely available from <https://github.com/ymatts/D3M/>.

Contact: ymatsui@med.nagoya-u.ac.jp or shimamura@med.nagoya-u.ac.jp

Supplementary information: [Supplementary data](#) are available at *Bioinformatics* online.

1 Introduction

DNA methylation is an epigenetic chemical alternation in which a methyl group is attached to a 5-carbon of a cytosine (C) base. It is closely related to gene expression, silencing and genomic imprinting, including oncogenesis. Typically, methylation is explained as occurring in cytosine-phosphate-guanine (CpG) sites in mammals. The methylation of promoter regions, in particular, silences cancer suppressor genes (Baylin, 2005; Kulis and Esteller, 2010).

We focus on the methylation data from Illumina's Infinium HumanMethylation450 BeadChip. One of the key issues for

methylation analysis is to detect differentially methylated sites, i.e. a significant difference in methylation levels between the case and the control groups at a site. When comparing groups, we often summarize (or aggregate) data in summary statistics, such as mean and variance, and then investigate the difference between the groups. For example, IMA (Wang *et al.*, 2012) detects the differentially methylated sites by *T*-test, Empirical Bayes (EB) method or by Mann–Whitney–Wilcoxon (MWW) test. DiffVar (Phipson and Oshlack, 2014) detects the sites by testing significant difference of variance. Other nonparametric approaches exist, such as the Kolmogorov–Smirnov test (KS) or

kernel-based approaches, such as maximum mean discrepancy (MMD) (Gretton *et al.*, 2012). In particular, since KS and MMD consider the underlying distribution structure, they are better suited for use with complicated distributions than methods based on summary statistics.

These approaches are insufficient from some perspectives; underlying distributions are complicated by being skewed, heavy-tailed and multimodal. In particular, since cancer cells include heterogeneities, measurements of the methylation levels potentially include complex distribution shapes. This observation indicates that we need to consider the underlying structure. The disadvantage of KS and MMD is infeasible interpretability of results because they measure the maximum and kernel distances of distributions, respectively, which are difficult to interpret corresponding to the actual difference of underlying distributions.

We develop a method to detect differential methylation sites with distribution-valued data (Irpino and Verde, 2014). Distribution-valued data are an example of symbolic data analysis (Diday, 1989). This framework can treat complex data such as functional (Ramsey and Silverman 2005), tree (Wang and Marron, 2007), set, interval and histogram values (Bock and Diday, 2000). The proposed method describes case and control groups using distribution values. We measure the differences between distributions using the Wasserstein metric. We detect the differential methylation sites using a statistical test of significant differences of distribution functions.

2 Methods

Our method is aimed at a distribution-based comparison of methylation levels in two groups at a single cytosine level. We construct distribution functions representing the two groups at each site. Next, we compare the groups using a dissimilarity measure and test statistical significance by resampling approach at a site. We adopt an L_2 -Wasserstein metric (Rueschendorff, 2011) as a dissimilarity measure, a distribution function-based measure of statistical distance. The advantage of this distance is the interpretability of results because the distance can be decomposed into three components, i.e. mean, variance distribution shape. This fact leads to visualization of results using a Q - Q plot to interpret the detected distribution difference including hypo- or hyper-methylation status. Methylation level can be represented by a beta value that is the ratio between the methylated probe intensity and the overall intensity (Du *et al.*, 2010). The definition of i th site beta value in Illumina methylation assay is as follows:

$$\text{Beta}_i = \frac{\max(z_{i,\text{methy}}, 0) + \alpha}{\max(z_{i,\text{unmethy}}, 0) + \max(z_{i,\text{methy}}, 0) + \alpha} \quad (1)$$

where $z_{i,\text{methy}}$ and $z_{i,\text{unmethy}}$ are the intensity measured by i th methylated and unmethylated probes, respectively and α is constant. β values have biologically direct interpretation (Du *et al.*, 2010), as opposed to M value (i.e. logit transformation of β value). IMA and DiffVar use β value and M value, respectively. In the following, we use β value as the input of D³M.

2.1 Construction of objects

$X(s_i)$ and $Y(s_i)$ ($i = 1, 2, \dots, S$) represent the beta values in a case group (e.g. cancer subjects) and a control group (e.g. normal

subjects) at a CpG site s_i . We represent the data as distribution values by

$$\begin{aligned} F_i(x) &= \Pr\{X(s_i) \leq x\}; x \in [0, 1]. \\ G_i(y) &= \Pr\{Y(s_i) \leq y\}; y \in [0, 1]. \end{aligned} \quad (2)$$

In practice, let the beta value observations be $x_j(s_i); j = 1, 2, \dots, n$ and $y_j(s_i); j = 1, 2, \dots, m$ following $F_i(x)$ and $G_i(y)$, respectively, where n and m are the respective numbers of observations at s_i . From the data, we construct the empirical distribution functions;

$$\begin{aligned} \hat{F}_i(x) &:= \frac{1}{n} \sum_{j=1}^n \mathbf{1}(x_j(s_i) \leq x) \\ \hat{G}_i(y) &:= \frac{1}{m} \sum_{j=1}^m \mathbf{1}(y_j(s_i) \leq y) \end{aligned} \quad (3)$$

where

$$\mathbf{1}(a \leq b) = \begin{cases} 1 & \text{if } a \leq b \\ 0 & \text{otherwise.} \end{cases} \quad (4)$$

2.2 Dissimilarity measure for distributions

The Wasserstein metric is defined by the following:

$$d^q(F_i, G_i) := \int_0^1 |F_i^{-1}(u) - G_i^{-1}(u)|^q du \quad (5)$$

where $1 \leq q \leq 2$ and $F_i^{-1}(x)$ and $G_i^{-1}(y)$ indicate quantile functions.

In particular, in the case of $q=2$, the metric can be decomposed into three components that describe the distribution characteristics, i.e. mean, variance shape (Irpino and Verde, 2014):

$$\begin{aligned} d^2(F_i, G_i) &= \int_0^1 |F_i^{-1}(u) - G_i^{-1}(u)|^2 du \\ &= (\mu_i - \mu_i')^2 + (\sigma_i^2 - \sigma_i'^2) + 2\sigma_i\sigma_i'(1 - \rho_i) \end{aligned} \quad (6)$$

where μ_i and σ_i^2 (respectively, μ_i' and $\sigma_i'^2$) are mean and variance of $F_i(x)$ (respectively, $G_i(y)$), and ρ_i is the correlation index of the points in the Q - Q plot of F_i and G_i .

The empirical estimator of the Wasserstein metric is given by

$$d^q(\hat{F}_i, \hat{G}_i) = \int_0^1 |\hat{F}_i^{-1}(u) - \hat{G}_i^{-1}(u)|^q du. \quad (7)$$

Technically, we use quantiles to compute the approximation of the (7) for reducing computational costs. Let $(Q_{i,1}, Q_{i,2}, \dots, Q_{i,K})$ and $(Q'_{i,1}, Q'_{i,2}, \dots, Q'_{i,K})$ be k -quantiles of $F_i(x)$ and $G_i(y)$. We calculate $d^2(\hat{F}_i, \hat{G}_i) \approx \sum_{l=1}^K (Q_{i,l} - Q'_{i,l})^2$ in the case of $q=2$, instead of evaluating the integral in (7). Here, we simply write $d_i := d(\hat{F}_i, \hat{G}_i)$.

2.3 Detection of differential methylation sites

We use the metric to investigate whether two distributions are significantly different. We pose statistical hypotheses as follows.

$$\begin{aligned} \text{Null hypothesis : } & F_i = G_i \\ \text{Alternative hypothesis : } & F_i \neq G_i \end{aligned} \quad (8)$$

We use resampling to construct a null distribution. From the null hypothesis (8), we jointly permute the observations $(x_1(s_i), x_2(s_i), \dots, x_n(s_i))$ and $(y_1(s_i), y_2(s_i), \dots, y_m(s_i))$ to obtain the new distribution

Table 1. Simulation models of eight cases μ_i , σ_i , and s_i indicate mean, variance, and shape (unimodal or bimodal) of the distributions in the case ($i=1$) and control ($i=2$) groups, respectively

	$F_1 = F_2$	$F_1 \neq F_2$						
	Case 1	Case 2	Case 3	Case 4	Case 5	Case 6	Case 7	Case 8
$\mu_1 = \mu_2$	T	T	T	F	T	F	F	F
$\sigma_1 = \sigma_2$	T	T	F	T	F	T	F	F
$s_1 = s_2$	T	F	T	T	F	F	T	F

functions $\hat{F}_i^*(x)$ and $\hat{G}_i^*(y)$. Next, we obtain the new distance $d_i^* = d^2(\hat{F}_i^*, \hat{G}_i^*)$ according to (7).

Let $D_i^* = (d_{i,1}^*, d_{i,2}^*, \dots, d_{i,B_{all}}^*)$ be all possible distances for the permutation process. Then P -value is

$$P_{all}(d_i) = \frac{\sum_{b=1}^{B_{all}} \mathbf{1}(d_{i,b}^* \geq d_i)}{B_{all}}. \quad (9)$$

Approximation of (9) uses the subset of D_i^* , $\tilde{d}_{i,1}^*, \tilde{d}_{i,2}^*, \dots, \tilde{d}_{i,B}^*$ where $B \leq B_{all}$:

$$P_{sub}(d_i) = \frac{\sum_{b=1}^B \mathbf{1}(\tilde{d}_{i,b}^* \geq d_i)}{B}. \quad (10)$$

In the simulation of Section 3 and data analysis in Section 4, we set $B = 10,000$.

The number of permutations B is closely related to the accuracy of the P -value. However, resolution of P_{sub} is limited to $1/B$, if we need the very small P -values. One solution is to perform a large number of permutations, but it is computationally expensive. A semi-parametric estimation of P -value is proposed by Knijnenburg et al. (2009) to obtain more accurate P -values.

We use an exponential distribution to estimate the distribution tail as follows,

$$P(d_i) = \begin{cases} \frac{1}{B} \sum_{j=1}^B \mathbf{1}(\tilde{d}_{i,j}^* \geq d_i) & \text{for } d_i < d_i^{(\min)} \\ \exp(-\lambda_i(d_i - d_i^{(\min)})) & \text{for } d_i \geq d_i^{(\min)} \end{cases} \quad (11)$$

where λ_i is a scale parameter and $d_i^{(\min)}$ is a threshold that we set to 99th percentile of null distributions. We estimate λ_i using data above the threshold. Technically, we perform the semi-parametric estimation only if $P_{sub}(d_i)$ reaches to zero.

2.4 Graphical representation of results

The graphical interpretation of the statistical test result is important. One approach is to plot all the distribution (density) functions of candidate sites, but this is infeasible for hundreds of sites. We use a Q - Q plot with two distributions. It enables us to visualize many pairs of distributions at a time, with the directions being easy to interpret. In the actual example shown in Section 4, we plotted 1000 pairs of differentially methylated distributions (Fig. 3B). We can see the hyper-methylation with the most significant 1000 sites (blue lines in Fig. 3B).

3 Simulation

3.1 Simulation setting

We evaluated the proposed method with simulated datasets with focus on single cytosine level in the case and control group. Our

simulation is intended for the detection of differential methylation sites when there is cancer heterogeneity. Here, the cancer heterogeneity is described by the multiple modes of distributions. We conduct a statistical test for $H_0 : F_i = G_i \leftrightarrow H_1 : F_i \neq G_i$ under significance levels 5%, and we compare the results to those of the other methods, i.e. DiffVar, MMD, KS, MWW and EB. We used several packages; missMethyl (Phipson and Oshlack, 2014), limma (Ritchie et al., 2015) written in R and mmd (Gretton et al., 2006) (<http://www.kyb.tuebingen.mpg.de/bs/people/arthur/code/mmd.zip>) written in MATLAB. The setting of missMethyl is default and mmd with options $\alpha = 0.5$ and $\text{MMD_METHOD} = \text{'approxmoments'}$.

We describe the outline of the simulation as follows. The details are described in Supplementary file S1 and R codes are described in Supplementary file S2. The data are generated by using two types of distributions. The control and case groups are represented by normal and normal mixture distributions, respectively. In each case (case 1–case 8 in Table 1), there are 80 samples; 40 samples for case and control groups, respectively (Fig. 1). We performed the statistical test with each method for 50 times in every case (i.e. cases 1–8) and evaluate averages of a type I error and power. Besides we repeat this process for 100 times to assess the variances of the averages.

3.2 Simulation results

The results are shown in Table 2. In the first case, it is shown that error rates of D³M, MMD, DiffVar, KS, EB and MWW are close to the significance levels, which indicates that they effectively control type I errors.

Furthermore, we investigate the power for cases 2–8. The KS test and MMD show relatively good performance in case 2 where only the shape parameters of the distributions differ. However, they cannot capture the feature of case 8 where the majority of the two groups are overlapped with each other although 15% of minority of distribution exists. Besides, KS testing fails to detect case 6 where 15% of minority distribution is hyper-methylated and MMD fails to detect case 4 where only mean is different. DiffVar shows high power for cases where explicit hypothesis is tested, however, it might capture the other distribution features for the cases with equal variances (case 6), leading to uninterpretable results. EB can appropriately distinguish only the mean difference that is the explicit hypothesis. MWW can detect case 4 and 7, but cannot detect cases in which the means differ under non-normality. The proposed D³M provides preferable results over the eight cases.

On one hand, from this simulation, the differential distributions in each case can be detected with D³M with the sufficient power, which indicates that it works well for the given situations and it can be applied to various cases flexibly. On the other hand, a small sample situation should also be examined. We conduct the simulations in Supplementary file S1 considering several situations such as; small sample size (25 samples) with the same setting of Table 2 (Supplementary Tables S.3 and S.6), sample sizes are extremely different between the case and control groups (Supplementary Table

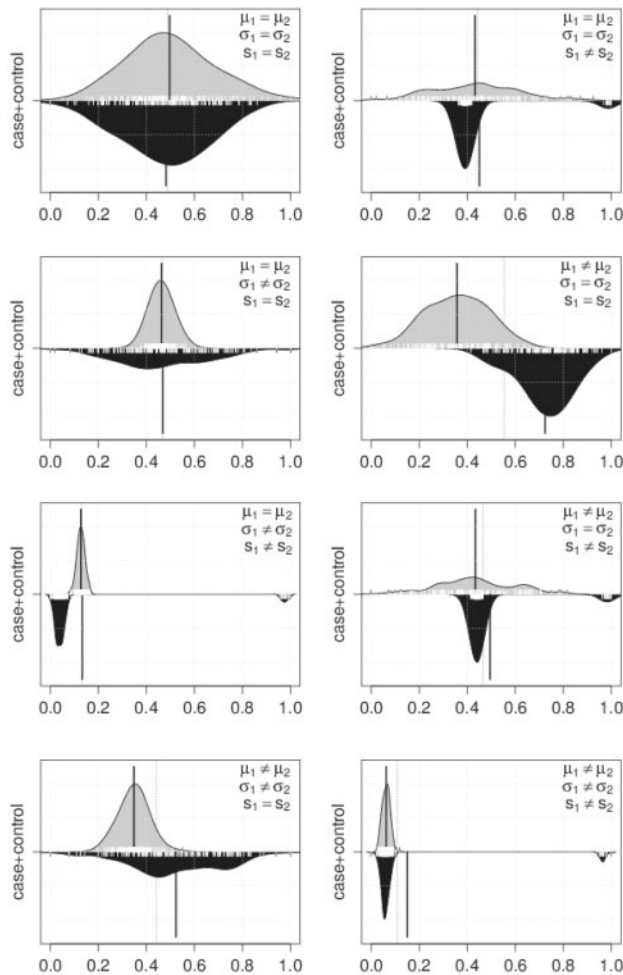


Fig. 1. The beanplot of eight cases

S.5). In general, the power of D³M is decreasing as well as the other methods.

4 Actual example

4.1 Datasets

We apply our method to methylation data of glioblastoma multiforme (GBM) and lower grade glioma (LGG) from The Cancer Genome Atlas (TCGA). GBM is the primary brain tumor that progresses with malignant invasion destroying normal brain tissues (TCGA, 2008). In this analysis, we compare the methylation levels in the LGG and GBM groups, and then specify the differential methylation sites. We focus on mean, variance and shape differences using EB, DiffVar, KS and D³M and compare the results.

Here, we briefly describe the datasets and preprocessing as follows. All the samples are hybridized to Illumina Infinium HumanMethylation450K arrays, including 485 577 CpG sites, which is downloadable from TCGA portal sites. Each CpG site contains 145 samples and 530 samples in GBM and LGG, respectively. We remove CpG sites on the X and Y chromosomes and SNP control probes(rs1–rs65). We also use *HumanMethylation450 v1.2 SNP Update Table* to remove SNP related probes by *minor allele frequency* $\leq 1\%$. As a result, we get 351 932 probes. Missing values in both groups are inferred using R package *pcaMethods* (Stacklies, 2007) with *pca* functions since DiffVar does not accept missing

values and the number of the probes including missing values cannot be ignored. We remove the batch effect using the ComBat function in SVA package (Leek *et al.*, 2012) with default settings. We use batches having more than 2 samples. The details are described in R codes in S3. We finally obtain 141 GBM and 530 LGG samples.

4.2 Analysis results

Significant differential methylation sites were identified as those having false discovery rate (*q*-value) (Benjamini and Hochberg, 1995) less than 1%. The Venn diagram (Fig. 2A) shows the number of detected probes from the perspective of mean, variance and shape difference of distributions using EB, DiffVar, KS and D³M, and 279 008, 191 050, 297 493 and 255 317 sites are totally detected, respectively. From Figure 2A, most of detected sites with the shape difference are overlapped with the sites by the mean and the variance differences, respectively. However, when focusing on top 1000 significant sites (Fig. 2B), the overlaps between them become very small. This suggests that the ‘signal’ of differentially methylated sites in terms of *q*-value is quite different from each other. The ‘signal’ is important for the further analysis, such as the pathway analysis since we often use the filtered gene set, e.g. using top 1000 significant sites. We also calculate the Spearman’s rank correlation between the four methods (D³M, EB, KS and DiffVar) resulting from *q*-values (Supplementary Table S.7).

Among the detected sites with D³M, we investigated sites with the smallest 1000 *q*-value. Heat map and Q–Q plots of the top 1000 sites are shown in Figure 3A and B. Comparing heat maps and Q–Q plots, the methylation levels are easy to interpret in the latter. From the Q–Q plot, we could see that the top 1000 sites tend to be hyper-methylated in LGG (with the reverse in GBM).

Among the top significant 1000 pairs of distributions of GBM and LGG, we can observe that there are mainly two patterns in terms of the Q–Q plot. Then, we cluster the 1000 curves in Q–Q plot into two classes; we consider input data as $Z_i = (Q_{i,1}, Q_{i,2}, \dots, Q_{i,K}, Q'_{i,1}, Q'_{i,2}, \dots, Q'_{i,K}) (i = 1, 2, \dots, 1000)$ and define the dissimilarity as $\sqrt{(Z_i - Z_j)^2 (i \neq j)}$ where $Q_{i,k}$ and $Q'_{i,k} (k = 1, 2, \dots, K)$ is *k*th quantile in LGG and GBM group, respectively. Then we apply the standard hierarchical clustering (Clusters 1 and 2 in Fig. 3A). Typical distribution examples in each cluster are shown in Figure 3C. Clusters 1 and 2 contain 715 and 209 sites, respectively.

Next, we perform enrichment analysis on gene sets in clusters 1 and 2. We used ingenuity pathway analysis (IPA) for 444 and 209 genes in clusters 1 and 2, respectively, and significantly enriched pathways in each cluster using Fisher’s exact test. Table 3 shows five pathways and related genes with *q*-values $\leq 10^{-3}$ in each cluster. Other pathways are described in Supplementary files S4 (cluster 1) and S5 (cluster 2).

The pathways in clusters 1 and 2 include significant pathways in GBM, which have been previously reported even though we do not include any information on GBM. The axonal guidance signaling pathway in cluster 1 has been suggested as prompting the cell invasion of GBM (Hoelzinger *et al.*, 2007) and ERK/MAPK Signaling is reported to be up-regulated in GBM (Liu *et al.*, 2013). The enrichments of Caveolar-mediated Endocytosis Signaling and Calcium Signaling are studied in (Dong *et al.*, 2010; Polisetty *et al.*, 2012). The remaining pathways might be explained elsewhere. Our prediction using D³M provides a hypothesis that DNA methylation in these pathways might cause the phenotypical difference between GBM and LGG.

Table 2. Hypothesis testing in each case (hypotheses are described under the method names)

		Type I Case 1	Power Case 2	Case 3	Case 4	Case 5	Case 6	Case 7	Case 8
D ³ M	Mean	5.12	94.32	95.42	84.52	98.60	84.88	99.12	100.00
$H_0 : F_1 = F_2 \ H_1 : F_1 \neq F_2$	sd	0.03	0.04	0.03	0.48	0.03	0.05	0.01	0.00
DiffVar	Mean	1.80	4.92	93.86	1.30	97.86	11.84	68.26	86.32
$H_0 : V(F_1) = V(F_2) \ H_1 : V(F_1) \neq V(F_2)$	sd	0.02	0.03	0.04	0.02	0.04	0.04	0.06	0.15
MMD	Mean	4.96	100.00	90.58	67.16	98.30	81.24	95.98	82.54
$H_0 : F_1 = F_2 \ H_1 : F_1 \neq F_2$	sd	0.07	0.00	0.10	0.10	0.14	0.08	0.11	0.17
KS	Mean	2.56	99.94	32.06	67.06	76.40	60.22	89.62	75.74
$H_0 : F_1 = F_2 \ H_1 : F_1 \neq F_2$	sd	0.02	0.00	0.07	0.213	0.07	0.04	0.04	0.23
EB	Mean	4.94	0.76	5.18	87.92	0.08	78.50	87.36	96.98
$H_0 : E(F_1) = E(F_2) \ H_1 : E(F_1) \neq E(F_2)$	sd	0.03	0.01	0.03	0.04	0.00	0.05	0.04	0.05
MWW	Mean	4.68	12.00	5.76	84.96	55.42	55.50	83.94	72.76
$H_0 : M(F_1) = M(F_2) \ H_1 : M(F_1) \neq M(F_2)$	sd	0.03	0.06	0.03	0.05	0.14	0.07	0.05	0.17

We focus on a probe and generate beta values, 40 samples each in a case and a control group respectively. We apply six methods to this data and evaluate the *P*-value at the significance level, i.e. 5%. This process is repeated 50 times and returns the averages of type I error (%) and power (%). We represent the power in italics within top three. Furthermore, we also evaluate the standard deviations of the averages; we repeated the process of the evaluating average of type I error and power a 100 times and then obtained the standard deviations of the averages.

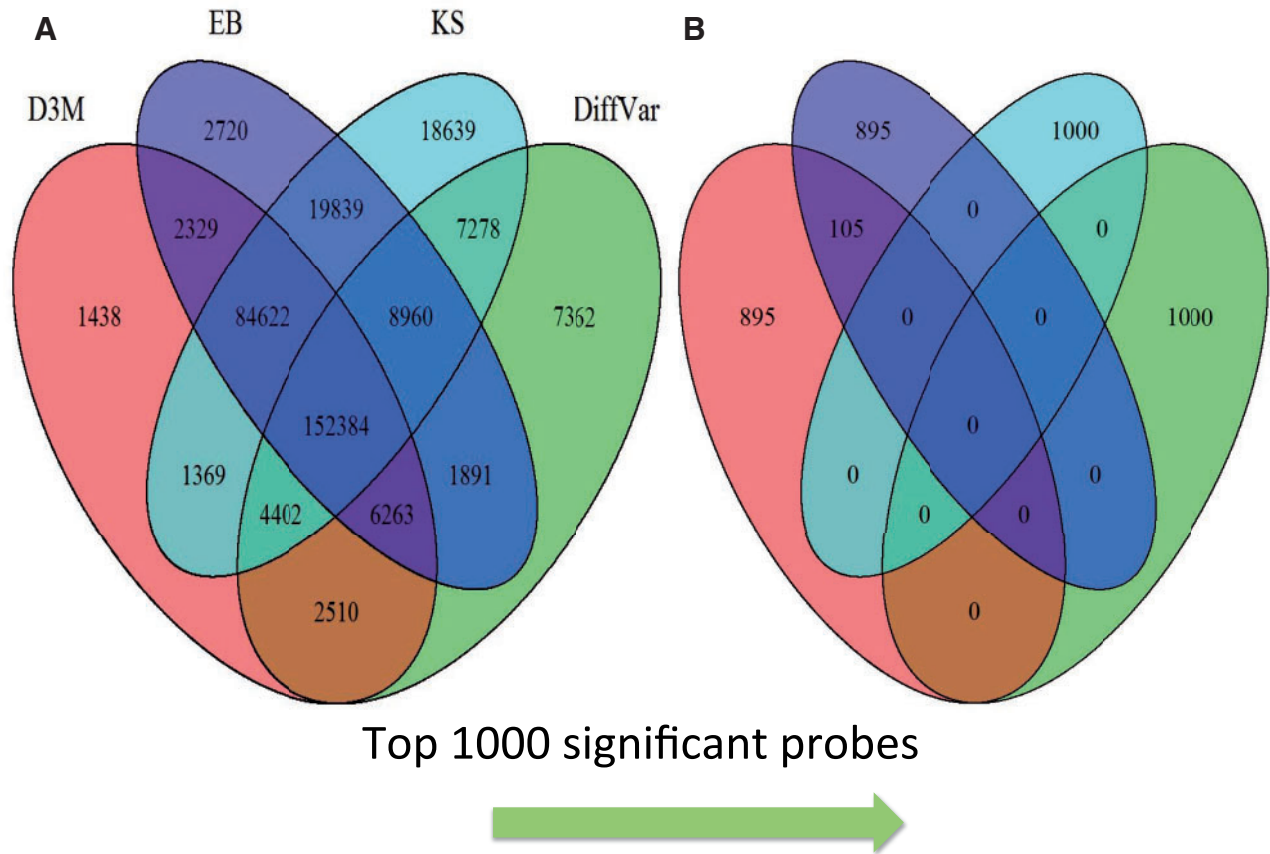


Fig. 2. Venn diagram of detected sites at the significance level 1% (A) and of top significant 1000 sites (B)

We further focus on Wnt in Human Embryonic Stem Cell Pluripotency pathway, and then compare the ranking based on *P*-value by D³M with those by other methods. The activation of Wnt family is closely related to cell differentiation of GBM (Rampazzo et al., 2013). In our analysis, there are 18 Wnt genes on the probes. Among them, six Wnt genes (Wnt2, Wnt2b, Wnt3, Wnt4, Wnt7a and Wnt 9a) are included in top 1000 significantly differentially methylated probes using D³M.

We investigate the enrichment of the six Wnt genes in top 1000 significant probes using Fisher's exact test and we confirm the enrichment of the genes (the details are described in Section S-3-3 in Supplementary file S1). The six Wnt genes are included in both the clusters 1 and 2, and the majority of the distribution is hypo-methylated and the minority is hyper-methylated in GBM, vice versa in LGG. This suggests that demethylation of Wnt might trigger the activation of Wnt family and prompt cell differentiation.

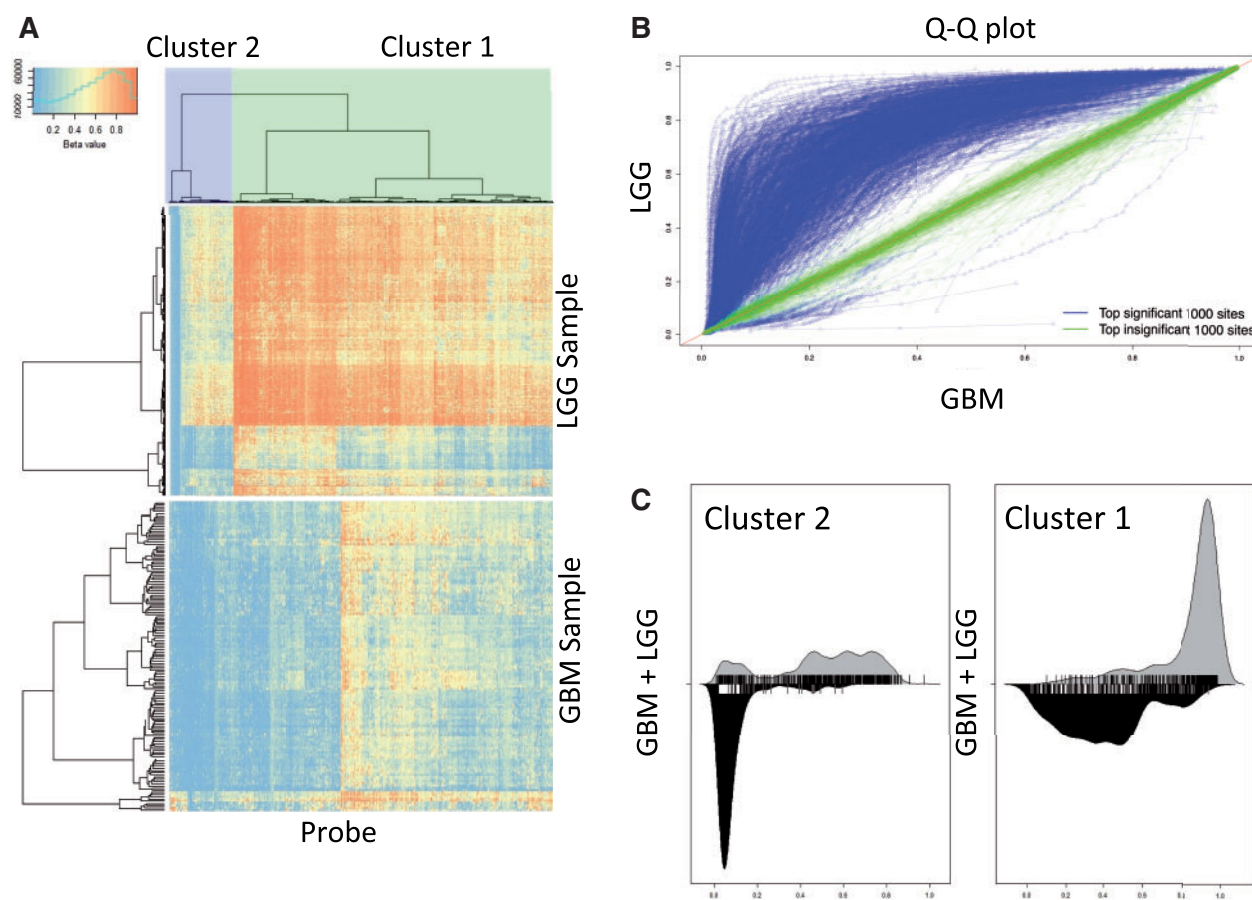


Fig. 3. Distributions of the top 1000 significant probes. (A) The heat map of beta values with 141 LGG and 530 GBM samples at each probe respectively. We jointly clustered quantiles of LGG and GBM probes into the two clusters based on the Wasserstein metric (i.e. we clustered curves in Q-Q plots of the top significant probes into two clusters). (B) Q-Q plots of beta values between the LGG and GBM probes (blue lines). D³M detects the probes with hyper-methylated LGG probes compared with GBM probes. Green lines are the top 1000 insignificant probes (as negative controls). (C) The instances of distributions with Clusters 1 and 2 shown in (A-1) and (A-2), respectively. The distributions appear to be highly heterogeneous especially in LGG probes. D³M recognizes strongly heterogeneous distributions and detects such probes as top significant probes

The ranking of Wnt genes in D³M, EB, KS and DiffVar is shown in Table 4.

5 Discussion

Here, we summarize the advantages and disadvantages of D³M, DiffVar as well as MMD. These methods are designed for detecting differential methylation levels focusing on cancer heterogeneity, which is caused by epigenetic instability and diversity. Cancer heterogeneity can often be confused with outliers. For example, DiffVar fails to detect simulation case 6 as differential methylation, even though we set the variance, but not the mean and the shapes, to be the same for the two groups. This is because DiffVar deals with minority distributions as outliers and evaluates only those in the majority.

In general, the significance of an outlier depends on the context of analysis (Aggarwal, 2013). When an outlier arises from measurement error not relevant to signals of interest, we must remove them prior to analysis. In contrast, when an outlier arises from an unusual event including new findings that we seek, we use them for further analysis. In this case, cancer heterogeneity could be regarded as an abnormal event compared with normal cases, and thus must be included in the analysis.

MMD was originally developed as a distribution-free two sample test based on a kernel. Both the Wasserstein metric and MMD are family of integral probability metrics (Müller, 1997) designed to capture not only such basic difference as those between the means but also such higher-order differences as those between the distribution shapes of the groups. From the simulation in Section 3.1, we see that MMD shows good detection performance in cases except cases 4 and 8. D³M shows good performance over all eight cases. Specifically, D³M shows excellent performance in cases 7 and 8.

In the study of sample size effect in Section S-2-1, using the same simulation model as in Section 3.1, the result indicates that the superiority of D³M to other methods are reduced compared to large sample case, although all the methods decrease in power (Supplementary Table S.3). On the contrary, MMD retains its power in cases 2 and 6, and so does D³M in cases 7 and 8.

We also investigate the power of D³M with other simulation models using beta distributions considering small sample size in Section S-2-2 and S-2-3 in Supplementary file S1. The result indicates that MMD detects differences well in the small samples ($n=28$) and so do D³M and KS in the moderate sample sizes ($n=36$). In case of unbalanced sample sizes between the case and control groups, D³M shows preferable results (Supplementary Table S.5).

Table 3. Pathways detected with the proposed method

Cluster	Pathway	-log (P-value)	Genes
Cluster 1	Axonal guidance signaling	5.85	SLIT3, ITGB1, SEMA3G, MYL10, WNT3, FES, NRP2, BMP8A, UNC5B, WNT2B, PIK3R5, GNAI1, ITGA5, KEL, WNT7A, MAG, ADAM12, NTRK1, PRKAG2, WNT4, BMP7, NTN3, PRKCB
	ERK/MAPK signaling	4.09	ITGB1, DUSP9, PLA2G4B, PRKAG2, ITGA5, PIK3R5, RAPGEF4, ESR1, PPP2R5A, KSR1, JMJD7-PLA2G4B, PRKCB
	Caveolar-mediated-Endocytosis signaling	3.72	ITGB1, INS, CD48, ITGA5, ITGB8, JMJD7-PLA2G4B, PRKCB
	Human embryonic stem-Cell pluripotency	3.35	WNT7A, WNT3, BMP8A, SMAD3, NTRK1, WNT2B, PIK3R5, WNT4, BMP7
Cluster 2	TREM1 signaling	3.55	STAT5A, MPO, CASP1, NOD1, CCL3, ITGAX
	Calcium signaling	3.28	GRIN2A, TNNT3, ITPR2, CHRN1, NFATC4, CAMKK2, PPP3CA, CAMKK2, PPP3CA, MEF2B

Table 4. The ranking of the six Wnt gene probes in 351 932 probes

	Wnt2	Wnt2b	Wnt3	Wnt4	Wnt7a	Wnt9a
D ³ M	342 (0.10)	549 (0.16)	591 (0.17)	829 (0.24)	641 (0.18)	887 (0.25)
EB	22 541 (6.40)	12 914 (3.67)	3246 (0.92)	1811 (0.51)	135 (0.04)	20 747 (5.90)
KS	4754 (1.35)	91 240 (25.93)	50 174 (14.26)	41 589 (11.82)	41 066 (11.67)	84 248 (23.94)
DiffVar	56 950 (16.18)	155 621 (44.22)	147 817 (42.00)	146 033 (41.49)	138 928 (39.48)	242 648 (68.95)

For each method, the upper values are the absolute ranking among 351 932 of the genes, and lower are their percentages.

D³M can be flexibly applied to differential methylation problems. Simulation results indicate that D³M can detect not only shape differences but also summary statistics differences as effectively as EB and DiffVar, i.e. natural results from the decomposition (6). This suggests that if we cannot obtain sufficient power using a simple summary statistics approach, we have other options to add shape information. If we would like to simultaneously check the results of variance and shape differences, we remove the mean from the data using $X - E[X]$ in each group at a site before applying D³M. This option is provided with R package D³M; in the function D3M::d3m, the logical parameters of rm.mean and rm.var exists. If we would like to see the variance and shape differences at the same time, we just set the rm.mean = T and rm.var = F.

The statistical test of D³M relies on resampling and requires computational time to calculate *P*-values. However, we could reduce the resampling time using a semi-parametric approach (Knijnenburg *et al.*, 2009).

A current limitation of D³M is that it deals with univariate distributions. In a case of the study of large sample sizes, we can deal with covariates, such as age and gender, but our method currently does not incorporate them into the model, and the user needs to remove the effects of covariates; for example, using residuals with regression analysis, prior to the analysis. The extension of D³M to multivariate distribution relies on the estimation of the Wasserstein metric between the empirical multivariate distributions. The approximation of the Wasserstein metric between the multivariate distributions has been studied recently in (Applegate *et al.*, 2011), and we could derive the null distributions based on that approximation.

D³M does not support a spatial correlation of methylation levels within the CpG island in the current form (Eckhardt *et al.*, 2006; Hansen *et al.*, 2012; Irizarry *et al.*, 2008). The spatial information

will enhance the power of detection for differential methylated sites. This could be accomplished by spatially weighted average of Wasserstein distance over a fixed range of locus. These extensions of D³M will be covered in a future study.

6 Conclusion

In this study, we proposed a novel method, D³M, for detecting differential methylation sites based on distribution-valued data. We showed that distribution shape includes interesting information other than that found using mean- and variance-based methods. A simulation study indicated that D³M is capable to detect various situations.

In the application to the GBM and LGG dataset in the TCGA cohort, we identified 1000 sites with the smallest *q*-values. Most of the sites detected by D³M show strong heterogeneity and tend to be hyper- and hypo-methylated in LGG and GBM, respectively, as found in previous studies.

Since the GBM and LGG dataset contains a large number of significantly different sites, including 279 008, 191 050, 297 493 and 255 317 sites for EB, DiffVar, KS, and D³M respectively; at the FDR ≤ 1%, it is difficult to understand the methylation levels at these sites. In the future, it would be of interest to develop a method that describes the diversity of methylation levels.

Acknowledgements

I am grateful to the reviewers. They suggested many useful comments to improve the manuscript.

Conflict of Interest: none declared.

References

- Aggarwal,C.C. (2013) *Outlier Analysis*. Springer, New York.
- Applegate,D. *et al.* (2011) Unsupervised Clustering of Multidimensional Distributions Using Earth Mover Distance. In: *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 636–644. ACM.
- Baylin,S.B. (2005) DNA methylation and gene silencing in cancer. *Nat. Rev. Clin. Oncol.*, **2**, S4–S11.
- Benjamini,Y. and Hochberg,Y. (1995) Controlling the false discovery rate: a practical and powerful approach to multiple testing. *J. R. Stat. Soc. B*, **57**, 289–300.
- Bock,H.H. and Diday,E. (2000) *Analysis of Symbolic Data*. Springer, Berlin, Heidelberg.
- Diday,E. (1989) Introduction a l'analyse des donnees symboliques. RR-1074, inria-00075485.
- Dong,H. *et al.* (2010) Integrated analysis of mutations, miRNA and mRNA expression in glioblastoma. *BMC Syst. Biol.*, **4**, 163.
- Du,P. *et al.* (2010) Comparison of beta-value and M-value methods for quantifying methylation levels by microarray analysis. *BMC Bioinformatics*, **11**, 587.
- Eckhardt,F. *et al.* (2006) DNA methylation profiling of human chromosomes 6, 20 and 22. *Nat. Genet.*, **38**, 1378–1385.
- Gretton,A. *et al.* (2012) A Kernel two-sample test. *J. Mach. Learn. Res.*, **13**, 723–773.
- Gretton,A. *et al.* (2006) A Kernel Method for the Two-Sample-Problem. *NIPS 2006*.
- Hansen,K.D. *et al.* (2012) BSmooth: from whole genome bisulfite sequencing reads to differentially methylated regions. *Genome Biol.*, **13**, R83.
- Hoelzinger,D.B. *et al.* (2007) Autocrine factors that sustain glioma invasion and paracrine biology in the brain microenvironment. *J. Natl. Cancer Inst.*, **99**, 1583–1593.
- Irizarry,R.A. *et al.* (2008) Comprehensive high-throughput arrays for relative methylation (CHARM). *Genome Res.*, **18**, 780–790.
- Irpino,A. and Verde,R. (2014) Basic statistics for distributional symbolic variables: a new metric-based approach. *Adv. Data Anal. Classif.*, **9**, 143–175.
- Knijnenburg,T.A. *et al.* (2009) Fewer permutations, more accurate *P*-values. *Bioinformatics*, **25**, i161–i168. ISMB (2009).
- Kulis,M. and Esteller,M. (2010) DNA methylation and cancer. *Adv. Genet.*, **70**, 27–56.
- Leek,J.T. *et al.* (2012) The sva package for removing batch effects and other unwanted variation in high-throughput experiments. *Bioinformatics*, **28**, 882–883.
- Liu,T. *et al.* (2013) Transcriptional signaling pathways inversely regulated in Alzheimer's disease and glioblastoma multiform. *Sci. Rep.*, **3**, doi:10.1038/srep03467.
- Müller,A. (1997) Integral probability metrics and their generating classes of functions. *Adv. Appl. Prob.*, **29**, 429–443.
- Phipson,B. and Oshlack,A. (2014) DiffVar: a new method for detecting differential variability with application to methylation in cancer and aging. *Genome Biol.*, **15**, 465.
- Polisetty,R.V. *et al.* (2012) LC-MS/MS analysis of differentially expressed glioblastoma membrane proteome reveals altered calcium signalling and other protein groups of regulatory functions. *Mol. Cell Proteomics*, **11**, M111.013565.
- Rampazzo,E. *et al.* (2013) Wnt activation promotes neuronal differentiation of Glioblastoma. *Cell Death Dis.*, **4**, 500e.
- Ramsay,J.O. and Silverman,B.W. (2005) *Functional Data Analysis*. 2nd ed. Springer-Verlag, New York.
- Ritch,P.A. *et al.* (2003) Neuregulin-1 enhances motility and migration of human astrocytic glioma cells. *J. Biol. Chem.*, **278**, 20971–20978.
- Ritchie,M.E. *et al.* (2015) limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.*, **43**, e47.
- Rueshendoff,L. (2011) Wasserstein metric. *Encyclopedia of Mathematics*. <http://www.encyclopediaofmath.org/index.php?title=About&coldid=20414>.
- Stacklies,W. *et al.* (2007) pcaMethods – a Bioconductor package providing PCA methods for incomplete data. *Bioinformatics*, **23**, 1164–1167.
- The Cancer Genome Atlas Research Network. (2008) Comprehensive genomic characterization defines human glioblastoma genes and core pathways. *Nature*, **455**, 1061–1068.
- Wang,H. and Marron,J.S. (2007) Object oriented data analysis: sets of trees. *Ann. Stat.*, **35**, 1849–1873.
- Wang,D. *et al.* (2012) IMA: an R package for high-throughput analysis of Illumina's 450K Infinium methylation data. *Bioinformatics*, **28**, 729–730.