



Smooth modeling of covariate effects in bisulfite sequencing-derived measures of DNA methylation

Kaiqiong Zhao

McGill University

Supervisors: Celia Greenwood & Karim Oualkacha

BIRS workshop

New Statistical Methods for Family-Based Sequencing Studies

Banff Centre – August 5 - 10, 2018

Epigenetics and DNA Methylation



- ▶ change gene expression without changing DNA sequence
- ▶ can be altered by age, diet, stress and environmental exposures



- ▶ **interest:** identify **genomic regions** where DNA methylation patterns demonstrate alterations associated with disease phenotypes (DMRs).

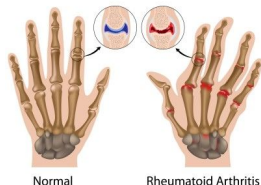
Motivating Dataset

Methylation profiles of Rheumatoid Arthritis (RA) patients and controls
(provided courtesy of Dr. Marie Hudson, McGill)



- ▶ **Samples:** cell-separated blood samples of 22 RA patients and 21 healthy individuals from either T cells or monocytes.

	MONO	TCELL
RA	10	12
Control	8	13



- ▶ **Methylation:** targeted custom captured bisulfite sequencing[†]
 - Prior selection of predefined genomic regions of interest
 - ~ 400,000 regions in the genome
- ▶ This presentation focuses on **one targeted region**
 - Chr4: 102,711,629 – 102,712,832** (near gene *BANK1*)
 - consists of 123 CpG sites

[†] Allum et al (2015) *Nature*.



Position	Unmeth counts	Meth counts	Total counts	Sample ID	Sample-level covariates
102711629	2	2	4	1	...
102711630	15	0	15	1	
102711649	15	0	15	1	
102711650	8	0	8	1	
102711850	15	0	15	2	
102711851	4	9	13	2	

Sample-level covariates: disease status, cell type composition, age, smoking status...



► Read-depth variability

- The total number of reads varies at different CpG sites. Modeling the proportion treats noisy measurements the same way as accurate ones ($\frac{5}{10} = \frac{50}{100}$)
- missing values occur frequently

► Variability in cell-type mixture proportions

→ **adjusting for multiple covariates**

- Methylation levels vary substantially across different cell types, which can confound the association of interest.

► Experimental errors

- Sequencing errors: more mis-alignment of unmethylated reads after bisulfite sequencing
- Bisulfite conversion errors: incomplete C-T conversion; or over-treatment with bisulfite leading to conversion of methylated C to T



Challenges	BSmooth (Hansen et al., 2012)	BiSeq (Hebestreit et al., 2013)	SMSC (Lakhal-Chaieb et.al., 2017)
Variable read-depth	✓		✓
Experimental error			✓
Mixture of cell types		✓	
Multiple covariates		✓	

- ▶ Most of the existing methods are of **two-stage** nature
 - (1) smooth the methylated proportions **for each sample**, and
 - (2) fit model (t-test or beta regression) to the **smoothed methylation data**.
- ▶ **Motivation**: extend the work in Lakhal-Chaieb et.al., 2017 to enable an integrated analysis of **multiple samples** that allows for
 - cell type mixtures**, and
 - multiple **covariates** in the model.

Notation & Model



Let (i, j, k) index sample, CpG sites and reads respectively.

$$i = 1, 2, \dots, l; \quad j = 1, 2, \dots, n_i; \quad k = 1, 2, \dots, X_{ij}, \quad N = \sum_{i=1}^l n_i.$$

Example: CpG site j for Individual i

(i, j)	t_{ij} Position	Y_{ij} Meth	X_{ij} Total	Sample ID	$Z_{1i}, Z_{2i}, \dots, Z_{Pi}$ covariates
(1, 1)	114354051	2	4	1	...
(1, 2)	114354052	0	15	1	
...	114354053	0	15	1	
(1, n_1)	114354054	0	14	1	
(2, 1)	114354056	0	15	2	
(2, 2)	114354057	9	13	2	

	Observed methylation	True methylation
Read 1	M $Y_{ij1} = 1$	S_{ij1}
Read 2	U $Y_{ij2} = 0$	S_{ij2}
Read 3	U $Y_{ij3} = 0$	S_{ij3}
Read 4	M $Y_{ij4} = 1$	S_{ij4}
	$Y_{ij} = 2$	$S_{ij} = \sum_k S_{ijk}$

(p₀ p₁)



- ▶ Assume **known error rates** ρ_0 and ρ_1 ,

$$\rho_0 = \mathbb{P}(Y_{ijk} = 1 \mid S_{ijk} = 0)$$

$$\rho_1 = \mathbb{P}(Y_{ijk} = 1 \mid S_{ijk} = 1).$$

(Lakhal-Chaieb et.al., 2017)

- ▶ We specify model

$$S_{ij} \mid \mathbf{Z}_i, X_{ij} \sim \text{Binomial}(X_{ij}, \pi_{ij})$$
$$\theta_{ij} = \log \left\{ \frac{\pi_{ij}}{1 - \pi_{ij}} \right\} = \beta_0(t_{ij}) + \beta_1(t_{ij})Z_{1i} + \beta_2(t_{ij})Z_{2i} + \dots + \beta_P(t_{ij})Z_{Pi}.$$

- ▶ Use splines to parametrize smooth covariate effects:

$$\beta_p(t_{ij}) = \sum_{l=1}^L \alpha_{pl} B_l(t_{ij}) \text{ for } p = 0, 1, \dots, P.$$

- ▶ Smoothing parameters $\{\lambda_0, \lambda_1, \dots, \lambda_P\}$ for controlling the smoothness of $\beta_p(\mathbf{t})$

$$\mathcal{L}^{\text{Penalization}} = \sum_{p=0}^P \lambda_p \int (\beta_p''(t))^2 dt = \sum_{p=0}^P \lambda_p \alpha_p^T \mathbf{A} \alpha_p$$



► Penalized EM algorithm

Initialization: $\alpha^* \lambda^*$;

repeat

1. E-step: calculate the conditional expected outcomes $\mathbb{E}(S_{ij} | Y_{ijk})$
2. M-step: $(\hat{\alpha}, \hat{\lambda}) = \arg \max Q(\alpha, \lambda | \alpha^*)$. (Q is the binomial likelihood replacing S_{ij} with its expectations)

repeat

- 2.1 P-IRLS iteration.
- 2.2 Smoothing parameters estimated by REML[‡].

until estimates converge;

$\hat{\alpha} \rightarrow \alpha^*$ and $\hat{\lambda} \rightarrow \lambda^*$

until estimates converge;

⇒ Estimates of the smooth functions of covariates effects $\widehat{\beta}_1(\mathbf{t}), \widehat{\beta}_2(\mathbf{t}) \dots \widehat{\beta}_p(\mathbf{t})$.

► Inference of smooth covariate effects taking account of the uncertainty in both E step and M step

⇒ Variance of the smooth estimates $\widehat{\text{Var}}(\widehat{\beta}_p(\mathbf{t}))$

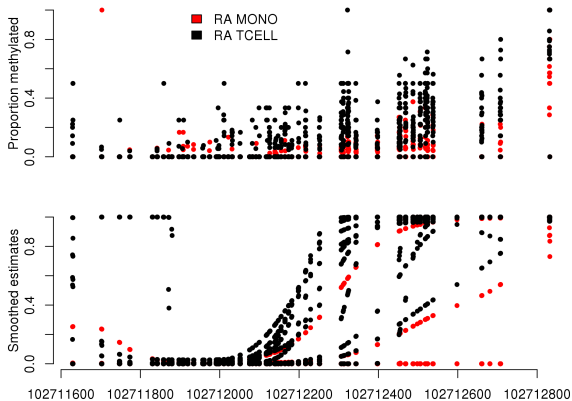
⇒ A region-wise test statistic (p-value) for $H_0 : \beta_p(\mathbf{t}) = \mathbf{0}$
 $\widehat{\alpha}_p \{ \widehat{\text{Var}}(\widehat{\alpha}_p) \}^{-1} \widehat{\alpha}_p^T \sim \chi_{edf}^2$ where
 $edf = \text{trace}(2\mathbf{H} - \mathbf{H}\mathbf{H}^T)$

$$\mathcal{H}(\alpha) = \left\{ \frac{\partial^2 Q(\alpha | \alpha^*)}{\partial \alpha \partial \alpha^T} + \frac{\partial^2 Q(\alpha | \alpha^*)}{\partial \alpha \partial \alpha^{*T}} \right\} \Bigg|_{\alpha^* = \alpha} \quad \ddagger$$

[‡] Oakes, D. (1999) Direct calculation of the information matrix via the EM. JRSSB

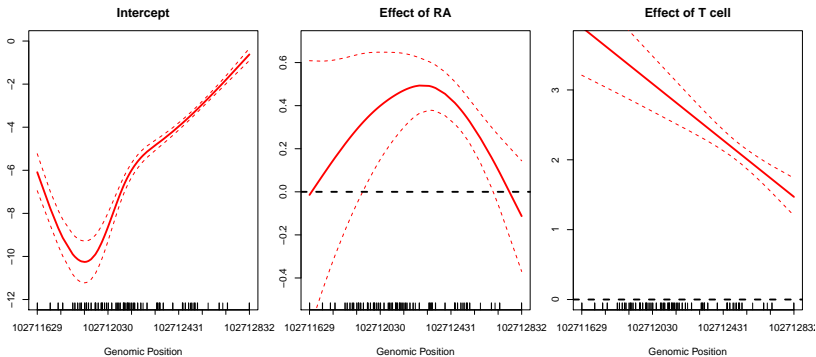
Data example

Raw data & per-sample smoothed estimates



Genomic position (in bp)

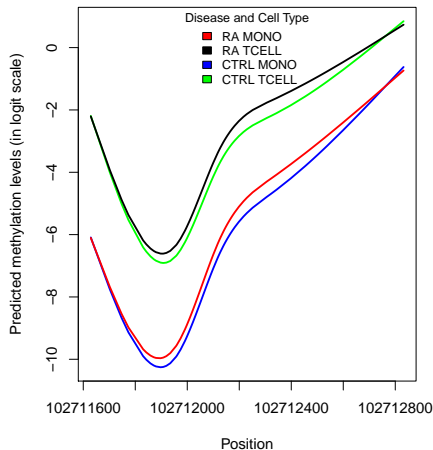
Inference of the smooth covariate effects



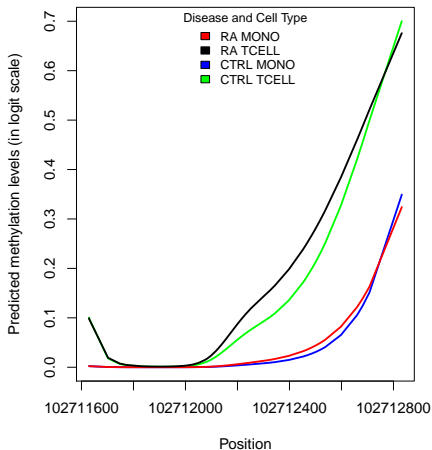
Predicted methylation levels



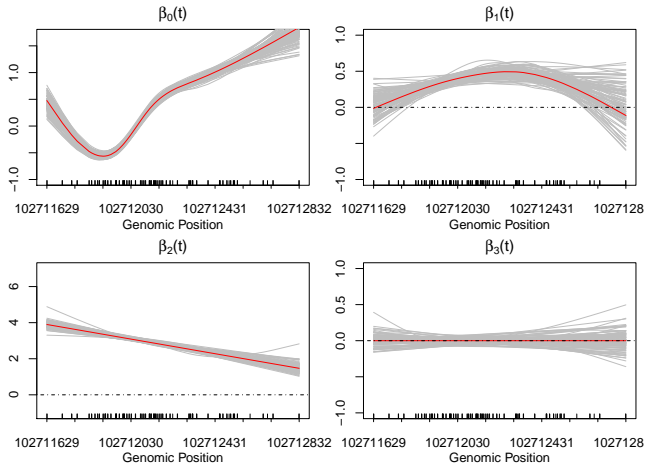
Logit scale



Proportion scale

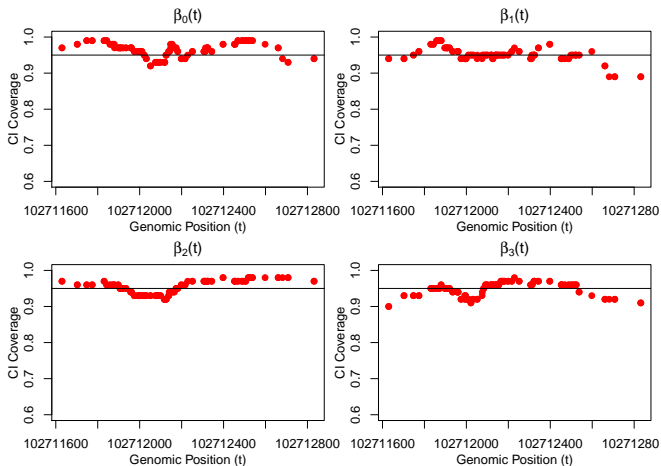


- $\beta_3(\mathbf{t}) = 0$ and error parameters $\rho_0 = 0.003$ and $1 - \rho_1 = 0.1$ [‡]

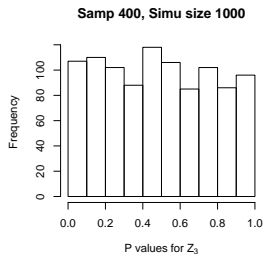
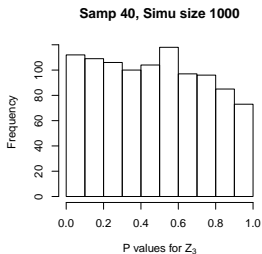
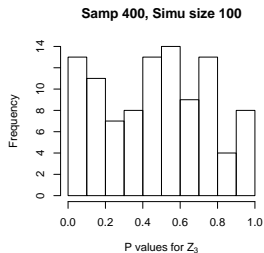
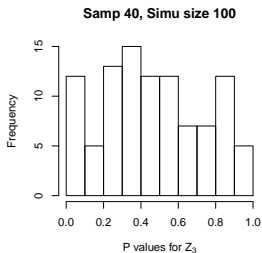


[‡] Prochenka et al. (2015) *Bioinformatics*.

Coverage probability of confidence intervals



Region-based p-values under null





- ▶ a **model** characterizing bisulfite sequencing data from multiple samples, which copes naturally with variable read depth, experimental errors and test samples with a mixture of cell types.
- ▶ a **smoothed EM method** to **make inference about smooth exposure/covariate effects**.
- ▶ the method is shown to be capable in **capturing the major underlying patterns** in the data.
- ▶ **Next step plans**
 - estimating error rate p_0 by calculating C-T conversion rate at non-CpG Cs
 - adding subject-specific random effects to account for the heterogeneity in methylation profiles that cannot be explained by covariates Z_i s in the model
 - correlated samples

Acknowledgement



Dr. Celia Greenwood (McGill)
Dr. Karim Oualkacha (UQAM)
Dr. Lajmi Lakhal-Chaieb (Laval)
Dr. Kathleen Klein (JGH)
Dr. Marie Hudson (McGill)



**Fonds de recherche
Santé**

Québec 



CIHR IRSC



Canadian Institutes of Health Research | Instituts de recherche en santé du Canada



Thanks

Questions & Comments