

University of Memphis

University of Memphis Digital Commons

Electronic Theses and Dissertations

7-21-2017

Epigenetic Marker Identification and Assessment of Methods on Cell Type Compositions at the Epigenome-Scale

Akhilesh Kaushal

Follow this and additional works at: <https://digitalcommons.memphis.edu/etd>

Recommended Citation

Kaushal, Akhilesh, "Epigenetic Marker Identification and Assessment of Methods on Cell Type Compositions at the Epigenome-Scale" (2017). *Electronic Theses and Dissertations*. 1709.
<https://digitalcommons.memphis.edu/etd/1709>

This Dissertation is brought to you for free and open access by University of Memphis Digital Commons. It has been accepted for inclusion in Electronic Theses and Dissertations by an authorized administrator of University of Memphis Digital Commons. For more information, please contact khggerty@memphis.edu.

EPIGENETIC MARKER IDENTIFICATION AND ASSESSMENT OF METHODS ON
CELL TYPE COMPOSITIONS AT THE EPIGENOME-SCALE

by

Akhilesh Kaushal

A Dissertation

Submitted in Partial Fulfillment of the

Requirements for the Degree of

Doctor of Philosophy

Major: Epidemiology

The University of Memphis

August 2017

Preface

Chapter 2 of this dissertation has been published as **Kaushal, A., Zhang, H., Karmaus, W. J., Everson, T. M., Marsit, C. J., Karagas, M. R., ... & Wang, S. L. (2017). Genome-wide DNA methylation at birth in relation to in utero arsenic exposure and the associated health in later life. *Environmental Health*, 16(1), 50.** I performed all the statistical analyses and drafted the manuscript along-with Zhang, H. Wang, S.L. conceived the study and collected all the data, Karmaus, W.J., provided guidance on epigenome and clinical aspects and T. M., Marsit, C. J., Karagas, M. R performed the replication study.

Chapter 4 of this dissertation has been published as **Kaushal, A., Zhang, H., Karmaus, W. J., Ray, M., Torres, M. A., Smith, A. K., & Wang, S. L. (2017). Comparison of different cell type correction methods for genome-scale epigenetics studies. *BMC bioinformatics*, 18(1), 216.** I performed all the analyses and drafted the manuscript along with Zhang, H. Karmaus, W. J. motivated the analyses and contributed to the manuscript, Ray, M. provided code for simulation scenario one and edited the manuscript . All authors were involved in editing and revising the manuscript.

Acknowledgements

I would like to express my deep appreciation and gratitude to my advisor, Dr. Hongmei Zhang, for providing me with the opportunity to perform doctoral research in her supervision. My graduate training under her guidance has been a great learning experience for me. She always encouraged me to develop independent thinking and research skills. Her guidance, mentorship, and support throughout this research journey have been invaluable.

I would like to express sincere thanks to my dissertation committee members, Drs. Wilfried JJ Karmaus and Shu-Li Wang for their valuable suggestions and their precious time.

I am also thankful to other members in my committee: Drs. Prateek Banerjee, and Meredith Ray for all their support and help from the start to the end of my graduate study journey in the School of Public health at University of Memphis.

I acknowledge all the funding sources that have sponsored my graduate studies (NSC98-2314-B-400-001-MY3, MOST103-2314-B-400-006, R01 AI091905 (PI: Wilfried Karmaus), R21AI099367 (PI: Hongmei Zhang), and 1R01AI121226 (PIs: Hongmei Zhang, John Holloway).

I am forever indebted to my parents for everything.

Abstract

Kaushal, Akhilesh. Ph.D. The University of Memphis. August 2017. Epigenetic Marker Identification and Assessment of Methods on Cell Type Compositions at the Epigenome-Scale. Major Professor: Dr. Hongmei Zhang.

Epigenetics is the study of heritable changes in genes which are caused by chemical compounds derived from natural and man-made sources. DNA methylation an epigenetic phenomenon, is most vulnerable to environmental factors during embryogenesis, which is a period of rapid cell division and epigenetic remodeling. Given the recent increase in the incidence of childhood diseases, it is crucial to understand the role of environmental factor through epigenetic study in causing adverse health effects. This dissertation revolves around three major hypotheses. In the first hypothesis we evaluated the association between in utero arsenic exposure and genome-wide DNA methylation in cord blood from the birth cohort data of Taiwan. The identified CpG sites were replicated in an independent birth cohort (New Hampshire birth cohort study; NHBCS) and further assessed longitudinal associations of DNA methylation with disease biomarkers measured at later ages in our cohort from Taiwan. In the second hypothesis we assessed the association between Immunoglobulin E (IgE) production and DNA methylation at birth via cord blood in a longitudinal study. The study was conducted from the birth cohort data of Taiwan and the findings were replicated in an independent birth cohort (Isle of Wight; IoW), and further the stability of identified CpG sites was assessed based on intra-class (ICC) correlation measure. In the third hypothesis we assessed the confounding effect in epigenome wide association study due to underlying cell composition and evaluated several methods and algorithms proposed to adjust for this confounding effect.

List of Abbreviations

CpG: 5'-cytosine-phosphate-guanine-3'; CpGs: Multiple CpG

LDL: Low density lipoprotein

tAs: Total arsenic obtained by adding inorganic arsenic (iAs), mono-methylated arsenic (MMA), di-methyl arsenic (DMA)

coeff: Coefficient; coeff.m: coefficient for main effect; coeff.int: coefficient for interaction effect.

DNA: deoxyribonucleic acid

DNA-M: DNA methylation

TSS: Transcription start site; TSS1500: within 1500 base pairs of a TSS; TSS200: within 200 base pairs of TSS.

FDR: False discovery rate

IoW: Isle of Wight

IgE: Immunoglobulin E

DAVID: Database for Annotation, Visualization and Integrated Discovery

KEGG: Kyoto Encyclopedia of Genes and Genomes

NHBCS: New Hampshire Birth Cohort Study

SVA: Surrogate variable analysis

GO: Gene ontology

ROS: Reactive oxygen species

DMRs: Differentially methylated regions

SNPs: Single nucleotide polymorphism; plural

Table of Contents

Section	Page
Abstract	iv
List of Abbreviations	v
List of Tables	x
List of Figures	xi
1 Introduction	1
1.1 Epigenetics.....	2
1.2 Arsenic.....	3
1.3 Immunoglobulin E (IgE).....	4
1.4 Underlying Cell compositions in cord blood.....	4
1.5 Contribution.....	5
1.6 Organization of the dissertation.....	6
2 Epigenetic marker identification at birth associated with in utero arsenic exposure	7
2.1 Abstract	7
2.2 Background	8
2.3 Methods	10
2.3.1 Data collection and pre-processing of birth cohort data from Taiwan.....	10
2.3.2 Data collection and pre-processing of birth cohort data from NHBCS.....	13
2.3.3 Correction for cell mixture proportion.....	15
2.3.4 DAVID.....	15
2.3.5 GeneMANIA.....	16
2.3.6 Statistical analyses.....	16
2.3.6.1 Statistical analyses in NHBCS.....	16
2.4 Results	18
2.5 Discussion	30
2.6 Conclusion	33
3 Epigenetics markers at birth longitudinally associated with Immunoglobulin E (IgE)	34

3.1 Abstract	34
3.2 Background	35
3.3 Methods	36
3.3.1 Data collection and pre-processing of birth cohort data from Taiwan.....	36
3.3.2 Data collection and pre-processing of birth cohort data from IoW.....	37
3.3.3 Functional annotation and pathway analysis.....	39
3.3.4 Statistical analyses	39
3.3.4.1 Statistical Analyses in IoW.....	40
3.4 Results	40
3.5 Discussion	45
3.6 Conclusion	46
4 Assessment of methods for cell type correction in epigenome wide association study	47
4.1 Abstract	47
4.2 Background	48
4.3 Method	50
4.3.1 Reference-based methods.....	51
4.3.2 Reference-free methods.....	52
4.3.3 Three real data sets used to compare the approaches.....	53
4.3.4 Simulated data sets to compare the methods.....	54
4.3.5 Statistical analyses.....	55
4.4 Results	57
4.3.1 Findings from prenatal arsenic exposure and DNA-methylation data.....	57
4.3.2 Findings from example data.....	62
4.3.3 Findings from breast cancer status and DNA-methylation data.....	64
4.3.4 Findings from simulated data.....	64
4.5 Discussion	70
4.6 Conclusion	72
5 Summary	73

Appendices	82
Table A1.1	82
Table A1.2	89
Table A2.1	91
Table A2.2	97
Figure A1.1	99
Figure A1.2	100
Figure A3.1	101
Figure A3.2	102
Figure A3.3	103
R_code_thesis	104
SAS code for epigenome wide mixed modeling	110
IRB Approval	115

List of Tables

2.1 Characteristics of mothers and their newborns by newborn sex in Taiwan during 2000-2001 (n=64)	19
2.2 Distribution of creatinine-adjusted concentrations of urinary arsenic species (iAs, MMA, and DMA) (n=64)	20
2.3 KEGG pathways identified using DAVID that are more specific to arsenic exposure based on data from n=64 pregnant women from the maternal infant cohort in Taiwan	27
3.1 Distribution of IgE across different time points	41
4.1 Number of significant CpG sites with and without cell type correction and overlap with the SVA method (data on prenatal arsenic exposure and DNA methylation)	61
4.2 Number of significant CpG sites with and without cell-correction methods and overlap of CpG sites with those from the SVA method (example data from FaST-LMM-EWASher package)	63
4.3 Summary of sensitivity, specificity of FaST-LMM-EWASher, RefFreeEWAS, RefFreeCellMix, ReFACTor, and SVA for 100 simulated data across three settings	68

List of Figures

2.1 Subject recruitment and preprocessing of DNA methylation data in Taiwanese birth cohort	21
2.2 The flow of analyses performed in the study	22
2.3 Manhattan plot for Genome-wide DNA methylation associated with creatinine adjusted urinary arsenic concentration	24
2.4 Association of arsenic exposure with the DNA methylation based on M-values of the 58 CpG sites mapped to 56 genes	25
2.5 Heatmap of the correlations between cord blood DNA methylation and LDL across different ages (2, 5, 8, 11, 14 years)	29
3.1 Manhattan plot for the longitudinal association of Genome-wide DNA methylation with log ₁₀ Immunoglobulin E (IgE)	42
3.2 Flow of analysis	43
3.3 Longitudinal association of the residual of DNA methylation with log ₁₀ Immunoglobulin E (IgE) of the 124 CpG sites mapped to 89 genes	44
4.1 Venn diagram illustrating the overlap of identified CpG sites that are associated with prenatal arsenic exposure at FDR level of 0.05 after incorporating estimated cell type compositions by different methods for the association study of prenatal arsenic exposure with DNA-methylation	59
4.2 Venn diagram illustrating the overlap of identified CpG sites that are associated with cancer status at FDR level of 0.05 after incorporating estimated cell type compositions by different methods for the association study of cancer status with DNA-methylation	60
4.3 Plots of sensitivity v.s. 1-specificity and estimated ROC curves	67

1 Introduction

The dissertation focuses on the identification of epigenetic markers from cord blood DNA methylation at birth associated with in utero arsenic exposure and immunoglobulin E (IgE).

There is increasing evidence that in-utero arsenic exposure causes adverse health effects later in life [1, 2]. Arsenic is a potent human toxicant and carcinogen, but knowledge of the mechanism through which it exerts long term adverse health effects is limited. Inorganic arsenic and its methylated metabolites can easily cross the placenta and thus producing arsenic concentrations in cord blood similar to maternal blood [3]. The study of epigenetic changes such as DNA-methylation alterations that can affect gene activity may provide insight into mechanism through which arsenic exerts its adverse effects [4].

Immunoglobulin E (IgE) is known to play a major role in many of the allergic diseases such as asthma, atopic dermatitis (eczema) and hay fever. IgE production leads to type I hypersensitivity, which manifests various allergic diseases. However, the mechanism underlying IgE production is poorly understood. There is evidence that DNA methylation is associated with total IgE. An epigenome wide association study has a potential to shed more lights on the production of IgE.

Epigenome wide association studies are known to be influenced by the underlying cell compositions. Several algorithms and methods have been proposed to estimate or adjust for these underlying cell compositions. However, it is unknown which of the methods and algorithms works best.

1.1 Epigenetics

Epigenetics is the study of heritable changes in genes, which are caused by chemical compounds derived from natural and man-made sources. These heritable changes regulate genome by: (a) methylating DNA in genome, (b) modifying the histone, a protein that enables DNA to form long molecules. The chemical compounds that cause heritable changes in genome are known as epigenome. Much of the epigenome is reset when parents pass their genomes to their offspring; however, under some circumstances, some of the chemical tags on the DNA and histones of eggs and sperm may be passed on to the next generation [5-7]. When cells divide, often much of the epigenome is passed on to the next generation of cells, helping the cells remain specialized [7].

DNA methylation is an epigenetic phenomenon wherein the nucleotides in DNA is modified by addition of methyl group through covalent bond. In DNA methylation a methyl group is attached to the 5th atom in the 6-atom ring of cytosine leading to 5 methyl cytosine (5mC) or at 6th position of adenine ring leading to 6-methyl adenine (6mA). The term CpG refers to the base cytosine (C) linked by a phosphate bond to the base guanine (G) in the DNA nucleotide sequence. In the human genome, it predominantly occurs at cytosine–guanine dinucleotide (CpG) sites, and serves to regulate gene expression and maintain genome stability DNA [8, 9]. Methylation is heritable and stable from one cell to another during cell division and thus leads to formation of epigenetic memory [10].

The epigenome is most vulnerable to environmental factors during embryogenesis, which is a period of rapid cell division and epigenetic remodeling [11-13]. Given the recent increase in the incidence of childhood immune-based diseases, it is crucial to understand the role of environmental stressors [14]. A stressor is a chemical or biological agent, environmental condition, external stimulus or an event that causes stress to an organism. Environmental studies

have shown that DNA methylation could be altered under environmental stress, by overall genome-wide reduction in DNA methylation content (global hypo methylation). This alteration in DNA methylation can alter the expression of underlying gene.

Arsenic (As) is one such environmental stressor whose exposure is known to alter DNA methylation both globally and in the promoter regions of certain genes [15, 16]. Upon entering the human body, inorganic As is methylated for detoxification. This detoxification process uses S-adenosyl methionine (SAM), which is a universal methyl donor for methyltransferases including DNA methyltransferases (DNMTs) that determine DNA methylation.

1.2 Arsenic

Arsenic is a metalloid found in numerous minerals usually in combination with sulfur and other metals. In reducing and oxygenated conditions, arsenite (AsIII), and arsenate (AsV), are the main oxidation states, respectively. Compounds of arsenic are divided into three major groups.

- a) Inorganic arsenic compounds (arsenic trioxide, sodium arsenite, arsenic trichloride, arsenic acid, and arsenic pentoxide),
- b) Organic arsenic compounds (arsanilic acid, methylarsonic acid, dimethylarsinic acid, arsenobetaine and arsenosugars), and
- c) Arsine gas.

Arsenite and arsenate are the most common inorganic forms in water. Arsenic is mainly transported in the environment via water from both natural and anthropogenic sources. In some regions of the world, groundwater (used for drinking water) is naturally contaminated with arsenic due to arsenic rich geological formations. These areas include Bangladesh, China, Taiwan, West Bengal (India), and some parts of Argentina, Chile, Mexico, Vietnam, Australia, and the USA. In unaffected areas, the levels of arsenic are only a few micrograms per liter of

ground water, whereas, in affected areas, the levels may range from tens to thousands of micrograms per liter [17].

Ingestion of arsenic contaminated water is the primary route of inorganic arsenic exposure for the general population. Approximately hundred million individuals world-wide are at risk of elevated arsenic exposure, mainly via drinking water.

Most of the ingested inorganic arsenic is absorbed in the gastrointestinal tract and then reduced in the blood. In humans, inorganic arsenic is metabolized through the conversion of AsV to AsIII, followed by methylation to monomethylated and dimethylated arsenicals (MMA and DMA, respectively) [18].

1.3 Immunoglobulin E (IgE)

Immunoglobulin E (IgE) is one of five isotypes of human immunoglobulins and is produced by plasma cells. Immunoglobulin E (IgE) is known to play a major role in many of the allergic diseases such as asthma, atopic dermatitis (eczema) and hay fever. IgE production leads to type I hypersensitivity, which manifests various allergic diseases. Thus understanding the mechanism leading to the IgE production is a key to understanding the pathophysiology of various allergic disease.

1.4 Underlying cell compositions in cord blood

Cord blood is the blood that remains in the vein of the umbilical cord and placenta at the time of birth. Umbilical cord blood consists of various cell types such as nucleated red blood cells, granulocytes, monocytes, natural killer cells, B cells, CD4⁺ T cells, and CD8⁺ T cells. Thus, DNA methylation measured in cord blood represents weighted averages of these cell-type specific methylation levels, with weights corresponding to the proportion of the different cell types in a cord blood sample. Thus, epigenome wide association study assessing the association

between cord blood and an exposure of interest could be confounded by cellular heterogeneity. Identifying and sub setting each cell types is not practical in larger epidemiological studies. Thus, several algorithms have been developed to measure and adjust for cellular heterogeneity in whole blood.

1.5 Contribution

Epigenome wide association study can identify epigenetic markers that will help reveal the adverse developmental effect of in utero arsenic. Also the longitudinal study predicting the production of IgE associated with cord blood DNA methylation can provide useful insight into the developmental origin of immunity based disease. The main contributions of the work presented in this dissertation can be summarized as below

1. Identification of epigenetic markers at birth associated with in utero arsenic exposure.
2. Identification of epigenetic markers at birth linking in utero arsenic exposure to cardiovascular disease.
3. Identification of epigenetic markers at birth predicting the production of immunoglobulin E (IgE) at later ages.
4. Best method to adjust for the underlying cell compositions for epigenome wide association study.

1.6 Organization of the dissertation

The dissertation is organized as self-explanatory chapters related to the epigenetic marker identification. Chapter 2 presents the role of in utero arsenic exposure on fetal developmental programming and its adverse influence in later life. Chapter 3 discusses the role of DNA methylation at birth in predicting the IgE production at later ages. Chapter 4 compares several methods and algorithms to adjust for underlying cell compositions in epigenome wide association study. Finally, in Chapter 5, I summarize the work presented here.

2 Epigenetic marker identification at birth associated with in utero arsenic exposure

2.1 Abstract

Background

In utero arsenic exposure may alter fetal developmental programming by altering DNA methylation, which may result in a higher risk of disease in later life. We evaluated the association between in utero arsenic exposure and DNA methylation (DNAm) in cord blood and its influence in later life.

Methods

Genome-wide DNA methylation in cord blood from 64 subjects in the Taiwanese maternal infant and birth cohort was analyzed. Robust regressions were applied to assess the association of DNA methylation with in utero arsenic exposure. Multiple testing was adjusted by controlling false discovery rate (FDR) of 0.05. The DAVID bioinformatics tool was implemented for functional annotation analyses on the detected CpGs. The identified CpGs were further tested in an independent cohort. For the CpGs replicated in the independent cohort, linear mixed models were applied to assess the association of DNA methylation with low-density lipoprotein (LDL) at different ages (2, 5, 8, 11 and 14 years).

Results

In total, 579 out of 385,183 CpGs were identified after adjusting for multiple testing (FDR=0.05), of which ~60% were positively associated with arsenic exposure. Functional annotation analysis on these CpGs detected 17 KEGG pathways (FDR=0.05) including pathways for cardiovascular diseases (CVD) and diabetes mellitus. In the independent cohort, about 46% (252 out of 553 CpGs) of the identified CpGs showed associations consistent with those in the study cohort. In total, 12 CpGs replicated in the independent cohort were in the pathways related

to CVD and diabetes mellitus. Via longitudinal analyses, we found at 5 out of the 12 CpGs methylation was associated with LDL over time and interactions between DNA methylation and time were observed at 4 of the 5 CpGs, cg25189764 (coeff=0.157, p-value=0.047), cg04986899 (coeff. for interaction [coeff.int]=0.030, p-value=0.024), cg04903360 (coeff.int=0.026, p-value=0.032), cg08198265 (coeff.int = -0.063, p-value=0.0021), cg10473311 (coeff.int = -0.021, p-value=0.027).

Conclusion

In utero arsenic exposure was associated with cord blood DNA methylation at various CpGs. The identified CpGs may help determine pathological epigenetic mechanisms linked to in utero arsenic exposure. Five CpGs (cg25189764, cg04986899, cg04903360, cg08198265 and cg10473311) may serve as epigenetic markers for changes in LDL later in life.

2.2 Background

Arsenic, a widespread element in the environment, poses a serious threat to human health. Millions of people around the globe are exposed to arsenic from drinking water that exceeds the safe limit of 10 ppb as recommended by World Health Organizations [19]. Arsenic is known to easily pass through the placenta in humans and other mammals, producing arsenic concentrations in cord blood similar to maternal blood [3]. Epidemiological studies have reported that gestational arsenic exposure is associated with increased risk of non-cancerous and cancerous diseases in adulthood [20, 21]. For instance, a number of studies have shown that early life arsenic exposure is associated with later cardiovascular diseases (CVDs) [22-24]. In animal studies, in utero exposure to low level arsenic in the womb and in adulthood was found to be associated with diabetes mellitus [25].

The mechanisms through which in utero exposure to arsenic may result in a higher risk of various diseases are not well understood. However, harmful effects such as the generation of reactive oxygen species (ROS), which causes oxidative DNA damage, binding and inhibition of arsenic metabolites to enzymes, and perturbation of key signaling pathways, are thought to play certain roles in disease development [26]. In addition, clinical and epidemiological studies have observed that environmental exposure in early life can affect the risk of disease later in life through a phenomenon known as developmental programming [21, 27, 28]. The study of epigenetic changes such as DNA-methylation alterations that can affect gene activity may provide insight into developmental programming [4].

Studies found that chronic arsenic exposure in adults is associated with increased DNA methylation extracted from whole blood leukocytes [29, 30]. Experimental studies in animals have also shown that intra-uterine exposure to arsenic alters DNA methylation in offspring [31]. Some studies examined the association of genome-wide DNA methylation in cord blood with in utero arsenic exposure. Most of them were based on cohorts established in the United States and Bangladesh. These studies did not identify any statistically significant CpGs at the whole epigenome level and thus focused on the top 100 [32] or 500 CpG sites [16] potentially associated with in utero arsenic exposure, while the study by Kile et al [33] investigated the association of CpG sites in *p16*, *p53*, LINE-1 and Alu repetitive elements.

Our study, based on data from a prospective birth cohort study established in Taiwan, aimed to comprehensively assess genome-wide DNA methylation in cord blood in association with in utero arsenic exposures (using maternal urinary arsenic concentrations), identify CpG sites showing such statistically significant associations after adjusting for multiple testing by controlling false discovery rate (FDR), and examine possible pathways of genes involving the

identified CpGs. Additionally, we attempted to replicate our finding in an independent birth cohort (New Hampshire birth cohort study; NHBCS) and further assessed longitudinal associations of DNA methylation with disease biomarkers measured at later ages in our cohort from Taiwan. The findings will contribute to an improved understanding of the adverse mechanisms of in utero arsenic exposure on genome-wide epigenetic variation and whether epigenetic markers in cord blood can influence children's diseases risk later in life.

2.3 Methods

2.3.1 Data collection and pre-processing of birth cohort data from Taiwan

Taiwanese maternal infant and birth cohort description

The data resulted from the Maternal and Infant Cohort Study in Taiwan investigating various in utero and postnatal factors considered to affect child health outcomes [21]. All pregnant women participating in this study signed informed consent forms explaining the benefits and risks of participation. This study was approved by Human Ethical Committee of the National Health Research Institutes in Taiwan. Pregnant women who received medical care at a local medical center were invited to join this study between December 2000 and November 2001. At the beginning, 430 of 610 pregnant women volunteered to participate, on average at 8 weeks gestation. Of the 430 pregnant women, 127 were excluded due to non-compliance of providing samples. Thus, urine samples were obtained from 313 pregnant women during the third trimester (28-38 weeks of gestation). Five newborns could not be included due to loss to follow up. Of all mother-newborn pairs, 299 have cord blood samples collected and DNA methylation data was available for 64 cord blood samples with sufficient amount of good quality DNA for this epigenome assay.

Assessment of arsenic exposure

Participants provided a spot urine sample at the time of enrollment in this study (at eight weeks of gestation). Urine was frozen at -20 °C in a 10-ml polypropylene tube. Arsenite (As^{III}), arsenate (As^{V}), monomethylarsonic acid (MMA), and dimethylarsinic acid (DMA) were quantified by high-performance liquid chromatography/inductively coupled plasma mass spectrometry (HPLC-ICP-MS) and anion exchange columns (Hamilton PRP X-100 [10 mm particle size, 250 mm \times 6.1 mm]). Total urinary arsenic (TUA) was calculated by adding iAs (As^{III} + As^{V}) + MMA + DMA. The limitations of detection (LOD) for the various species were 0.09 mg/L for As^{III} , 0.05 mg/L for As^{V} , 0.05 mg/L for MMA, and 0.04 mg/L for DMA. Creatinine was measured by the Beckman Synchron LX20 auto-system (Beckman Coulter, Brea, CA, USA) in the central lab of Chung-Ho Memorial Hospital of Kaohsiung Medical University using a spectrophotometric method with picric acid as the reactive at 520 nm. We used total arsenic (tAs) as the sum of inorganic arsenic (As^{III} + As^{V}) and organic arsenic (MMA and DMA) divided by urinary creatinine; this ratio was used in the subsequent analyses.

Assessment of creatinine

Creatinine was measured by the Beckman Synchron LX20 auto-system (Beckman Coulter, Brea, CA, USA) in the central lab of Chung-Ho Memorial Hospital of Kaohsiung Medical University using a spectrophotometric method with picric acid as the reactive at 520 nm.

Assessment of low density lipoprotein (LDL)

The LDL Cholesterol Direct (DLDL) method was used to measure LDL cholesterol from the serum and plasma of the participants using ADVIA® Chemistry systems.

DNA methylation

DNA was isolated from cord blood samples using buffy coat isolated from EDTA-treated blood (Gentra Puregene; Qiagen, Hilden, Germany) and bisulfite converted using the EZ DNA Methylation kit. Samples were randomized across several plates for epigenome-wide DNA methylation assessment using the Illumina Infinium Human Methylation 450 BeadChip which simultaneously profiles the methylation status of > 485,000 CpG sites with single-nucleotide resolution across the human genome.

DNA methylation is measured using beta values, calculated as $M / (M + U + \epsilon)$, where M is the methylation signal of the target CpG, U is the unmethylated signal and $\epsilon = 100$, is a constant to protect division by zero. Thus, average beta-value (β) represents the percent methylation of the target CpG site and its value ranges between 0 and 1.

Quality control

The raw data of DNA methylation were pre-processed to achieve high quality for data analyses. The function *PreprocessSWAN* in the Bioconductor package *minfi* was used for normalization, background correction and peak correction. The function *preprocessSWAN* uses subset within array normalization (SWAN) technique, which normalizes Infinium type I and type II probes together within a single array. This technique reduces technical variability between arrays by accounting for differences in the comparison of the two probe types between arrays [34]. CpG sites located on sex chromosome and annotated probe SNPs within 10bp of the target CpGs were dropped to eliminate bias caused by subject's gender and bias due to genetic variability.

2.3.2 Data collection and pre-processing of birth cohort data from NHBCS (Replication Study sample)

NHBCS Cohort Description:

The New Hampshire Birth Cohort Study (NHBCS) is an ongoing prospective study that began in 2009 and includes over 1500 women from two regions of New Hampshire, USA, enrolled between 24-28 weeks gestation. Mothers were recruited into the cohort if they were literate in English, mentally competent, between 18–45 years old, and reported using a private, unregulated well as the primary source of home drinking water. Infants included in the cohort were singleton pregnancies. Pre- and post-delivery questionnaires were administered to collect self-reported sociodemographic, lifestyle, and medical history data, and a structured medical records review was employed to collect information from the pregnancy and delivery. Cord blood samples are collected on >80% of all deliveries. This study consisted of the first participants born in the study with available cord blood samples for DNA methylation analysis and mothers that were not missing for urinary arsenic or any of the covariate data (n=109).

Arsenic in maternal urine:

Measures of maternal urinary arsenic have been described thoroughly elsewhere [1]. Briefly, spot urine samples were collected between 24–28 weeks gestation with 30 μ L of 10 mM diammonium diethyldithiocarbamate, then samples were frozen at -80°C until analysis. High-performance liquid chromatography inductively coupled plasma mass spectrometry (ICP-MS) system measured individual arsenic species. Samples with values below the limit of detection (LOD) were assigned a value equal to the LOD divided by the square root of two [2]. Total maternal urinary arsenic concentrations (U-As) were calculated as the sum of inorganic arsenic

(As^{III} & As^V), monomethylarsonic acid (MMA^V) and dimethylarsinic acid (DMA^V), which was then log₁₀-transformed prior to analyses.

Cord blood DNA-M processing and QA/QC:

DNA was bisulfite converted using the EZ DNA Methylation kit and subsequently subjected to epigenome-wide DNA methylation assessment using the Illumina Infinium®

HumanMethylation450 BeadChip (Illumina, San Diego, CA) at the University of Minnesota Genomics Core Facility (Minneapolis, MN) following standardized protocols. Post-array processing was conducted in the ‘minfi’ package in R. Samples in which >2% of probes had poor detection p-values were excluded; probes with detection p-values > 0.01 in at least one sample were also removed. Functional normalization (funNorm) and ComBat were utilized to remove technical variations in the data; removal of batch effects was confirmed with principal components analysis. The normalized and batch-corrected beta-values were then transformed into M-values via $\log_2(\beta/(1-\beta))$ prior to statistical analyses.

Covariates:

Maternal age, maternal BMI, and estimated cell proportions were included as continuous covariates. Child gender and mother’s education were included as dichotomous covariates.

Maternal education was defined as those with at least a college degree vs. those without a college degree. Cell type proportions were estimated with the current gold standard method [3] via the minfi package in R. In regression models, 5 of the 6 cell types were included as covariates (NK cells were excluded) to account for overall cellular heterogeneity.

Replication Sample

In total, 109 cord blood samples for which DNA-M had been obtained and complete arsenic and covariate data were available. Of the 579 CpG sites identified in the cohort in Taiwan, 553 were available for replication analyses within the NHBCS.

2.3.3 Correction for cell mixture proportion

Blood is a mixture of functionally and developmentally distinct cell populations [35]. Adjusting for this cell type will remove potential confounding effects of cell heterogeneity in DNA methylation in blood samples [36]. Cell type composition of the blood sample was calculated using function *estimateCellCounts* in the R package *minfi* [37, 38]. IDAT files from 450k Illumina DNA methylation were used to estimate the proportion of 6 cell types: CD8T, CD4T, NK, Monocyte, Granulocyte and B-cell. Cell type proportions are provided in supplemental Table A1.2.

2.3.4 DAVID

Illumina Infinium 450K Human Methylation Beadchip array version 1.2 was used to map the significant CpG sites to USCS reference genome and identify genes associated with these CpG sites. Functional enrichment and pathway analysis of resulting genes was carried out using DAVID gene functional classification tool [39]. DAVID is a large gene-centered knowledgebase which integrates the diverse annotation resources in a centralized location. DAVID knowledgebase is based upon single-linkage algorithm called the DAVID Gene Concept, and it serves as a gene/protein IDs database. The DAVID gene functional classification tool aggregates a list of genes or associated biological terms into organized classes of related genes or biology. For a given gene list, DAVID identifies enriched functional-related gene groups and provides a visualization of genes on pathway maps such as KEGG [40]. KEGG pathway maps are

collection of manually drawn pathways based upon knowledge gained from experiments on functions of the cell and its metabolism. Genetic interaction in KEGG pathway represents the network of molecular interaction and reactions of gene products.

2.3.5 GeneMANIA

GeneMANIA [41] is used to build the network between query genes, based upon genes that are functionally similar. This connection between genes is established based upon their similar expressions and functional association across different conditions via published data.

GeneMANIA uses the publicly available genomics and proteomics data, as well as organism-specific functional genomics data sets. Six organisms are currently supported by GeneMANIA (*Arabidopsis thaliana*, *Caenorhabditis elegans*, *Drosophila melanogaster*, *Mus musculus*, *Homo sapiens*, and *Saccharomyces cerevisiae*).

2.3.6 Statistical Analyses

The dataset consists of 64 samples from cord blood specimen with DNA methylation (DNA-M) data for 485,577 CpG site. After quality control using Bioconductor package *minfi* 385,183 CpG sites were retained for statistical analysis. The pre-processed DNA-M data in beta values were transformed to M values, approximated as $\log_2 [\beta/(1-\beta)]$, in order to ensure a better fit to statistical model assumptions used in our analyses.

To identify CpG sites whose DNA-M is influenced by in utero arsenic exposure (tAs), robust regressions (lmFit function R-package *limma*) [42] were applied to model the association of DNA-M with urinary creatinine-adjusted total arsenic (tAs). Child's sex, batch effect, mother's age, BMI, and education level, and estimated blood cell proportions (CD8T, CD4T, NK, and B-cells, monocytes and granulocytes [37, 38]) were included as covariates. Robust regressions in *limma* package use an empirical Bayes approach to estimate sample variances

which provides stable inference when the number of arrays is small [43]. In the robust regression analyses, multiple testing is adjusted by controlling FDR of 0.05. For the replication analyses we reproduced the statistical models described above in the NHBCS sample (detailed description is given below). CpGs with regression coefficients are in the same directions were considered to be successfully replicated, and we attempted to control for multiple testing via FDR of 0.05.

To assess the association of DNA methylation at CpGs of genes in some of the identified pathways with longitudinal (2, 5, 8, 11 and 14 years) low-density lipoprotein (LDL), a biomarker for CVD and diabetes, we applied linear mixed models. \log_{10} LDL concentrations at different ages were the dependent variable and residuals of DNA methylation, age, as well as interaction between age and DNA methylation were included in the model as predictors, and child's age, sex, and birth weight were treated as covariates. A statistical significance level was set at 0.05. The residuals of DNA methylation were obtained by regressing DNA-methylation at each of 12 CpG sites on proportions of each of the six cell types (CD8T, CD4T, NK, and B-cells, monocytes and granulocytes) and batch.

2.3.6.1 Statistical Analyses in NHBCS:

For all 553 regressions, we tested the linear relationship between maternal total urinary As and cord blood DNA-M M-values while adjusting for confounders. Since batch effects were removed via ComBat during data processing, no batch variable was included in these analyses. The following model, consistent with the model used in the study based on data from Taiwan, was fit using robust regression via the `lmFit` function in the `limma` package in R (version 3.2.2), confidence intervals were extracted using the `confint=TRUE` option.

M-values = \log_{10} (U-As) + Child Gender + Urine Creatinine + Mother's Age + Mother's BMI + Mother's Education + Estimated Cell Proportions

2.4 Results

The data are from a birth cohort study examining multiple in utero and postnatal factors in relation to child health outcomes as part of the nationwide Taiwan Maternal and Infant Cohort Study established in Taiwan in 2000-2001 [21]. In total, 64 subjects with genome-wide DNA methylation in cord blood, level of maternal urinary arsenic exposure, urinary creatinine, along with a child's sex, gestational age, maternal age, maternal pre-pregnancy body mass index (BMI) and the mother's educational level were available and utilized in the study. Table 2.1 presents the characteristics of pregnant women and newborns by sex. Of the 64 newborns, 38 (59.4%) were male. Maternal characteristics are comparable between male and female newborns, and there is no statistically significant difference in gestational ages between sexes of newborns.

Table 2.1. Characteristics of mothers and their newborns by newborn sex in Taiwan during 2000-2001 (n=64)

Characteristics	Sex of the infant			p-value ^b
	All (n=64) ^a	Male (n=38) ^a	Female (n=26) ^a	
Pregnant Women				
Age (years)	28.9±4.8	28.6±4.1	29.5±5.7	0.492
Pre-pregnant BMI (Kg/m ²)	20.5±2.6	20.2±2.4	21.0±2.9	0.244
Urinary Creatinine (mg/dL)	63.6±41.7	70.9±46.0	53.0±32.9	0.078
Maternal Education				0.303
high school + 2 years	25(39)	13(34)	12(48)	
≥high school + 4 years	39(61)	25(66)	14(52)	
Newborns				
Gestational Age (weeks)	39±1.2	39±1.1	39±1.4	0.791

^aPresented as the mean±SD or number (percentage).

^bp-value for difference between male and female newborns using t-test for continuous variables and χ^2 or Fishers Exact Test for categorical variable

The levels and distribution of arsenic metabolites in maternal urine after adjusting for creatinine levels are shown in Table 2.2, distinguishing between mono-methylated arsenic (MMA), di-methylated arsenic (DMA), inorganic arsenic (iAs), and the sum of the three (total arsenic or tAs). Concentrations of each urinary arsenic species showed a large variation among the 64 mothers. We focused on tAs to represent overall arsenic exposure. The distribution of tAs

is severely skewed with a median of 23.19 $\mu\text{g g}^{-1}$ crea ($\mu\text{g g}^{-1}$ crea [creatinine]), and 5th and 95th percentiles being 3.76 $\mu\text{g g}^{-1}$ crea and 76.02 $\mu\text{g g}^{-1}$ crea, respectively (Table 2.1 and Supplemental Figure A1.2). The results reported in this article are based on log10-transformed total arsenic concentration.

Table 2.2. Distribution of creatinine-adjusted concentrations of urinary arsenic species (iAs, MMA, and DMA) (n=64)

Exposure variables ^a \ Percentile ^b	Min	5th	25th	50th	75th	95th	Max
As metabolites ($\mu\text{g g}^{-1}$ crea)							
MMA	0.06	0.08	0.19	0.40 (1.10)	1.67	6.14	28.5
DMA	0.07	3.09	11.27	20.73 (14.58)	29.5	70.75	129.1
iAs	0.11	0.19	0.41	0.83 (0.61)	1.33	4.74	6.55
tAs	0.34	3.76	12.09	23.19 (16.29)	33.29	76.02	137.5

^aAbbreviations: iAs represents the sum of As^{3+} and As^{5+} ; MMA: methylarsonic acid; DMA: dimethylarsinic acid; tAs: the sum of iAs, MMA, and DMA; $\mu\text{g g}^{-1}$ crea: μg per g creatinine

^bThis study. Pregnant women from Maternal Infant cohort in Taiwan (n = 64)

LOD of detection for As^{3+} is 0.09 $\mu\text{g/L}$, As^{5+} is 0.05 $\mu\text{g/L}$, for MMA it is 0.05 $\mu\text{g/L}$ and for DMA it is 0.04 $\mu\text{g/L}$

The values inside parenthesis are the average value of unadjusted arsenic expressed as $\mu\text{g/L}$.

The standard deviation for unadjusted tAs is 16.22 $\mu\text{g/L}$ and interquartile range for adjusted tAs is 21.21 $\mu\text{g g}^{-1}$ cre.

After pre-processing the DNA methylation data as depicted in Figure 2.1, 385,183 CpG sites were analyzed. The flow for the analyses is depicted in Figure 2.2.

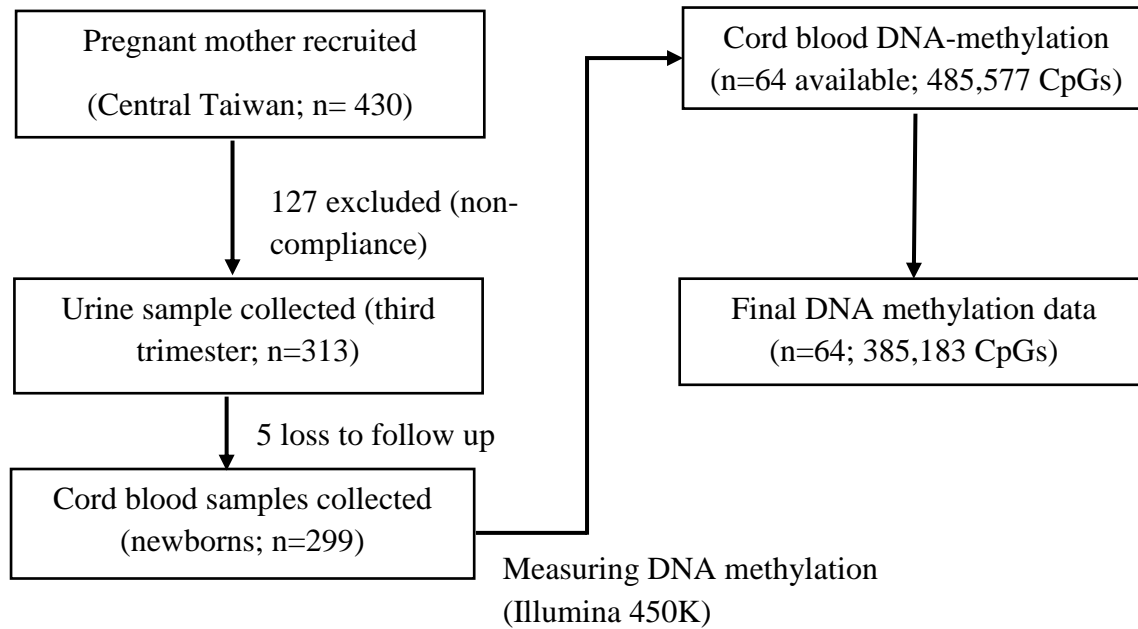


Figure 2.1. Subject recruitment and preprocessing of DNA methylation data in Taiwanese birth cohort

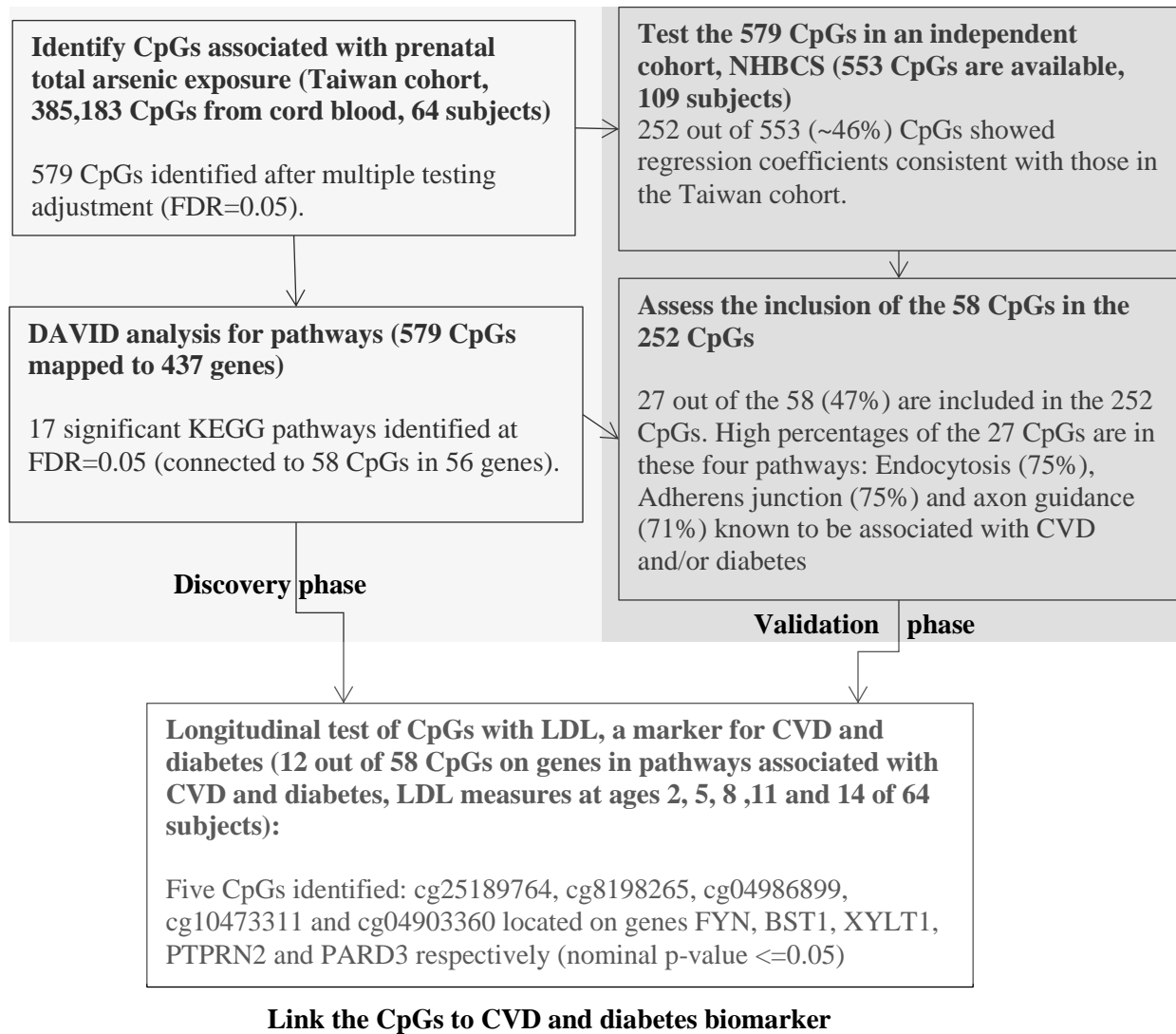


Figure 2.2. The flow of analyses performed in the study.

Epigenome-wide assessments of statistical associations between \log_{10} creatinine-adjusted maternal urinary arsenic level and logit transformed DNA methylation (also noted as M values) were conducted via robust regressions. Covariates included in robust regressions were child's sex, batch of DNA methylation analyses, mother's age, mother's pre-pregnancy BMI, mother's education level, and estimated proportions of six blood cell-types (Appendix Table A1.2, related methods are in the Methods section). Figure 2.3 shows the Manhattan plot of p-values for testing on the 385,183 CpG sites, with a dashed blue line indicating the p-value threshold corresponding to FDR of $p=0.05$ [44]. In total, 579 CpG sites showed statistically significant associations at FDR of 0.05. Supplemental Table A1.1 lists these 579 CpG sites along with their regression coefficients, p-values, and corresponding chromosomes, locations on the chromosomes, corresponding genes, and location on the genes. About 60% of these 579 CpGs showed a positive association between DNA methylation and in utero tAs. The majority of the CpG sites located in the North shore regions of the CpG Island had higher DNA methylation associated with higher in utero tAs and about 39% of these CpG sites were located upstream of transcription start site (TSS1500, TSS200) or 1st Exon (Appendix Table A1.1).

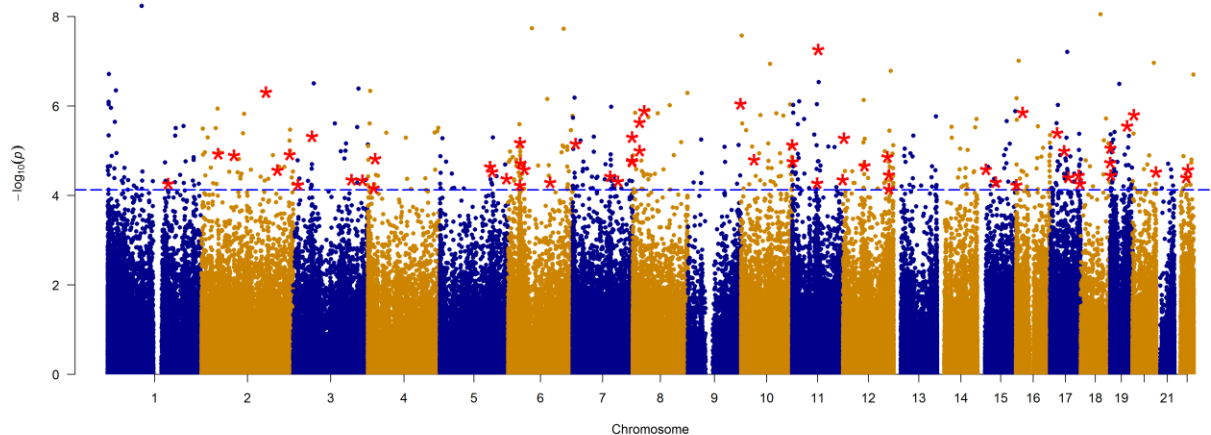


Figure 2.3. Manhattan plot for Genome-wide DNA methylation associated with creatinine adjusted urinary arsenic concentration. The horizontal dashed blue line corresponds to the significance threshold $p = 7.51E-05$ (FDR Adjusted p-value ≤ 0.05), red color stars represent the CpG sites corresponding to genes enriched in KEGG pathways from DAVID analysis. Blue and golden colors are used to differentiate the chromosomes.

The 579 CpG sites were mapped to 437 genes (Appendix Table A1.1), which were further analyzed using the bioinformatics tool DAVID [45, 46]. This analysis led to 17 significantly enriched KEGG pathways (at FDR=0.05) and 58 CpGs were within the genes involved in these pathways, including pathways connected to CVDs and diabetes [47] (e.g., Type I and Type II diabetes mellitus, focal adhesion, calcium signaling pathway, adherens junction, and chondroitin sulfate biosynthesis [48]), pathways linked to neurological and cognitive abilities (Alzheimer’s disease and amyotrophic lateral sclerosis [ALS]), and pathways in cancer (the 58 CpG sites involved in these pathways are marked by red stars in Figure 2.3). The network, constructed using GeneMANIA [41] based on the genes enriched in DAVID analysis, indicated inter-connections among the genes (Supplemental Figure A1.3) via co-expression or shared pathways. Among these 58 CpG sites corresponding to the genes enriched

in KEGG pathways, most of them are located in the body region of a gene (Figure 2.4). Majority of these 58 CpGs are located in the island region (~57%) or north Shore (~22%). Furthermore, in approximately 55% out of the 58 CpG sites, we found that higher in utero tAs were linked to higher DNA methylation in cord blood, as indicated by positive regression coefficients in Figure 2.4. The strongest association between in utero tAs and cord blood DNA methylation occurred at CpG cg23767840, which is in the 5'UTR region of gene *EPN2* (coding for the Epsin-2 protein).

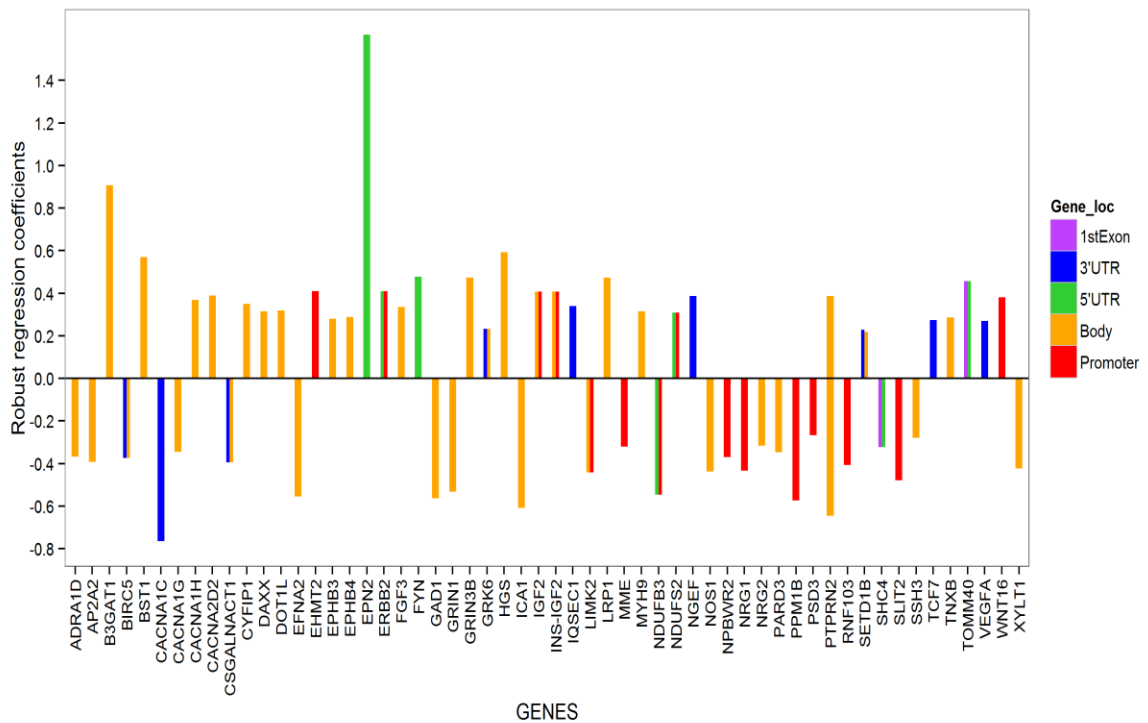


Figure 2.4. Association of arsenic exposure with the DNA methylation based on M-values of the 58 CpG sites mapped to 56 genes. The x-axis has the 56 genes enriched in KEGG pathways at FDR level of $p=0.05$, while the y-axis shows the estimates of total arsenic coefficients related to 58 CpG sites from robust regression. Adjusting factors include cell counts, child's sex, batch effect, mother's age, mother's BMI and mother's education level. M-values are defined as $\log_2 [\beta/(1-\beta)]$. Different colors indicate the location of the CpGs on a gene.

The resulting 579 CpG sites from our study were further tested in the independent New Hampshire Birth Cohort Study (NHBCS) (n=109). Of the 579 CpG sites 553 were available for analyses in NHBCS. We applied robust regression models with covariates comparable to those included in our study to assess the association of tAs with cord blood DNA methylation at these 553 CpG sites. At 46% of the 553 CpG sites (252 CpGs), the associations of in utero tAs with cord blood DNA-methylation levels were consistent with those found in our study in terms of direction of regression coefficients, although none survived multiple testing. In addition, 27 of these 252 CpGs are in the list of 58 CpGs (27/58= \sim 47%) noted earlier that are involved in the enriched KEGG pathways (Table 2.3). Genes corresponding to these 27 CpGs are more often linked to pathways involved in endocytosis, adherens junction, axon guidance (a neural developmental process in which neurons send out axons to reach the correct targets) and chondroitin sulfate biosynthesis.

Table 2.3. KEGG pathways identified using DAVID that are more specific to arsenic exposure based on data from n=64 pregnant women from the maternal infant cohort in Taiwan.

KEGG-Pathways	Genes	Adjusted p-value (FDR=0.05)
Calcium signaling pathway	NOS1*, BST1, ERBB2*, GRIN1, CACNA1G, CACNA1H, CACNA1C, ADRA1D	0.000011
Endocytosis	PARD3*, AP2A2*, RNF103, PSD3, GRK6*, HGS*, IQSEC1*, EPN2	0.0000010
Axon guidance	NGEF*, LIMK2, FYN*, EFNA2*, EPHB3*, EPHB4, SLIT2*	0.000015
Alzheimer's disease	NDUFB3, NOS1*, LRP1, GRIN1, MME, CACNA1C, NDUFS2*	0.000033
MAPK signaling pathway	CACNA1G, CACNA1H, PPM1B*, CACNA1C, CACNA2D2, DAXX*, FGF3	0.00057
Regulation of actin cytoskeleton	LIMK2, INS-IGF2, SSH3, CYFIP1*, IGF2, MYH9*, FGF3	0.0019
Type I diabetes mellitus	ICA1*, INS-IGF2, PTPRN2*, IGF2, GAD1	0.0018
Amyotrophic lateral sclerosis (ALS)	NOS1*, GRIN1, TOMM40, DAXX*	0.0037
Adherens junction	TCF7, PARD3*, FYN*, ERBB2*	0.0087
Pathways in cancer	WNT16*, TCF7, ERBB2*, VEGFA, BIRC5, FGF3	0.0091
ErbB signaling pathway	ERBB2*, NRG1, NRG2, SHC4	0.01
Focal adhesion	TNXB*, FYN*, ERBB2*, VEGFA, SHC4	0.0097
Chondroitin sulfate biosynthesis	CSGALNACT1*, B3GAT1, XYLT1*	0.01
Neuroactive ligand-receptor interaction	PARD3*, GRIN1, NPBWR2, GRIN3B*, ADRA1D	0.017
Lysine degradation	DOT1L*, SETD1B, EHMT2*	0.034
Type II diabetes mellitus	INS-IGF2, CACNA1G, IGF2, CACNA1C	0.038
Huntington's disease	NDUFB3, AP2A2*, GRIN1, NDUFS2*	0.045

* CpG sites of these genes were consistently associated (in terms of regression coefficient) with total urinary arsenic exposure in an independent cohort NHBCS.

Given the connection of arsenic exposure with CVDs and diabetes [24, 25, 49, 50], findings from the pathway analyses, and findings in the replication study, we further investigated the CpG sites of the genes enriched in KEGG pathways that are potentially linked to cardiovascular diseases and diabetes in our Taiwan cohort. In particular, 12 CpGs (located on 11 genes, Appendix table A1.1) were included in this analysis and these 12 CpGs were among the 27 CpGs replicated in the NHBCS cohort. We assessed the association of cord blood DNA methylation at these CpGs with a biomarker of CVDs and diabetes, plasma low density lipoprotein (LDL). LDL was measured at multiple ages of the children (at 2, 5, 8, 11, and 14 years). Plasma LDL concentration is the most stable in humans, with or without fasting, among blood lipids such as triglycerides. Among the 12 CpGs, cord blood DNA methylation at some CpGs showed a pattern of positive correlations with LDL at each age. While some were negatively correlated with LDL at age 2 and positively correlated at later ages (Figure 2.5), for most CpGs, the strongest correlations (positive or negative) occurred at age 2. In particular, the heatmap (Figure 2.5) indicated that DNA methylation levels at two CpGs, cg06419180 and cg25189764, were positively correlated with the LDL at different ages, while the directions of correlations at the rest of the CpG sites seemed to change over time. Via linear mixed models, we tested the association of LDL with DNA methylation (with LDL at ages 2, 5, 8, 11 and 14 as the outcome, cell type compositions and batch-effect adjusted DNA methylation as the predictor, and child's age, sex of the child, and birth weight as covariates) as well as the interaction effect between DNA methylation and age. We found that CpG cg25189764 had a statistically significant association with LDL (coefficient=0.157, p-value=0.047). DNA methylation at another 4 CpG sites showed statistically significant interaction with time, cg08198265 (coefficient for main effect [coeff.m] = 0.498, coefficient for interaction effect [coeff.int] = -

0.063, p-value=0.002), cg04986899 (coeff.m = -0.36, coeff.int = 0.030, p-value=0.024), cg10473311 (coeff.m= 0.145, coeff.int= -0.021, p-value=0.027) and cg04903360 (coeff.m= -0.189, coeff.int=0.026, p-value=0.032). It is worth noting that DNA methylation at these 5 CpG sites was found to be stable across the life course. The stability was assessed using Accessible Resource for Integrated Epigenomic Studies (ARIES) explorer [51].

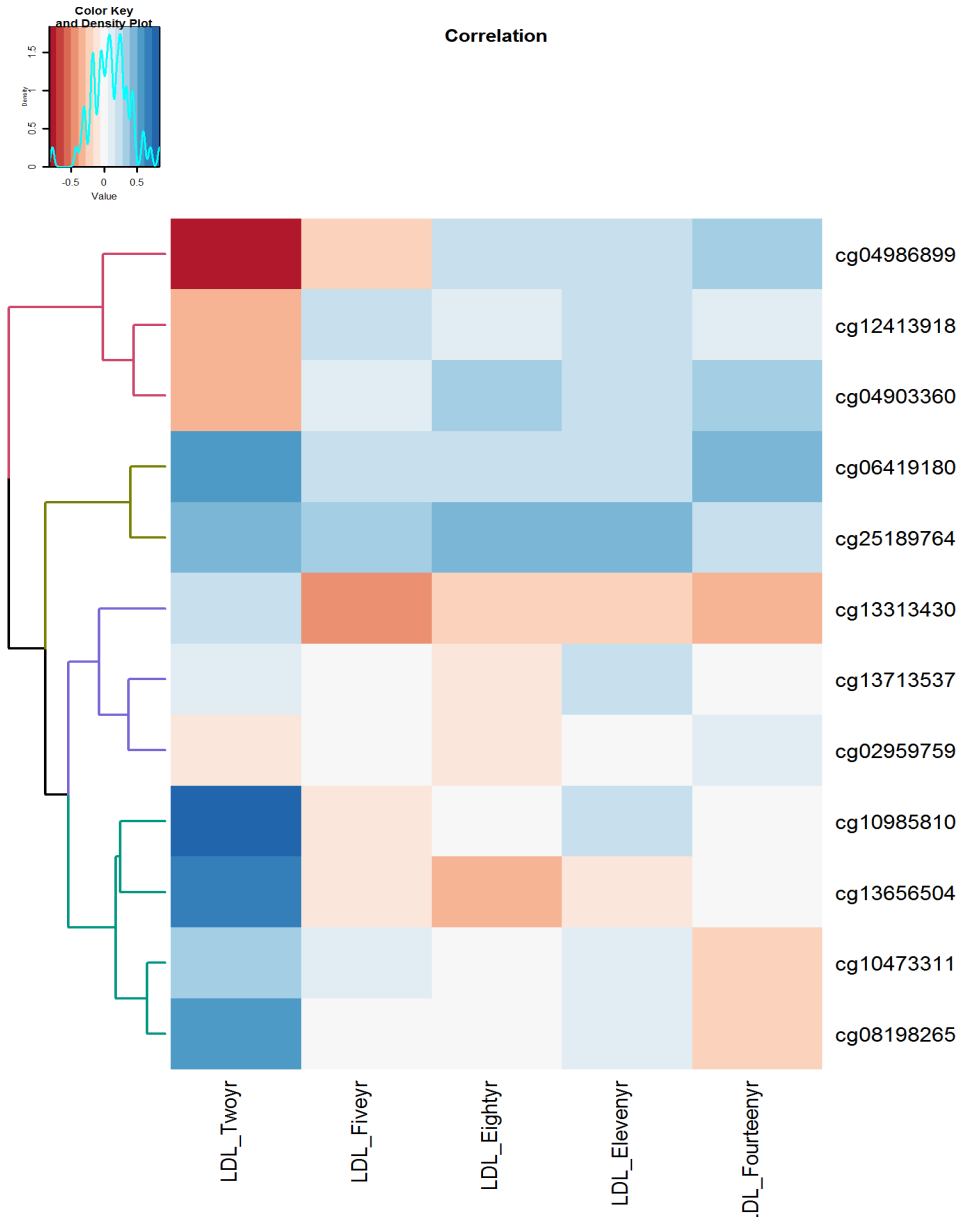


Figure 2.5. Heatmap of the correlations between cord blood DNA methylation and LDL across different ages (2, 5, 8, 11, 14 years).

2.5 Discussion

The overall aim of this study was to identify CpG sites that would represent biomarkers of possible adverse effects of arsenic in newborns and of future health outcomes. In total, at 579 CpGs identified from a cohort in Taiwan DNA methylation was associated with in utero arsenic exposure. To further understand the biological mechanisms of genes linked to these 579 CpG sites, a gene annotation analysis using DAVID was performed, which led to an identification of 17 statistically significant KEGG pathways. Genes corresponding to the identified CpGs are known to be involved in arsenic-associated diseases including neuronal [52-54], immune [55], cancer [56], cardiovascular and diabetes [25, 49, 50]. Experimental models have demonstrated a role of in utero acquired somatic epigenetic alternations in diseases [57-59]. Given the regulatory functionality of DNA methylation on different genes, the identified CpG sites may serve as epigenetic biomarkers of potential harmful effects of *in-utero* arsenic exposure among newborns.

Findings at 46% of the identified 579 CpG sites were replicated in an independent cohort, the NHBCS, with respect to directions of associations, though these did not survive multiple testing adjustments. However, the median tAs (without creatinine adjustment) in NHBCS was 2.8 µg/L with interquartile range (IQR) of 3.64 µg/L, which is substantially lower than that in the Taiwanese cohort (median= 11.51 µg/L and IQR= 16.80 µg/L). This difference, small sample sizes from both studies, differences in ancestry and unmeasured confounding may explain the limited agreement in the findings between the two cohorts.

The post hoc analysis on CpG sites replicated in the NHBCS cohort and related to genes enriched in KEGG pathways for cardiovascular disease and diabetes led to the identification of five CpG sites cg25189764, cg08198265, cg04986899, cg10473311 and cg04903360 located on genes *FYN*, *BST1*, *XYLT1*, *PTPRN2* and *PARD3*, respectively. *FYN* is an important regulator of

whole body metabolism and is known to be associated with insulin sensitivity in mice [60]. *BST-1* is a glycosyl-phosphatidylinositol (GPI) and is expressed in abundant in pancreatic islet cells [61]. Proteins containing a *GPI* anchor play key roles in a wide variety of biological processes [62]. *XYLT1* is involved in heparan sulfate (a type of glycosaminoglycan; GAG) biosynthesis [63, 64]. *GAGs* have been studied for their role as a potential target in treating CVDs [65, 66]. Protein encoded by *PTPRN2* (also known as *IAR*) is a known autoantigen in insulin-dependent diabetes mellitus [67]. *PARD3* has been identified as candidate gene for its association with type 2 diabetes in Mexican study [68]. Out of these five CpGs, cg25189764 is located in the 5'UTR of gene *FYN*, and the other four CpGs were located in the body of the genes. We observed that most CpG sites on genes enriched in KEGG pathways were located in the body region of a gene (Figure 2.4). The regulatory functionality of DNA methylation on genes at those CpG sites is likely to be different from the functionality at CpG sites in the promoter region [69, 70]. Further assessment on their associations with gene expressions will improve our understanding of their regulatory functionality.

The temporal stability in DNA methylation at the five CpG sites (cg25189764, cg08198265, cg04986899, cg10473311 and cg04903360) showing associations with LDL across different ages raised a possibility of long term consequences of DNA methylation, established in utero, on LDL at later life. More interestingly, for the four CpGs (cg08198265, cg04986899, cg10473311 and cg04903360) showing interactions with age, the turning point of DNA methylation effects (from negative to positive effects or from positive to negative effects over time) are always around the age of 8 years. Of interest, ages 11 and 14 are during adolescence, a period of significant changes, e.g., puberty, rapid growth, and often BMI increase.

A previous study *in utero* arsenic exposure in the NHBCS was reported by Koestler et al. [32]. The top 100 CpGs identified in Koestler et al. did not overlap with the 579 CpGs, although 25% of their 100 CpGs showed statistical significance at the 0.05 level in our study (not surviving multiple testing). The disagreement could have been driven by some key differences in the analytical methods. Koestler et al. categorized arsenic exposure levels into quartiles and applied analysis of covariance with tests for trends, while our study applied robust regressions to log₁₀-transformed arsenic concentrations to take into account possible outliers. By categorizing a continuous variable, statistical testing power for testing the associations might have been reduced. In addition, Koestler et al. did not adjust for maternal BMI, nor the cell type proportions estimated using the *minfi* R package [37, 38], though they did explore associations between urinary arsenic and estimated cell-type proportions in cord blood.

We also compared the findings from our study with another epigenome-wide study by Broberg et al [16]. The focus of that study also concentrated on the top CpG sites ranked by statistical significance on their association with in utero arsenic exposure, although none of the top CpG sites survived multiple testing corrections. The top CpG sites determined by Broberg et al. did not overlap with those identified in our study, nor overlapped with the top CpGs in Koestler et al. [32]. Broberg et al [16] utilized linear regression and did not adjust for cell type heterogeneity. In addition, some top CpG sites discussed in Broberg et al. included annotated probe-SNPs (single nucleotide polymorphisms) located within 10 base-pairs of the target CpG. They can result in biased methylation measurements, and were excluded from our analysis.

It is worth noting that the three studies we discussed herein (Koestler et al. [32], Broberg et al. [16], and ours) were conducted in different regions (United States, Bangladesh, and Taiwan, respectively) with vastly different medians in utero arsenic exposures which may have

limited replicability (for tAs, in Koestler et al., median=4.1 µg/L, in Broberg et al., median=66 µg/l, and in our study, median= 11.51 µg/L). It is also possible that ancestry, race/ethnicity or other regional differences may have contributed to the disagreement in the findings. In addition, all studies had small sample sizes (less than 200), so some of the findings are also likely to be false-positives. A large-scale study incorporating different races/ethnicities, with a wide exposure range, is well deserved. Our study had a benefit of replicating results using standard statistical approaches. Nonetheless, replicating DNA methylation analyses in additional populations, harmonizing, and comparing different DNA methylation studies on in utero arsenic exposure will help to assess the generalizability of the results. Future studies also should be directed at examining whether arsenic-related health outcomes are associated with cord blood DNA methylation in a long-term follow-up of the children in multiple cohorts.

2.6 Conclusion

We found that *in utero* arsenic exposure was associated with cord blood DNA methylation. The genes corresponding to the identified CpG sites were involved in various pathways including signaling pathways, Type I and Type II diabetes mellitus, and neuroactive ligand-receptor interactions. Cord blood DNA methylation at cg25189764, cg08198265, cg04986899, cg10473311 and cg04903360 were associated with low-density lipoprotein (LDL) at later life. These CpGs need to be studied further for their role in cardiovascular disease and diabetes in arsenic-exposed populations. Although larger studies are needed, results from this study contribute to a better understanding of epigenetic mechanism of diseases related to in utero arsenic exposure in infants.

3 Epigenetics markers at birth longitudinally associated with Immunoglobulin E (IgE)

3.1 Abstract

Background: Immunoglobulin E (IgE) is known to play a major role in allergic diseases.

Epigenetic markings acquired due to modification of DNA methylation in early life may have phenotypic consequences later in development through their role in transcriptional regulation with relevance to the developmental origins of diseases including allergy. However, epigenome-wide studies on the association of cord blood DNA methylation and IgE over time are lacking.

Method: A total of 64-cord blood samples from Taiwan Maternal and Infant Cohort Study were analyzed using the Infinium Human Beadchip to obtain DNA methylation at ~450K Cytosine-phosphate-Guanine (CpG) sites. Linear mixed models were implemented to assess the association between preprocessed, batch and cell type corrected cord blood DNA methylation at >380k CpG sites with IgE levels at 5,8 and 11 years of age, adjusting for cord blood IgE.

Identified statistically significant (at a false discovery rate, FDR, of 0.05) CpGs were replicated in an independent cohort, Isle of Wight (IoW) dataset. Gene ontology analysis was performed using DAVID to identify significantly enriched biological process of genes associated with resulting CpG sites. Stability assessment of the identified and replicated CpG sites was measured using ICC.

Results: DNA-methylation of 458 CpG sites were significantly ($FDR \leq 0.05$) associated with IgE levels at different ages. Among the identified CpG sites available in both cohorts ($n=241$), of which, about 50% of CpGs were replicated in the IoW cohort in terms of consistency in direction of association between DNA methylation and longitudinal IgE levels. Gene ontology analysis of 84 genes linked to 124 CpG sites led to the enrichment of statistically significant biological process: PI3K-Akt signaling pathway, pathways in cancer, metabolic pathways, polymorphism,

alternative splicing, phosphoprotein, disease mutation, glycoprotein, protein transport, transcription regulation. Further temporal stability assessment of the 124 CpG sites identified 59 CpGs with significant ICC values ($p\text{-values}\leq 0.05$) at least 0.5.

Conclusion: Biological finding combined with the temporal stability measures for 59 CpG sites suggest them as a potential epigenetic marker for predicting later IgE production.

3.2 Background:

The prevalence of allergic diseases is increasing worldwide and the severity of allergic diseases, including asthma, continues to increase in children and young adults. About fifty percent of school children are sensitized to one or more common allergens [71]. Allergic disease is known to be hereditary implying individual's susceptibility to the genetic factors [72].

The gene-environment interaction during critical periods of immune development is assumed to be one of the causes of this disease later in life. Epigenetic variation is postulated to be an important mechanism through which these interactions are mediated [73]. Epigenetic processes regulate gene expression during immune development, and evidence suggests disruption in these processes can modify disease risk in a manner analogous to single nucleotide polymorphisms (SNPs) [74]. DNA methylation is one such epigenetic process which is associated with gene silencing and with the patterning of gene expression that determines cellular types and functions.

Immunoglobulin E (IgE) is known to play a major role in many of the allergic diseases such as Asthma, atopic dermatitis (eczema) and hay fever. IgE production leads to type I hypersensitivity, which manifests various allergic diseases. However, the mechanism underlying IgE production is poorly understood. There is evidence that DNA methylation is associated with total IgE. For instance, epigenome-wide study using Illumina methylation 27k array have

identified CpG loci from peripheral blood associated with total serum IgE [75]. Cross sectional study using peripheral blood of 18 year old men and women was used to identify association between CpG loci and serum IgE [76]. However, longitudinal epigenome-wide studies on the association of DNA methylation at birth via cord blood and IgE over time are lacking. Thus, it is unknown whether and how DNA methylation in cord blood is associated with IgE during the course of early life.

In this study we aimed at assessing genome-wide DNA methylation with respect to their association with IgE levels of child at ages 5, 8 and 11 years via linear mixed models adjusting for multiple testing by controlling FDR, and examine possible pathways of genes involved in the identified CpGs. We tested the identified CpGs in an independent cohort, the Isle of Wight (IoW) birth cohort in the United Kingdom. The findings will contribute to an improved understanding of in utero epigenetic development (or mechanism) leading to IgE changes later in life.

3.3. Method

3.3.1 Data collection and pre-processing of birth cohort data from Taiwan

In this section, we focus on the assessment of IgE. For other contents of this cohort including information related to DNA methylation measurement, quality control, and cell type compositions, please refer to Chapter 2, the Method Section.

Assessment of Immunoglobulin E

To determine serum IgE levels in children at 5, 8 and 11 years, blood samples (0.5 mL) obtained via venipuncture were centrifuged and the sera stored at -20°C prior to analysis. Serum total IgE (tIgE) levels were measured using the ADVIA Centaur chemiluminescence immunoassay system (Siemens Healthcare Diagnostics; Deerfield, IL, USA). The assay range was approximately

1.5–3000 IU/mL.

3.3.2 Data collection and pre-processing of birth cohort data from IoW

Participant selection and sample collection

The Isle of Wight (IOW) birth cohort was established to study the natural history of asthma and allergies in children born between January 1, 1989 and February 28, 1990 on the Isle of Wight, UK. The study was approved by the local research ethics committee (now named the National Research Ethics Service, NRES Committee South Central – Southampton B, 06/Q1701/34) and written informed consent was provided by the infants' parents (F1 generation). Details about the birth cohort have been described in detail elsewhere [32–34].

Immunoglobulin E (IgE) measurement:

Cord blood was taken from 1023 infants. Blood was collected from umbilical cord using fine needle in a specimen bottle containing dipotassium EDTA as anti-coagulant. Duplicate measurements of cord IgE were made using ULTRA EIA® kit (Pharmacia Diagnostics AB, Uppsala, Sweden) unmodified, designed to measure IgE between 0.2 and 50 ku/l on 0.1 ml of serum or plasma. Total IgE was measured in samples of serum collected at age 10 (n = 923) and age 18. IgE at age 10 and at age 18 were determined using PRIST® (Pharmacia Diagnostics AB, Uppsala, Sweden) designed to measure IgE between 2.0 to 1000 kU/L.

DNA sample collection and processing

Blood samples for the F1 generation were obtained from Guthrie card (n=34) of newborn babies. For methylation assays, DNA was extracted from whole blood using a standard salting out procedure [38]. DNA concentration was determined by the PicoGreen dsDNA quantitation kit (Molecular Probes, Inc., OR, USA). One microgram of DNA was bisulfite-treated for cytosine to thymine conversion using the EZ 96-DNA methylation kit (Zymo Research, CA, USA), following the manufacturer's standard protocol. This process converted unmethylated cytosines

into thymines (T) while leaving methylated cytosines (C) unaltered, allowing the array technology described below to distinguish between unmethylated and methylated sites by sequence (C vs. T) recognition [39]. Probes, 50 nucleotides in length, were developed to target specified CpGs throughout the genome with the CpG at the very end of the probe. Multiple copies of each probe were attached to a BeadChip, allowing for the interrogation of multiple DNA reads for each target CpG site [39,40]. After bisulfite conversion, the DNA samples underwent whole genome amplification, then each sample was washed over an array containing many BeadChips to identify the proportion of methylated probes for each target CpG.

Genome-wide DNA methylation was assessed using the Illumina Infinium HumanMethylation450K BeadChip (Illumina, Inc., CA, USA), which interrogates >484,000 CpG sites associated with approximately 24,000 genes. Arrays were processed using a standard protocol as described elsewhere [40], with multiple identical control samples assigned to each bisulphite conversion batch to assess assay variability and samples randomly distributed on microarrays to control against batch effects. The BeadChips were scanned using a BeadStation, and the Methylation Module of BeadStudio software calculated the methylation level for each queried CpG as beta (β) values. They represent the proportions of methylated (M) over methylated (M) and unmethylated (U) sites ($\beta = M/[c+M+U]$) with constant c introduced for the situation of too small $M+U$).

Quality control (QC) measures were employed to improve the reliability of data prior to analysis. In our study, the detection P -value reported by BeadStudio (Illumina software to process raw intensities) was used as a QC measure of probe performance, in which large P -values were deemed to be unreliable measures of methylation. Probes whose detection P -values > 0.01 in >10% of the samples were removed [41]. The methylation data were then preprocessed

using the Bioconductor IMA package for peak correction, background noise removal and batch effect correction [77]. The program for data cleaning was written in R (R Development Core Team, 2012). Some probes may overlap with known single nucleotide polymorphisms (SNPs), which may result in measurement errors [43]. SNPs at the target CpG [44,45] or within close proximity to the target CpG appear to be most likely to affect measurement of methylation, thus probes with SNPs within 5 nucleotides of the target CpG [46,47] and within the probe binding region were excluded from further analyses. Probes more than 5 nucleotides from the target CpG, but still within the binding region will be flagged with an indicator variable, and followed up if selected from the GWAS analysis.

3.3.3 Functional annotation and pathway analysis

DAVID analyses were performed in this study, and the method was discussed in Chapter 2. In the following, I focus on the statistical methods.

3.3.4 Statistical analyses

The final dataset consisted of 64 samples from cord blood specimen with DNA methylation data for 485,577 CpG site. After quality control using Bioconductor package *minfi* 385,183 CpG sites were retained for statistical analysis. The pre-processed DNA-M data in beta values were transformed to M values, approximated as $\log_2 [\beta/(1-\beta)]$, in order to ensure a better fit to statistical model assumptions used in our analyses. To identify CpG sites whose DNA-M could predict IgE levels in children at 5, 8 and 11 years of age, the analysis was performed in two stages. In stage 1 we obtained the residuals of DNA-M by regressing DNA-M of each CpG (385,183) on cell proportions and batch effect. In stage 2 we used residuals of DNA-M to predict the IgE in longitudinal setting using linear mixed model, while adjusting for cord blood IgE, birth weight and gender of the child. Linear mixed model was fitted using *proc mixed* in SAS

9.4. Covariance structure was determined by comparing the Schwarz's Bayesian information criterion (BIC) of four covariance structure: Compound symmetry, Toeplitz, Unstructured and Auto-regressive¹. To determine the best covariance structure 100 CpG sites were randomly selected to fit mixed model under each of the four covariance structure, and covariance structure with the smallest BIC was chosen for the final model. In all analyses, multiple testing is adjusted by controlling false discovery rate of $p=0.05$. Statistically significant CpG sites were further replicated in an independent cohort Isle of Wight (IoW), using the statistical method similar to our study.

An R package *irr* was used to estimate intra-class correlations of DNA-M at 124 CpG sites at birth via Guthrie card, at ages 10 and age 18 years in the IoW cohort.

3.3.4.1 Statistical analyses in the IoW

Analysis similar to described in 3.3.5 was performed for data from the IoW cohort. In the IoW, DNA methylation at birth was assessed from Guthrie cards. IgE was assessed at birth from Cord blood, at age 10 and age 18 years. Longitudinal association between residuals of DNA methylation and \log_{10} IgE was assessed using *proc mixed* in SAS 9.4.

3.4 Results

The data are from a birth cohort study examining multiple prenatal and postnatal factors in relation to child health outcomes as part of the nationwide Taiwan Maternal and Infant Cohort Study [78, 79] established in Taiwan in 2000-2001. In total, 64 subjects with genome-wide DNA methylation from cord blood, child's gender, batch effect, and birth weight were available and utilized in the study.

Table 2.1 as described in Chapter 2 presents the characteristics of pregnant women and newborns by sex. Of the 64 newborns, 38 (59%) were male and 26 (40%) were female. The levels

and distribution of cord blood and serum IgE for children's at ages 5, 8 and 11 years is provided in Table 3.1.

Table 3.1. Distribution of IgE across different time points

Outcome Variable\Percentile	Min	5th	25th	50th	75th	95th	Max
IgE in Cord blood (IU/ml)	0.03	0.03	0.06	0.215	0.885	22.9	61.4
IgE at 5 years (IU/ml)	5.7	7.59	17.95	61.4	125.13	303.65	524
IgE at 8 years (IU/ml)	6.38	7.78	11	96.2	169	695	921
IgE at 11 years (IU/ml)	6	6	15.25	55	185.25	616.35	946

Epigenome-wide assessments of statistical associations between \log_{10} IgE and residuals of DNA methylation (DNA-M) in cord blood at 385,183 CpG sites were conducted in a longitudinal setting via linear mixed modeling. This analysis was performed in two stages. In stage 1 we obtained the residuals of DNA-M by regressing DNA-M of each CpG (385,183) on cell proportions (CD4T, CD8T, NK, Mono, Bcell, Gran, Eos) and an indicator variables associated with different batches of DNA methylation data. In stage 2 we used residuals of DNA-M in cord blood to predict IgE at ages 5, 8, and 11 years using linear mixed models, adjusting for cord blood IgE, birth weight and gender of the child. In total, 458 CpG sites showed statistically significant associations with total IgE, after correcting for multiple testing by controlling FDR of 0.05. Figure 3.1 shows the Manhattan plot of p-values for testing on the 385,183 CpG sites, with a red line indicating the p-value threshold corresponding to FDR of 0.05 [44]. Appendix Table A2.1 lists the 458 CpG sites with their regression coefficients along with p-values, chromosomes they belong to, location on the chromosome, their corresponding genes and their location on the genes.

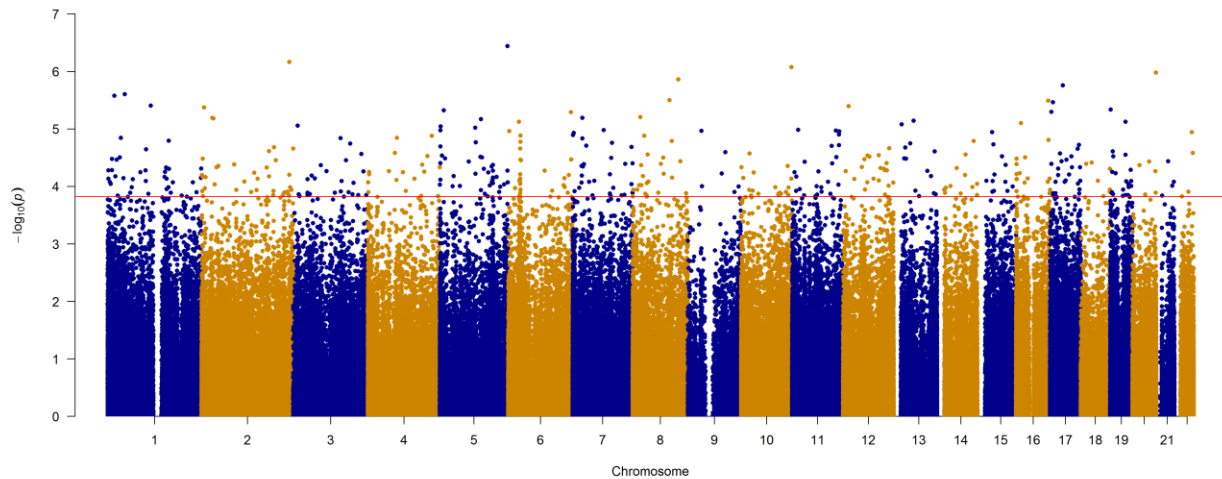


Figure 3.1 Manhattan plot for the longitudinal association of Genome-wide DNA methylation with \log_{10} Immunoglobulin E (IgE). The horizontal dashed red line corresponds to the significance threshold $p = 7.51E-05$ (FDR Adjusted p-value ≤ 0.05). Blue and golden colors are used to differentiate the chromosomes.

The resulting 458 CpG sites from our study were further tested in an independent cohort, the Isle of Wight (IoW) birth cohort. Of the 458 CpG sites 241 were available for analyses in IoW. The flow of this study is provided in Figure 3.2. We used linear mixed models to assess the longitudinal association between residuals of DNA-methylation of 241 CpG sites from Guthrie card blood samples with serum IgE of 30 children measured at ages 10 and 18 years as outcome. The analysis was performed in two stages similar to the main study and used the same additional covariates. At about 51% of the 241 CpG sites (124 CpGs), the longitudinal associations of Guthrie card blood DNA-methylation with IgE over time were consistent with those found in our study in terms of directions of regression coefficients, although none survived multiple testing. Majority of the 124 CpGs were located in body (~50%) and promoter regions (~48%) of the genes, Figure 3.3. In addition, 22 of these 124 CpGs were nominally significant with p-value ≤ 0.05 . Functional annotation analysis using DAVID of 84 genes corresponding to 124 CpG sites led to identification of the following statistically significant (FDR ≤ 0.05) KEGG pathways

and functional categories: PI3K-Akt signaling pathway, pathways in cancer, metabolic pathways, polymorphism, alternative splicing, phosphoprotein, disease mutation, glycoprotein, protein transport, transcription regulation.

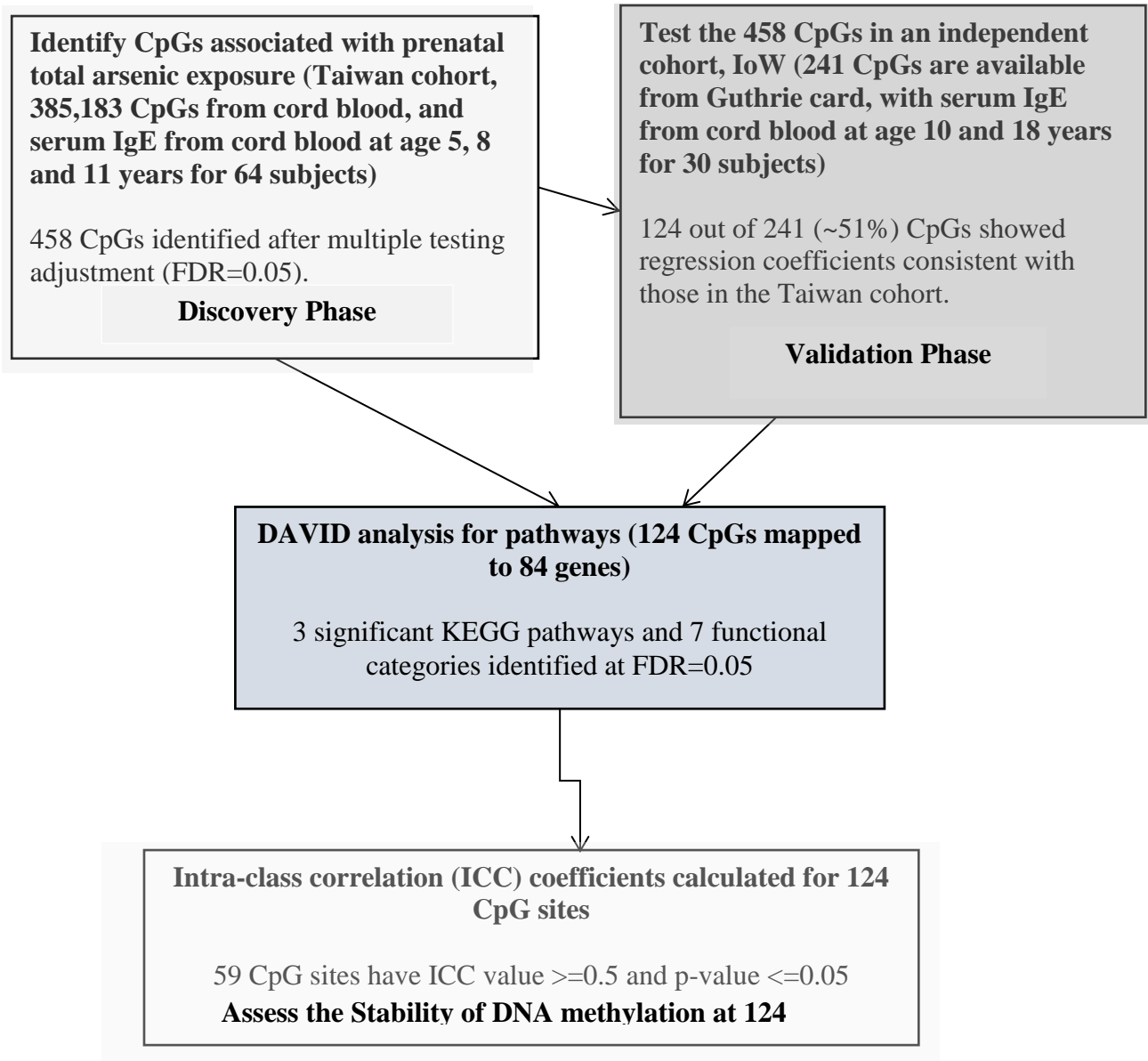


Figure 3.2 Flow of analysis.

The stability of these 124 CpG sites were inferred via intra-class correlation (ICC) based on DNA methylation data in the IoW cohort. In particular, DNA methylation of these 124 CpG sites at birth from Guthrie cards, at age 10 and age 18 years for five subjects were included in this assessment. Among the 124 CpGs, DNA-M at 59 CpG sites had ICC at least 0.5 ($p\text{-value} \leq 0.05$) (Appendix Table T2.2).

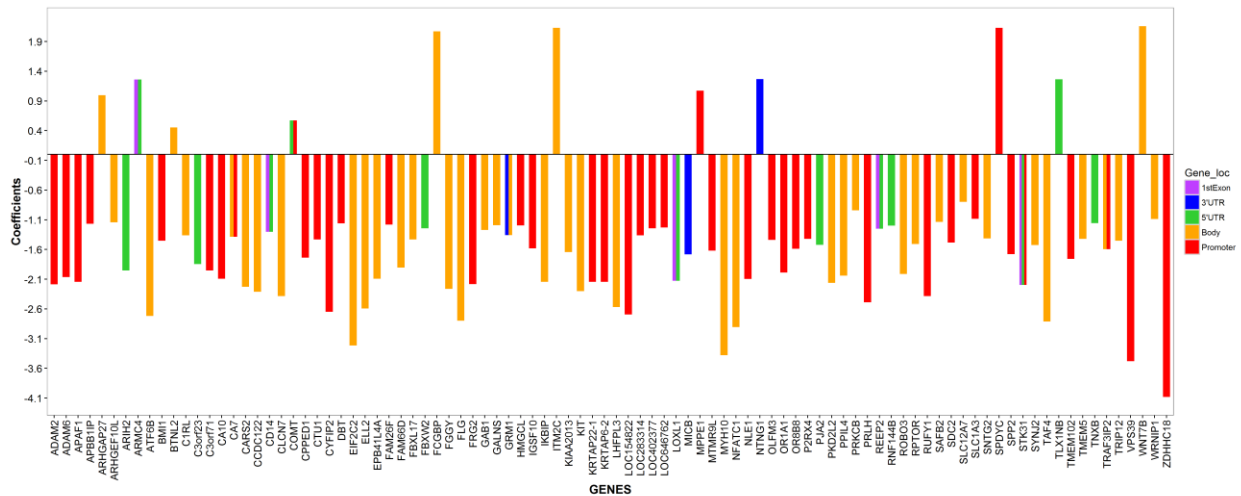


Figure 3.3 Longitudinal association of the residual of DNA methylation with \log_{10} Immunoglobulin E (IgE) of the 124 CpG sites mapped to 89 genes. Please refer to method section for the detail of analysis. Different colors indicate the location of the CpGs on a gene.

3.5 Discussion

The overall aim of this study was to identify CpG sites of which DNA methylation at birth measured in cord blood could potentially predict the level of IgE at later ages. We identified 458 CpG sites (at FDR cutoff of 0.05) longitudinally associated with IgE in discovery cohort, 241 of these 458 CpG sites were available in the replication IoW birth cohort, and 124 out of the 241

CpG sites showed results consistent with those from the discovery cohort in terms of direction of association with IgE. These 124 CpGs were on 84 genes.

Functional analyses indicated that a number of genes associated with the 124 CpG sites were involved in biological processes related to IgE production. Phosphoinositide 3-kinase (*PI3K*) signaling is known to play a crucial role in IgE production, blockade of *PI3K* enhances IgE levels [80]. IgE production is also known to be altered in cancer patients compare to control groups [81] and thus have been studied for their role in tumor surveillance and immunotherapy of cancer patients [82, 83]. Gene *CD14*, one of the 84 genes has been shown to elevate the IgE production [84]. Phosphodiesterase 11A (*PDE11A*) has been suggested for its role in asthma pathogenesis [85, 86], which indicates a potential role of this gene in IgE production.

Stability assessment of DNA methylation based on ICC at 124 CpG sites from birth to age 18 revealed that DNA methylation at 59 CpG sites is likely to be stable. SDA1 Domain Containing 1 (*SDADI*) gene associated with one these 59 CpG sites has been known to contribute towards development of seasonal allergic rhinitis in Japanese population [87]. Fc Fragment Of IgG Binding Protein (*FCGBP*) another gene with CpG site having stable DNA methylation is known to share significant sequence homology with the carboxyl terminal of IgE binding protein [88]. Similarly, Immunoglobulin Superfamily Member 10 (*IGSF10*) and Butyrophilin Like 2 (*BTNL2*) have been known to be associated with IgE [89] and they both have CpG sites with stable DNA methylation.

Combining this finding with their associations with IgE at different ages, these 59 CpG sites have the potential to serve as epigenetic biomarkers for IgE changes over time.

3.6 Conclusion

Genome wide longitudinal assessment of 385,183 CpG sites with IgE production led to identification of number CpG sites. Among the identified CpG sites, 124 were replicated in an independent cohort for their longitudinal association with IgE production at later ages. Genes associated with these CpG sites were enriched in DAVID pathways and categories known to influence the IgE production. Most importantly 59 of 124 CpG sites have stable DNA-methylation across different time points. Thus, the identified CpG sites have the potential to serve as an epigenetic biomarker for IgE changes over time and can serve as candidates for future studies related to IgE production.

4 Assessment of methods for cell type correction in epigenome wide association study

4.1 Abstract

Background

Whole blood is frequently utilized in genome-wide association studies of DNA methylation patterns in relation to environmental exposures or clinical outcomes. These associations can be confounded by the cellular heterogeneity. Several algorithms have been developed to measure or adjust for this heterogeneity. However, it is unknown whether these approaches are consistent and if not, which method(s) perform better.

Results

Methods: We compared eight cell-type correction methods including the method implemented in the minfi R package, the method by Houseman et al., the Removing unwanted variation (RUV) approach, the methods implemented in FaST-LMM-EWASher, ReFACTor, RefFreeEWAS, and RefFreeCellMix R programs, along with one approach utilizing surrogate variables (SVAs). In the first comparison, we evaluated the association of DNA methylation at each CpG across the whole genome with prenatal arsenic exposure levels and with cancer status, adjusted for estimated cell-type information obtained from different methods. We then compared CpGs showing statistical significance from different approaches. For the methods implemented in minfi and proposed by Houseman et al., we utilized homogeneous data with composition of some blood cells available and compared them with the estimated cell compositions. Finally, for methods not explicitly estimating cell compositions, we evaluated their performance using simulated DNA methylation data with a set of latent variables representing “cell types”.

Results: Results from the SVA-based method overall showed the highest agreement with all other methods except for FaST-LMM-EWASher. Using homogeneous data, minfi provided

better estimations on cell types compared to the originally proposed method by Houseman et al. Further simulation studies on methods free of reference data revealed that SVA provided good sensitivities and specificities, RefFreeCellMix in general produced high sensitivities but specificities tended to be low when confounding present, and FaST-LMM-EWASher gave the lowest sensitivity but highest specificity.

Conclusions

Results from real data and simulations indicated that SVA is recommended when the focus is on the identification of informative CpGs. When appropriate reference data are available, the method implemented in the minfi package is recommended. However, if no such reference data are available or if the focus is not on estimating cell proportions, the SVA method is suggested.

4. 2 Background

Whole blood is frequently utilized in genome-wide association studies of DNA methylation patterns in relation to environmental exposures or clinical outcomes. However, for DNA methylation assessed from whole blood, the association between DNA methylation and an exposure of interest could be confounded by cellular heterogeneity [90, 91]. In larger epidemiological studies, it is not feasible to isolate and profile every individual cell subset. Thus, several algorithms have been developed to measure and adjust for cellular heterogeneity in whole blood.

Houseman et al. proposed a method to infer the cell mixture proportions based on a regression calibration technique, which uses an external validation dataset to calibrate the model and correct for the bias [38]. Their approach was specifically designed for the Illumina 27k beadchip [92]. Jaffe and Irizarry [37] modified the Houseman et al.'s algorithm and tailored it to predict cell mixture composition of DNA-methylation profiles obtained from Illumina 450k

beadchip (450K array; Illumina, Inc., San Diego, CA, USA). This cell type correction method is implemented in Bioconductor [93] package minfi [94]. The above two approaches require external validation datasets and are designed to identify cell mixtures in tissues such as whole blood.

Apart from these two reference-based techniques, non-reference-based methods have also been developed. An advantage is that these non-reference based methods can also be applied to any other tissue in addition to blood. Zou et. al developed a non-reference-based method FaST-LMM-EWASher. This method is built upon linear mixed models with top principal components as the covariates. RefFreeEWAS [95] and its recently improved version (RefFreeCellMix)[96] are another two non-reference-based methods. They both utilize singular value decompositions (SVDs) and extract latent subject and cell-specific effects, but RefFreeCellMix incorporated additional constraints and utilities aiming to reduce the occurrence of false positives. Surrogate variable analysis (SVA) [97], based on SVDs of residuals in linear regressions, uses permutations to identify statistically significant eigen-vectors and consequently infer potential confounding factors (surrogate variables). A Bioconductor package is available to estimate surrogate variables using this approach [97]. ReFACTor [98] is another method that is free of reference database and it is based on principal component analyses on a set of potentially informative CpG sites.

Removing unwanted variation (RUV) is an approach designed to estimate cell type heterogeneity and built upon factor analyses. This approach utilizes reference CpGs inferred from reference database, based on which factor analyses are conducted. The factors are then included in subsequent analyses for the purpose of adjusting for cell type effects. Although

reference database is needed, this method does not estimate cell type proportions as in the minfi package and in the Houseman et al. method.

Among all these methods (eight approaches in total), it was unknown whether the existing methods were comparable to one another, and if not, which method(s) might perform better. To this end, we applied each cell type correction method (Houseman et al., minfi, RUV, FaST-LMM-EWASher, ReFACTor, RefFreeEWAS, and RefFreeCellMix) as well as the surrogate variable analyses (SVA) to real data sets. We evaluated the association between genome-scale DNA methylation and a variable of interest adjusting for cell type compositions. Then we compared the methods with regard to agreement on the number of CpG sites identified as being statistically significant. In addition, for the methods implemented in the minfi package and proposed by Houseman et al., we utilized homogeneous data with some blood cells composition available and compared these with the estimated cell compositions. For methods free of reference groups (FaST-LMM-EWASher, RefFreeEWAS, RefFreeCellMix, ReFACTor, and SVA), we further utilized simulated data generated under different scenarios to compare different methods, which, combined with findings from the real data, enabled us to demonstrate the quality of each method.

4.3 Method

In this section we briefly describe the existing techniques for estimating cell proportions or inferring latent variables due to cell compositions, data sets (real and simulated) to assess these methods, and statistical methods used in the analyses. All the analyses were programmed in R and a tutorial website including all the programs demonstrating the methods is available at <https://akhilesh362.wordpress.com/>.

Existing methods for cell compositions

4.3.1 Reference-based methods

Houseman et al. [38] developed a method for cell type correction that capitalizes on the idea that differentially methylated regions (DMRs) can serve as a signature for the distribution of different types of white blood cells. It uses these DMRs as a surrogate in a regression calibration based technique to identify the cell mixture distribution. Regression calibration technique can lead to bias estimate, thus external validation data is used to calibrate the model and to correct for the bias [99]. Their method was specifically for the Illumina 27k beadchip array.

The method by Jaffe and Irizarry [37] was adapted from the Houseman et al. [38] method and is tailored for Illumina450k along with 27k array. The algorithm in Houseman et al. identifies 500 CpG sites used to estimate cell mixture proportions from the Illumina 27k array. The modification of Jaffe and Irizarry was motivated because of the existence of probe SNPs in the 500 CpG sites and the inconsistency of CpG sites between the 27k and 450k arrays. In addition, the flow-sorted data of the six adult male subjects were used as references [35].

The method of removing unwanted variation (RUV) uses information from reference database, but it does not estimate cell type proportions. Instead, this approach bases on the information of negative control probes and performs factor analysis on these probes to identify factors due to unmeasured confounders. These factors are then included in subsequent analyses to adjust for cell type effects. The negative control probes were chosen as top 500 CpG sites from the reference databases of DNA methylation known to be correlated with the cell types [100].

4.3.2 Reference-free methods

In total, four commonly used or recently developed reference-free methods are implemented in our study, FaST-LMM-EWASher, RefFreeEWAS, RefFreeCellMix, and ReFACTor. These methods do not need any external validation datasets and have the potential to adjust for cell mixture arising from any other tissue in addition to blood. FaST-LMM-EWASher [101] applies the maximum likelihood (ML) approach in linear mixed models and optimize spectral decomposition to estimate cell types [102]. RefFreeEWAS utilizes singular value decomposition (SVD) to decompose the residuals of unadjusted linear model along with unadjusted linear coefficient estimates, and estimate latent subject and cell-specific effects. Bootstrap estimates for coefficient standard errors are used to account for the correlation in the error structure.

Surrogate variable analysis (SVA) estimates potential confounding factors from a singular value decomposition (SVD) of residuals and was initially applied to gene expression data [103]. SVA utilizes the concept of expression heterogeneity while estimating surrogate variables. Expression heterogeneity (EH) refers to certain plausible biological profiles of the subject, which may not be captured by the covariates in study. Compared to the method in RefFreeEWAS, SVA decomposes the residual matrix and utilize permutations to identify statistically significant eigen-vectors which serve as a representative of EH (the so-called eigengenes), and then infer surrogate variables based on theses “eigengenes”. Surrogate variables from SVA have the potential to cover information on cell types in DNA methylation from blood cells.

The method built in the R package RefFreeCellMix is improved from that in RefFreeEWAS. It uses a variant of non-negative matrix factorization to decompose the total methylation sites into CpG-specific methylation states for a pre-specified number of cell types

and subject-specific cell-type distributions [96]. Another approach in the R package, ReFACTor, a variant form of principal component analysis (PCA) to adjust for the cell type effects. This method assumes that a small number of methylation sites are affected by underlying cell mixtures. It filters out CpGs if the variation is not large enough (the default cutoff is standard deviation=0.02). To avoid too many CpGs filtered out, in our analyses, we excluded CpGs such that their standard deviations were in the lower 5th percentile. By default this method searches for top 500 most informative methylation sites and performs PCA with a fixed number of components on these CpG sites to obtain the components. These ReFACTor components can be used as a covariate in epigenome wide association study or can be added one at time to remove the inflation due to cell type composition [98].

4.3.3 Three real data sets used to compare the approaches

These three data sets include data on prenatal arsenic exposure and DNA methylation, an example data from FaST-LMM-EWASher, and data on breast cancer status and DNA methylation. The first two data sets were utilized to demonstrate each of the five methods for cell type compositions and their agreement in terms of identified CpGs potentially associated with a variable of interest. The third data set was used to assess the agreement between the estimated cell type proportions (using the Houseman et al. method and the method in minfi) and the physical counts of the cells. This data set served as a benchmark and was critical for the comparison between the Houseman et al. method and the method in minfi. The benchmark data used to compare reference-free methods were simulated data, as discussed in the next section.

Prenatal arsenic exposure and DNA methylation data: The data were from a birth cohort study examining multiple prenatal and postnatal factors in relation to child health outcomes, part of the nationwide Taiwan prenatal and infant cohort study [78, 79] established in Taiwan in 2000-2001.

In total, 64 subjects with genome-scale DNA methylation and level of prenatal arsenic exposure were included in our study. DNA methylation data were pre-processed including quantile normalization, probe-type correction, and probe SNPs exclusion. After pre-processing, in total, 385,183 CpG sites were included in the analyses. All the five methods were applied to this data set. This and the following example data set were used to compare the performance of the five methods.

An example data from FaST-LMM-EWASher: This is an example data provided by the FaST-LMM-EWASher package [104]. It was originally used to illustrate the method incorporated into FaST-LMM-EWASher. In total, 204 subjects with cancer status and DNA methylation from Illumina 27K array on 25,978 CpG sites are available.

Breast cancer status and DNA methylation data: This data set has been previously described [105] and has genome-scale DNA methylation and breast cancer status available on 61 subjects at baseline and 39 subjects at six month follow-up along with complete blood counts. After pre-processing, 484,489 CpG sites were included in the study. In this article, we focus on granulocytes, monocyte and lymphocytes cells since proportions of these cells can be estimated by use of the minfi package and the original Houseman et al. approach. In our study, proportions of these cells from the physical counts were compared to the cell proportions estimated by minfi and the Houseman et al. method.

4.3.4 Simulated data sets to compare the methods

To further evaluate the three reference-free methods (FaST-LMM-EWASher, RefFreeEWAS, RefFreeCellMix, ReFACTor, and SVA), we simulated DNA methylation data under different settings with “latent” variables representing “cell types”. These data sets served as benchmark

data for comparing reference-methods because the underlying truth was known. Two simulation scenarios were employed to evaluate the methods.

Scenario 1: We simulated DNA methylation data at 2,000 CpG sites across 600 samples, of which the first n CpG sites were associated with covariates of interest (e.g., level of arsenic exposure) and a set of latent variables, and the remaining CpG sites were only associated with the latent variables. The set of latent variables represent “cell types”. One covariate of interest was considered and generated from a Normal distribution with mean 1 and variance 1 ($N(0, 1)$), The coefficients of this covariate was set at 0.3 and the intercept in the regressions was 0.5. Five “latent” variables were used and generated from five different Normal distributions: $N(0,5)$, $N(3,1)$, $N(0,1)$, $N(2,4)$, $N(0,3)$, respectively. The association of DNA methylation and the latent variables was assumed linear and the coefficients were generated from $N(0.5, 0.01)$. The distribution of random errors in the linear regressions was assumed to be Normal with mean 0 and variance 1.2 for the n CpGs, mean 0 and variance of 1.2 for the next 100 CpGs, and mean 0 and variance 2 for the remaining CpGs. The last setting with larger variance in random errors was for situations that the influence of cell types on DNA methylation was weaker.

We took three values of n , $n=50$, 100, and 150, representing different sparsity levels (from high to low) of informative CpGs. In total, 100 data sets for each n were simulated. Note that under this scenario, the covariates and latent variables were generated separately and had no correlations.

Scenario 2: Latent variables generated under this scenario have potential confounding effects. The overall setting is the same as in Scenario 1, except that the covariate of interest and the five latent variables (6 variables in total) were correlated such that correlation is equal to 0.7^{i-1} .

j , $i, j = 1, 2, 3, 4, 5, 6$. For instance, the correlation of the continuous covariate with the first latent variable was 0.7, and with the second latent variable was $0.7^2=0.49$.

4.3.5 Statistical analyses

Linear regression-based analyses were used to assess the associations of DNA methylation with variables of interest with cell type heterogeneity adjusted using eight different methods. In the analyses of the two real data sets (the arsenic and DNA methylation data, and the FasT-LMM-EWASher example data), we recorded CpG sites showing statistically significant association with variables of interest (i.e., arsenic exposure and cancer) after implementing different cell type heterogeneity inference methods. We also inferred the number of statistically significant CpG sites without adjusting for cell type heterogeneity. To compare the eight cell type heterogeneity inference methods (Houseman et al., minfi, FaST-LMM-EWASher, RefFreeEWAS, RefFreeCellMix, ReFACTor, RUV, and SVA), we assessed the percentage of overlap between different methods in the number of identified CpG sites that showed statistical significance, and calculated a similarity index, Jaccard index (J-index) [106]. The percentage of overlap is calculated as the number of identified CpGs overlapped with that from SVA divided by the number of CpGs identified by SVA. We used Fisher exact test to assess the significance of overlap. Jaccard index measures the similarity between two finite sample sets. We used a Bioconductor package GeneOverlap to calculate this index. To assess whether the CpGs uniquely identified by the SVA approach are informative, we used the Database for Annotation, Visualization and Integrated Discovery (DAVID) [45, 46] to analyze the enrichment in Gene ontology (GO) [107] categories and Kyoto Encyclopedia of Genes and Genomes (KEGG) [40, 108] pathways.

As for each simulated data set, we calculated sensitivity and specificity of the selected CpG sites for each cell type heterogeneity inference method. They were calculated by comparing the detected CpGs with the truly important CpGs. For each of the five methods, median of sensitivity and specificity along with 95% empirical intervals across 100 data sets were recorded for each setting under each simulation scenario.

4.3 Results

4.3.1 Findings from prenatal arsenic exposure and DNA-methylation data

We used genome-scale DNA methylation data from a birth cohort study consisting of 64 cord blood samples examining multiple prenatal factors in relation to child health outcomes, pilot of the nationwide Taiwan Maternal and Infant Cohort Study [78, 79].

We assessed the association of DNA methylation at each CpG site across the whole genome with prenatal urinary arsenic exposure levels (a continuous measure), adjusting for cell-type effects with cell type information inferred from one of the eight methods. For each method, the number of CpGs was recorded showing statistically significant associations with prenatal urinary arsenic exposure after adjusting for multiple testing by controlling false discovery rate (FDR) at 0.05. ReFACTor identified the largest number of CpGs (~60,000) and no CpGs were detected by FaST-LMM-EWASher (Table 4.1). RefFreeCellMix also identified a large number of CpGs (~3000). SVA and RefFreeEWAS detected more CpGs compared to the remaining methods. (Table 4.1). Next, we assessed the number of identified CpGs that overlap between different methods. The diagram in Figure 4.1 shows the overlap of CpG sites from four approaches (Houseman et al., minfi, RefFreeEWAS, and SVA) as well as the analyses without adjusting for cell types. Results from SVA showed the best agreement with findings from the other four analyses (Figure 4.1). Two identified CpG sites cg06434480 and cg10662395 were

common to all these five analytical methods labeled in Figure 4.1. Further comparisons indicated that CpG site cg10662395 was also identified by RefFreeCellMix and RUV, and this is the only CpG site overlapped among all the seven analyses (Houseman et al., minfi, RefFreeEWAS, SVA, RefFreeCellMix and RUV, as well as the analyses without adjusting for cell types). Although ReFACTor identified the largest number of CpGs, they did not overlap with the joint findings from the aforementioned seven analyses. Overall, CpGs identified via SVA overlapped with those from the Houseman et al. method, minfi and RefFreeEWAS (p-value<0.0001, Table 4.1, Figure 4.1. The definition of percentage overlap is given in the Methods section). One of the two CpGs (cg06434480 and cg10662395), cg06434480 is located within 200 base pairs of transcription start site of gene *HMGCR* (3-hydroxy-3-methylglutaryl-CoA reductase) known to be associated with inorganic arsenic exposure [109]. While in a study conducted in humans Mono-methylated arsenic (MMA) it was found to downregulate the gene expression of *HMGCR*, a gene involved in cholesterol biosynthesis [110]. The other CpG cg10662395 is located in the body region of gene *HCN2* (hyperpolarization activated cyclic nucleotide gated potassium channel 2). This gene was not found to be directly associated with arsenic exposure in the literature, but *HCN2* has been known to regulate pacemaker activity in the heart and the brain of mouse and human [111, 112], and arsenic has been found to induce QT interval (i.e., time between initial deflection of QRS complex to the end of T wave) prolongation probably by altering potassium ion channel [113].

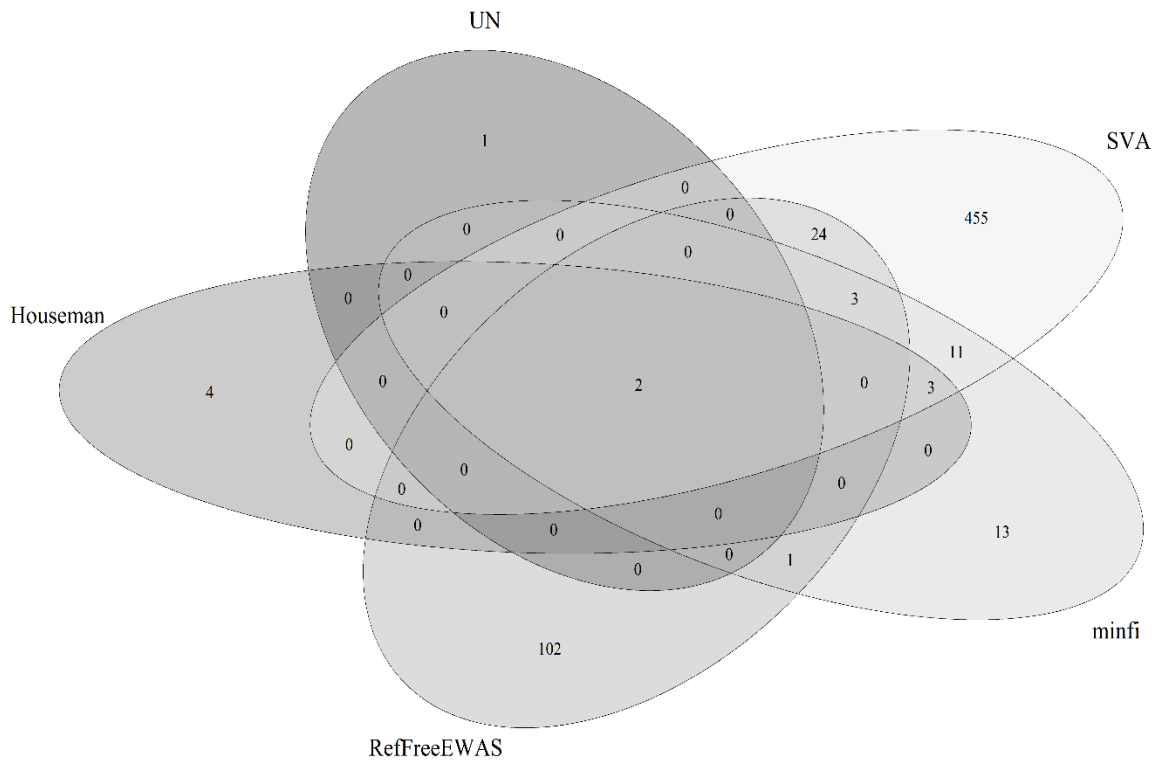


Figure 4.1. Venn diagram illustrating the overlap of identified CpG sites that are associated with prenatal arsenic exposure at FDR level of 0.05 after incorporating estimated cell type compositions by different methods for the association study of prenatal arsenic exposure with DNA-methylation. “UN”: results from an analysis without adjusting for cell type compositions.

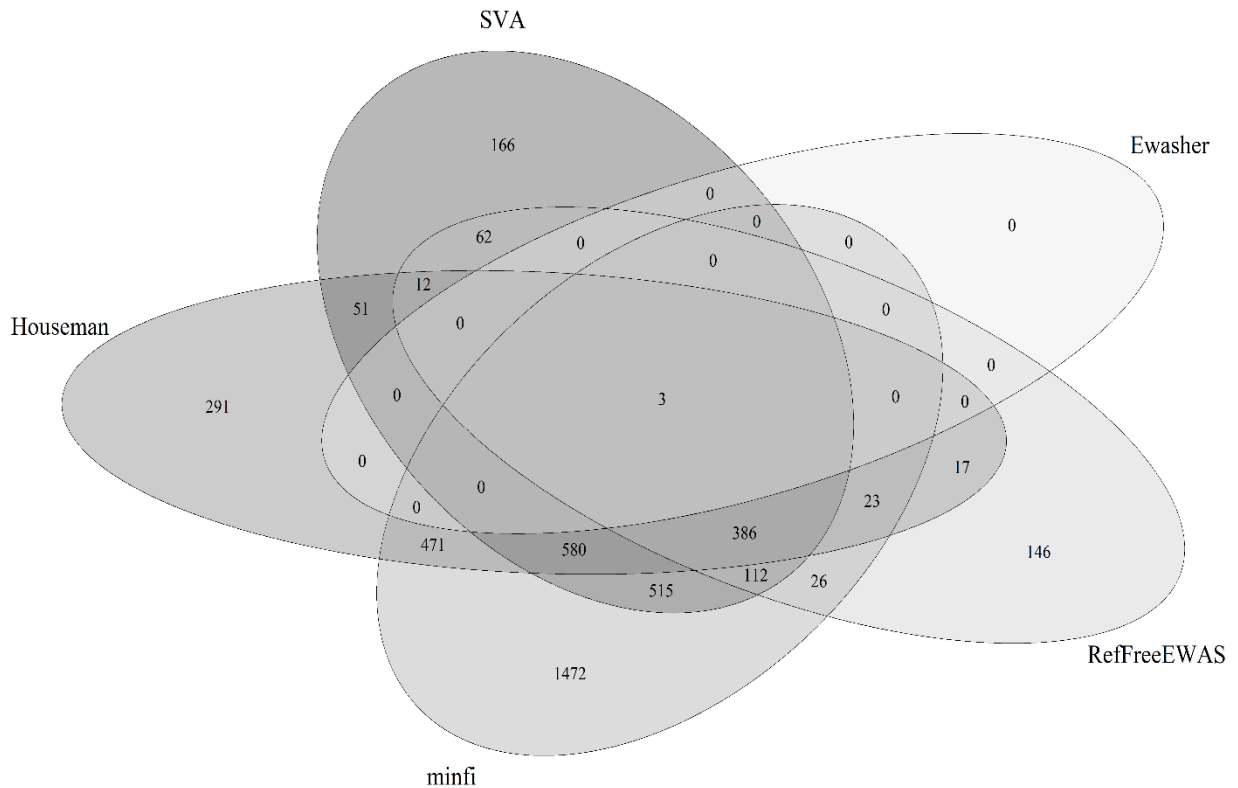


Figure 4.2. Venn diagram illustrating the overlap of identified CpG sites that are associated with cancer status at FDR level of 0.05 after incorporating estimated cell type compositions by different methods for the association study of cancer status with DNA-methylation.

The motivation of adjusting for cell types was due to the potential confounding effects of cell type compositions with respect to the association of arsenic exposure with DNA methylation, caused by the association of arsenic exposure with cell type compositions [114-117]. Our assessment on the correlations between total arsenic exposure and estimated cell type proportions also supported the potential confounding-role of cell types (Appendix Figure A4.1). To support the existence of such confounding effects, we assessed the associations with and without adjusting for cell type proportions at all CpG sites. We found that at more than 99% of all the CpGs the effects (regression coefficients) of prenatal arsenic exposure changed by more than

10% from cell type unadjusted (the median of the coefficients was 2.32 with 5th percentile of 0.40 and 95th percentile of 3.46) to cell types adjusted (the corresponding statistics were 0.080, 0.0073 and 0.25), indicating a need of adjusting for cell types.

Overall, the analysis based on SVA identified CpG sites that had better overlap with the CpGs identified by other methods. To acquire the biological relevance of CpGs uniquely identified by use of SVA, we implemented DAVID to perform Gene Ontology (GO) analysis and to identify KEGG pathways. The 455 (out of 498) significant CpG identified uniquely by SVA were mapped to genes using Illumina annotation file for 450K DNA methylation array. Of great interest, GO categories related to transcription and regulation of RNA metabolic process were enriched after controlling FDR at 0.05, as well as three KEGG pathways, endocytosis, cancer pathway and MAPK signaling pathway. A discussion on the connection of arsenic exposures and the identified GO categories and KEGG pathways is presented in the Discussion section.

Table 4.1. Number of significant CpG sites with and without cell type correction and overlap with the SVA method (data on prenatal arsenic exposure and DNA methylation).

Method	Identified CpGs (N) [#]	Overlap with SVA (%)	p-value ^{##}
Houseman et al.	10	1.20	<0.0001
minfi	57	4.62	<0.0001
SVA	498	---	---
RefFreeEWAS	133	6.01	<0.0001
RefFreeCellMix	2,932	0.60	1.0
ReFACTor	58,871	13.03	1.0
EWASher*	0	0.0	---
RUV	356	0.20	1.0
Unadjusted**	3	0.60	<0.0001

* The FasT-LMM-EWASher method.

** Unadjusted: cell type compositions were not included in the analyses.

[#]The selection of CpG sites is based on FDR-adjusted p-values (FDR is controlled at 0.05).

^{##}P-value is based on Fishers exact test for overlap with results from SVA. The null hypothesis is that there is no overlap with the CpGs identified based on SVA.

4.3.2 Findings from example data

We repeated the same analysis on an example data set provided by the FasT-LMM-EWASher package. A tutorial website for applying all the cell type composition inference methods to this example data is available at <https://akhilesh362.wordpress.com/>. This data set includes DNA methylation from the Illumina 27K array and measures of a binary variable (cancer status) for 204 subjects. In total, 7,648 CpGs were included in our study based on initial screening done by the FasT-LMM-EWASher package. In this example data, cell type proportions were likely to be different on average between subjects with cancer and those without cancer, based on two-sample t-tests applied to logit-transformed sample proportions, explaining the potential need to adjust for their confounding effects. Since Illumina 27K focuses more on cancer genes, DNA

methylation at a large number of CpG sites showed statistically significant associations with cancer status (Table 4.2). Some similar findings as in Table 4.1 were observed. ReFACTor identified a large number of CpGs, Fast-LMM-EWASher identified the least number of CpG sites, and SVA agreed nicely with minfi (Jaccard similarity index=0.4). A unique observation from this analysis is that RUV identified the largest number of CpGs (6,008 CpGs, close to the number of CpGs in the candidate pool, 7,648 CpGs). Since the original Houseman et al. method was designed specifically for Illumina 27K platform, it is reasonable that SVA also showed a large overlap with results from this approach (Jaccard similarity index=0.4). In total, 3 identified CpGs (cg22029275 located in the 1st Exon of *FAM123A* gene, cg07080358 located in 1st Exon of *CNRIP1*, and cg15202954 located within 200 base pair of transcription start site of *NALCN* gene) were common to all the eight cell correction methods as well as to the analyses without cell type composition adjusted. There is evidence that these three genes (*FAM123A*, *CNRIP1* and *NALCN*) are associated with the risk of colorectal cancer [118-120].

Table 4.2. Number of significant CpG sites with and without cell-correction methods and overlap of CpG sites with those from the SVA method (example data from FasT-LMM-EWASher package).

Method	Identified CpGs (N)[#]	Overlap with SVA (%)	p-value^{##}	J-index^{###}
Houseman et al.	1,835	54.71	<0.0001	0.40
minfi	3,589	84.59	<0.0001	0.40
SVA	1,888	---	---	---
RefFreeEWAS	788	30.51	<0.0001	0.30
RefFreeCellMix	1,006	18.38	<0.0001	0.10
ReFACTor	4,224	87.45	<0.0001	0.40
EWASher*	3	0.16	<0.0001	0
RUV	6,008	99.95	<0.0001	0.30
Unadjusted**	3,768	82.89	<0.0001	0.40

* The FasT-LMM-EWASher method.

** Unadjusted: cell type compositions were not incorporated into the analyses.

[#]The selection of CpG sites is based on FDR-adjusted p-values (FDR is controlled at 0.05).

^{##} P-value is based on Fishers exact test for overlap. The null hypothesis is that there is no overlap with the CpGs identified based on SVA.

^{###} J-index is Jaccard index.

DAVID analysis of genes associated with the significant CpGs identified uniquely by SVA led to the identification of three GO categories related to plasma membrane at FDR of 0.05 (integral to plasma membrane, intrinsic to plasma membrane, and plasma membrane part), as well as KEGG pathways such as pathways in cancer and signaling pathways, which indicates that genes corresponding to these CpG sites may play a role in the regulation of cancer.

4.3.3 Findings from breast cancer status and DNA-methylation data

This analysis uses a data set discussed in Smith et al. [105]. Breast cancer status, DNA-methylation, and cell counts for granulocytes, monocyte, and lymphocytes for 61 subjects at baseline and a subset of 39 subjects at six months follow up are implemented in the analyses. Among all the methods discussed, the method implemented in the minfi package and the original Houseman et al. method are able to estimate cell proportions. We used minfi and the Houseman et al. approach to estimate the proportions of granulocyte, monocyte and lymphocyte cells. Lymphocytes proportion were derived by adding the proportions of B cell, T cell and Natural Killer (NK) cells. For the three cell types (granulocyte, monocyte and lymphocyte), Pearson correlations between estimated (minfi) and true cell proportions were 0.85, 0.79, 0.88 at baseline and 0.84, 0.78, 0.87 at the six month follow up, respectively. For the correlations based on the Houseman et al. method, they were 0.84, 0.78 and 0.88 at baseline and 0.78, 0.73 and 0.83 at the six month follow up, respectively. All the correlations showed statistically significant difference from zero (p -value <0.05).

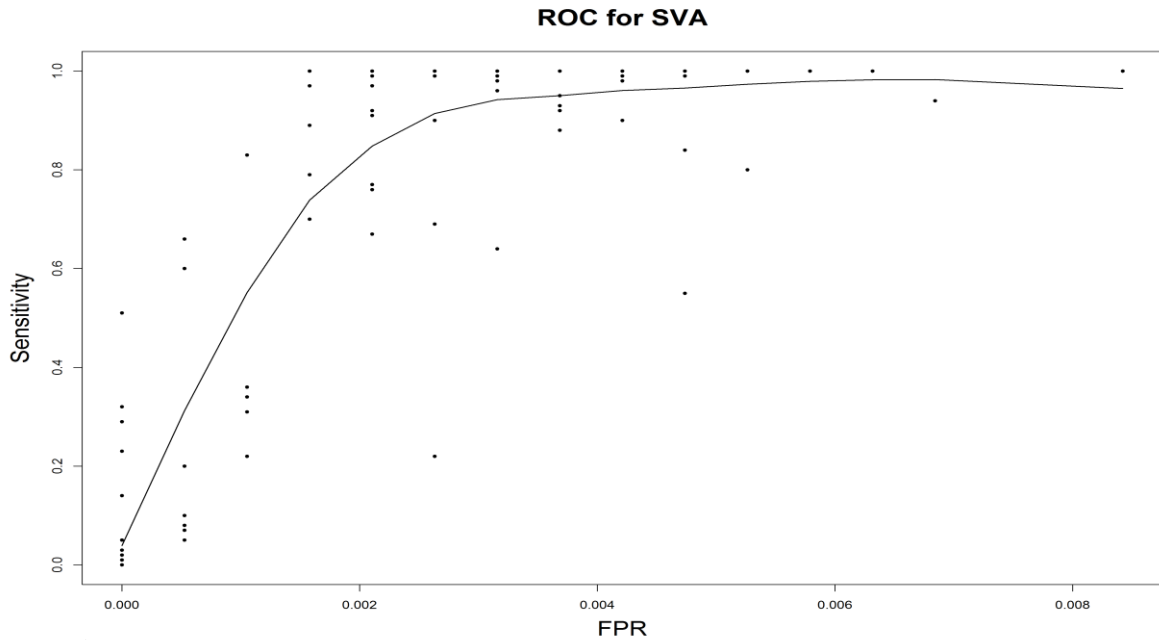
4.3.4 Findings from simulated data

We simulated data applying two scenarios with the first scenario focusing on latent variable effects (comparable to effects of cell composition), and the second focusing on latent variable effects with confounding (comparable to effects of cell composition as well as confounding effects). In total, 100 data sets were simulated under each scenario. Details of the simulation scenarios are given in the Methods section. The simulated data were used to evaluate the five methods that do not estimate cell proportions nor need reference databases, specifically, FaST-LMM-EWASher, RefFreeEWAS, RefFreeCellMix, ReFACTor, and SVA.

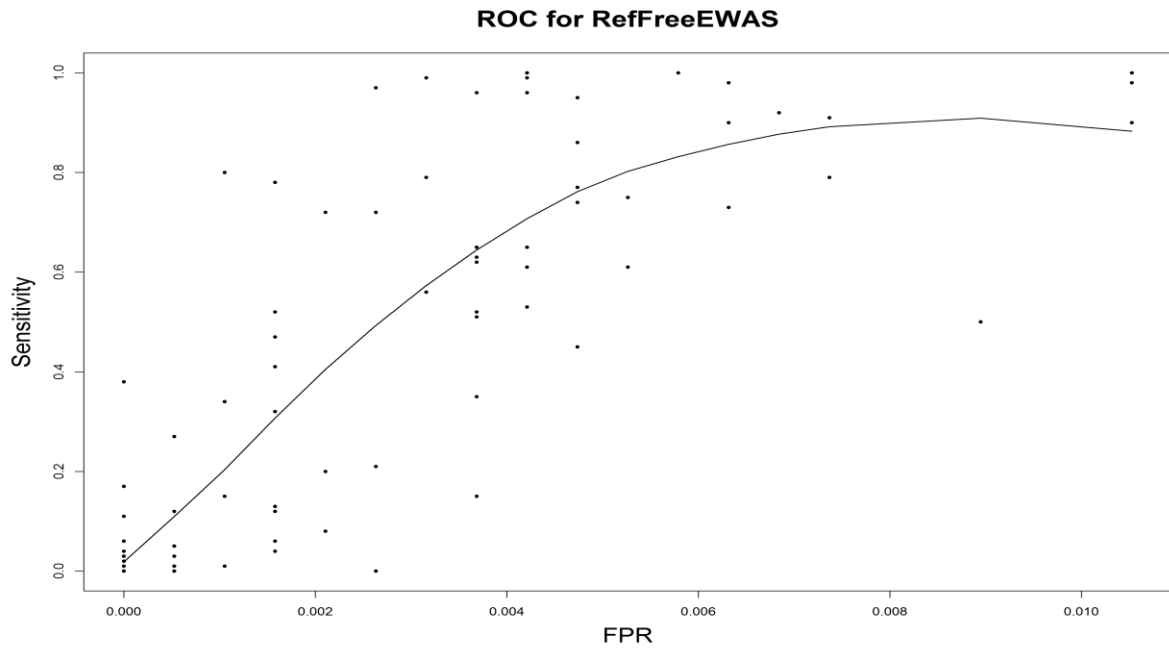
For data under all scenarios, we applied each of the five methods to each simulated data to draw information on cell compositions. We then incorporated the information to assess the associations of “DNA methylation” with the variable of interest at each pseudo CpG site, and compared each method by assessing the sensitivity and specificity of the selected CpG sites across all 100 data sets. Regardless of the number of important CpGs, FaST-LMM-EWASher resulted in the lowest sensitivity but the highest specificity for both scenarios, consistent with findings from real data (Table 4.3). Findings from RefFreeEWAS, RefFreeCellMix, ReFACToR, and SVA are, in general, comparable for data simulated under scenario 1, but SVA gives consistently higher sensitivity and specificity in all settings (Table 4.3). For data simulated under scenario 2 with high correlations ($\rho=0.7$), SVA outperformed FaST-LMM-EWASher, RefFreeEWAS, RefFreeCellMix and ReFACToR and had higher sensitivity and specificity. Compared with RefFreeEWAS, overall RefFreeCellMix outperformed RefFreeEWAS when confounding effects present, showing much higher sensitivities with relatively lower specificities. Results from ReFACToR indicated extremely low specificity under scenario 2, which is consistent with the rather large numbers of CpGs identified in real data. The performance of FaST-LMM-EWASher was similar between the two scenarios and was inferior to all other methods. On the other hand, the SVA method performed well under both scenarios, followed by RefFreeEWAS and RefFreeCellMix with RefFreeEWAS being weaker in capturing confounding effects. We also considered a situation with $\rho=0.3$, mimicking a situation of moderate confounding, and similar patterns observed as those from the relatively two extreme cases ($\rho=0$ and $\rho=0.7$).

In the above simulations, we fixed the regression coefficients of the important CpGs. To demonstrate the pattern of sensitivity and specificity, we implemented receiver operating

characteristic (ROC) plots. In total, 100 data sets were simulated under scenario 1 with regression coefficients for the variable of interest ranged from 0.01 to 0.3. For each data set, we calculated sensitivity and specificity of selected CpGs, based on which we estimated the ROC curves. Sensitivities from FaST-LMM-EWASher were substantially low and were not considered in this demonstration. The performance of RefFreeEWAS, RefFreeCellMix, and ReFACTor was comparable under scenario 1 (Table 4.3). We therefore only presented ROC curves for ReFreeEWAS and SVA for the purpose of comparison (Figure 4.2). The findings are consistent with what we observed from Table 4.3 for scenario 1, that is, SVA performed better than RefFreeEWAS. In addition, the results indicated that both SVA and RefFreeEWAS have high specificity regardless of the underlying regression coefficients, indicating the conservatism when selecting informative CpGs.



a)



b)

Figure 4.3. Plots of sensitivity v.s. 1-specificity and estimated ROC curves, a) SVA. b) RefFreeEWAS

Table 4.3. Summary of sensitivity, specificity of FaST-LMM-EWASher, RefFreeEWAS, RefFreeCellMix, ReFACTor, and SVA for 100 simulated data across three settings.

	Sensitivity (Median, 95% interval)		Specificity (Median, 95% interval)	
Number of Important CpGs =50				
	Scenario 1 ($\rho =$ 0)	Scenario 2 ($\rho =$ 0.7)	Scenario 1 ($\rho =$ 0)	Scenario 2 ($\rho =$ 0.7)
Ewasher ^a	0.00 (0.00, 0.00)	0.00 (0.00, 0.00)	1.00 (1.00, 1.00)	1.00 (1.00, 1.00)
RefEWASRef ^b	1 (0.96, 1.00)	0.00 (0.00,0.49)	1.00 (0.99,1.00)	0.58 (0.06,1.00)
CellMix ^c	1.00 (0.98, 1.00)	1.00 (1.00, 1.00)	1.00 (0.99, 1.00)	0.55 (0.20, 0.92)
ReFACTor	1.00 (0.96, 1.00)	1.00 (1.00, 1.00)	1.00 (0.83, 1.00)	0 (0.00, 0.00)
SVA ^d	1.00 (0.98, 1.00)	1.00 (0.96, 1.00)	1.00 (0.996, 1.00)	1.00 (0.996, 1.00)
Number of Important CpGs =100				
Ewasher ^a	0.00 (0.00, 0.00)	0.00 (0.00, 0.00)	1.00 (1.00, 1.00)	1.00 (1.00, 1.00)
Ref ^b RefEWAS	1.00 (0.97,1.00)	0.00 (0.00,0.40)	1.00 (0.99, 1.00)	0.52 (0.01,1.00)
CellMix	1.00 (0.98, 1.00)	1.00 (1.00, 1.00)	0.99 (0.97, 1.00)	0.21 (0.05, 0.53)
ReFACTor	1.00 (0.97, 1.00)	1 (1.00, 1.00)	0.99 (0.81, 1.00)	0.00 (0.00, 0.00)
SVA ^c	1.00 (0.99,1.00)	0.99 (0.97, 1.00)	1.00 (0.99, 1.00)	1.00 (0.99, 1.00)
Number of Important CpGs =150				
Ewasher	0.00 (0.00,0.00)	0.00 (0.00,0.00)	1.00 (1.00,1.00)	1.00 (1.00,1.00)
RefRefEWAS	0.99 (0.97,1.00)	0.00 (0.00,0.29)	0.99 (0.99,1.00)	0.50 (0.01,1.00)
CellMix	1.00 (0.98, 1.00)	1.00 (1.00, 1.00)	0.98 (0.93, 0.99)	0.10 (0.02, 0.29)
ReFACTor	1.00 (0.98, 1.00)	1.00 (1.00, 1.00)	0.99 (0.79, 1.00)	0.00 (0.00, 0.00)
SVA	1.00 (0.99,1.00)	0.99 (0.97, 1.00)	0.99 (0.99,1.00)	0.99 (0.99, 1.00)

Footnote: ρ = correlation between primary covariate and latent variables. $\rho = 0$ corresponds to data simulated from Scenario 1, while $\rho=0.7$ corresponds to data simulated from Scenario 2. a= FaST-LMM-EWASher, b= RefFreeEWAS, c= RefFreeCellMix, d=Surrogate variable analysis.

Discussion

We compared eight cell-type correction methods using real and simulated data. Based on DNA methylation in a cohort study, the methods in ReFACTor identified the largest number of CpGs (~60K CpGs), none of which overlapped with the common CpGs detected by other methods including the analysis without adjusting for cell type compositions (but excluding the method in FaST-LMM-EWASher). The method in FaST-LMM-EWASher did not identify any CpG sites. Except for ReFACTor and FaST-LMM-EWASher, at least one detected CpG was shared between all the other methods. More than 50% of CpGs identified by the Houseman et al. method and by the approach implemented in minfi were also detected by the SVA method; The overlap in CpGs was much less between these two methods and the remaining methods. The genes associated with CpGs uniquely identified by using SVA with prenatal urinary arsenic as primary exposure led to the enrichment of GO categories and KEGG pathways that were consistent with our understanding with respect to the effect of arsenic on DNA methylation. Arsenic exposure leads to generation of reactive oxygen species (ROS) which induces DNA damage [121]. This reactive oxygen species play a crucial role in signal transduction pathways, transcription factor regulation [122], and mitogen activated protein kinases (MAPKs) signal transduction pathway is one such pathway that is affected by ROS [123]. DAVID analysis of genes associated with the CpGs uniquely identified by SVA for FaST-LMM-EWASher example dataset led to enrichment of KEGG pathways in cancer. All these imply that the CpGs uniquely

identified by using SVA are potentially informative. Using the example dataset provided by FaST-LMM-EWASher method, we found that all methods except for FaST-LMM-EWASher identified a large number of CpG sites. This was likely due to the platform used to measure DNA methylation levels (Illumina 27K), which is centered more on cancer genes. However, CpGs identified based on ReFACTor and RUV were close to the number of CpGs in the pool of candidate CpGs, indicating possible inflations. On the other hand, results from minfi showed the greatest overlap with the SVA method (Table 4.2). Based on these two real data sets, results from the method in the minfi package and those from SVA were most agreeable. However, for real data, the underlying truth was unknown, which was the motivation of incorporating a data set with cell counts known and the use of a series of simulation studies. Findings from these data were further discussed in this section.

Using the available cell counts in the cancer status and DNA methylation dataset we observed agreements between cell types estimated by Houseman et al. and minfi, but minfi showed a better agreement. The Houseman et al. approach was designed for the Illumina 27K beadchip array, which may not fit the 450K array as noted in the literature [37]. The modification of the Houseman et al. approach implemented in the minfi package, on the other hand, is suitable for both 27K and 450K array. The reference data were from six adult white European males. It has been shown that DNA methylation patterns vary by sex, age and ancestry [124-128]. Generalizing the cell mixtures estimated by minfi to studies with both genders and non-Europeans of different age groups may potentially introduce bias.

Further simulations investigating reference-free methods supported the findings from real data. Regardless of the number of important CpGs, FaST-LMM-EWASher showed the lowest sensitivity, indicating low power to identify truly important CpGs if using that method to adjust

for cell type compositions. ReFACTor produced lowest specificity when confounding effects were present, supporting the rather low overlapping with findings from other methods. On the other hand, findings from ReFACTor, RefFreeEWAS, RefFreeCellMix and SVA were in general comparable for data simulated under scenario 1 (no-confounding effects), but SVA gave consistently higher sensitivity and specificity when confounding effects present.

The SVA approach does not provide estimates on cell type compositions; however, our ultimate goal was not to estimate cell counts. The goal was to identify an approach that best assesses DNA methylation differentiation due to exposure or diseases, corrected for a potential cell type bias. From this viewpoint and the findings from real data and the high sensitivities and specificities from simulations (under both scenarios, confounding and no confounding), using SVA to adjust for cell type compositions seems to be an appropriate method and may perform better than the existing methods. It is worth noting that information included in the surrogate variables produced by the SVA method may also include other information in addition to cell type compositions. There is a potential of over-adjustment by use of this approach. Furthermore, we would like to point out that all these reference-free methods can be directly applied to genome-wide bisulfite sequencing data and we expect similar findings in terms of their ability in inferring cell type compositions.

Conclusion

When appropriate reference data are available and if inferences on cell type compositions are needed, the method implemented in the minfi package is recommended. However, if no such reference data are available or if the focus is not on estimating cell proportions, the SVA method is suggested to correct for bias resulting from varying cell mixtures.

5 Summary

The work presented in this dissertation is distributed in the ongoing project related to epigenome wide association studies. The epigenetic markers identified in this study will benefit researchers in epigenetic studies in that they will help elucidate the pathophysiology of disease associated with in utero arsenic exposure, and provide insight into the production of immunoglobulin E. The replications of identified CpG sites in independent cohorts strengthen the validity of the findings. I summarize the highlights of the work as follows:

1. The 252 CpG sites identified and replicated in an independent cohort can serve as an epigenetic marker for the adverse health effect of in utero arsenic exposure on the newborn subjects.
2. Of the 252 CpG sites, 5 CpG sites were found to be longitudinally associated with low density lipoprotein measures of the subjects at ages 2, 5, 8, 11 and 14. The genes corresponding to these five CpGs (cg25189764, cg04986899, cg04903360, cg08198265 and cg10473311) also had literature support for their association with cardiovascular disease or diabetes. The DNA methylation measurements at the identified five CpG sites are known to be stable across different age group. Thus, these five CpG sites have the potential to serve as an epigenetic marker for the adverse effect of in-utero arsenic exposure and its association with LDL measures at later ages.
3. The 124 CpG sites longitudinally associated with IgE in the main cohort and replicated in an independent cohort can serve as an epigenetic marker predicting the production of IgE. Of 124 the DNA methylation measures at 59 CpG sites were found to be stable at birth, age 10 and age 18 in IoW cohort. Thus, these 59 CpG sites are more reliable to serve as an epigenetic marker explaining the production of IgE and could eventually help in revealing the pathophysiology of the developmental immune based disease.

4. In the assessment of the methods for cell type adjustment we found that Houseman's algorithm implemented in Bioconductor package "minfi" is the best choice if reference dataset is available. Although R package "sva" performed best compared to other methods in most of the situations, but given the possibility of over fit it is recommended in the situation where reference dataset is not available.

References:

1. Farzan SF, Karagas MR, Chen Y: **In utero and early life arsenic exposure in relation to long-term health and disease.** *Toxicology and applied pharmacology* 2013, **272**(2):384-390.
2. Farzan SF, Li Z, Korricks SA, Spiegelman D, Enelow R, Nadeau K, Baker E, Karagas MR: **Infant Infections and Respiratory Symptoms in Relation to in Utero Arsenic Exposure in a U.S. Cohort.** *Environmental health perspectives* 2016, **124**(6):840-847.
3. Guan H, Piao F, Zhang X, Li X, Li Q, Xu L, Kitamura F, Yokoyama K: **Prenatal exposure to arsenic and its effects on fetal development in the general population of Dalian.** *Biological trace element research* 2012, **149**(1):10-15.
4. O'Sullivan L, Combes AN, Moritz KM: **Epigenetics and developmental programming of adult onset diseases.** *Pediatric nephrology* 2012, **27**(12):2175-2182.
5. Youngson NA, Whitelaw E: **Transgenerational epigenetic effects.** *Annual review of genomics and human genetics* 2008, **9**:233-257.
6. Chong S, Whitelaw E: **Epigenetic germline inheritance.** *Current opinion in genetics & development* 2004, **14**(6):692-696.
7. Egger G, Liang G, Aparicio A, Jones PA: **Epigenetics in human disease and prospects for epigenetic therapy.** *Nature* 2004, **429**(6990):457-463.
8. Holliday R, Pugh JE: **DNA modification mechanisms and gene activity during development.** *Science* 1975, **187**(4173):226-232.
9. Yoder JA, Walsh CP, Bestor TH: **Cytosine methylation and the ecology of intragenomic parasites.** *Trends in genetics : TIG* 1997, **13**(8):335-340.
10. Bird A: **DNA methylation patterns and epigenetic memory.** *Genes & development* 2002, **16**(1):6-21.
11. Foley DL, Craig JM, Morley R, Olsson CA, Dwyer T, Smith K, Saffery R: **Prospects for epigenetic epidemiology.** *American journal of epidemiology* 2009, **169**(4):389-400.
12. Wang SC, Oelze B, Schumacher A: **Age-specific epigenetic drift in late-onset Alzheimer's disease.** *PLoS one* 2008, **3**(7):e2698.
13. Dolinoy DC, Weidman JR, Waterland RA, Jirtle RL: **Maternal genistein alters coat color and protects Avy mouse offspring from obesity by modifying the fetal epigenome.** *Environ Health Perspect* 2006, **114**(4):567-572.
14. Dietert RR, Dietert JM: **Potential for early-life immune insult including developmental immunotoxicity in autism and autism spectrum disorders: focus on critical windows of immune vulnerability.** *Journal of toxicology and environmental health Part B, Critical reviews* 2008, **11**(8):660-680.
15. Bailey KA, Wu MC, Ward WO, Smeester L, Rager JE, Garcia-Vargas G, Del Razo LM, Drobna Z, Styblo M, Fry RC: **Arsenic and the epigenome: interindividual differences in arsenic metabolism related to distinct patterns of DNA methylation.** *Journal of biochemical and molecular toxicology* 2013, **27**(2):106-115.
16. Broberg K, Ahmed S, Engstrom K, Hossain MB, Jurkovic Mlakar S, Bottai M, Grander M, Raqib R, Vahter M: **Arsenic exposure in early pregnancy alters genome-wide DNA methylation in cord blood, particularly in boys.** *Journal of developmental origins of health and disease* 2014, **5**(4):288-298.
17. **A review of human carcinogens. Part C: Arsenic, metals, fibres, and dusts / IARC Working Group on the Evaluation of Carcinogenic Risks to Humans.** *International Agency for Research on Cancer* 2012, **100C**.
18. Vahter M: **Effects of arsenic on maternal and fetal health.** *Annual review of nutrition* 2009, **29**:381-399.

19. Nordstrom DK: **Public health. Worldwide occurrences of arsenic in ground water.** *Science* 2002, **296**(5576):2143-2145.
20. Smith AH, Marshall G, Liaw J, Yuan Y, Ferreccio C, Steinmaus C: **Mortality in young adults following in utero and childhood exposure to arsenic in drinking water.** *Environmental health perspectives* 2012, **120**(11):1527-1531.
21. Chou WC, Chung YT, Chen HY, Wang CJ, Ying TH, Chuang CY, Tseng YC, Wang SL: **Maternal arsenic exposure and DNA damage biomarkers, and the associations with birth outcomes in a general population from Taiwan.** *PloS one* 2014, **9**(2):e86398.
22. Rosenberg HG: **Systemic arterial disease and chronic arsenicism in infants.** *Archives of pathology* 1974, **97**(6):360-365.
23. Hawkesworth S, Wagatsuma Y, Kippler M, Fulford AJ, Arifeen SE, Persson LA, Moore SE, Vahter M: **Early exposure to toxic metals has a limited effect on blood pressure or kidney function in later childhood, rural Bangladesh.** *International journal of epidemiology* 2013, **42**(1):176-185.
24. Yuan Y, Marshall G, Ferreccio C, Steinmaus C, Selvin S, Liaw J, Bates MN, Smith AH: **Acute myocardial infarction mortality in comparison with lung and bladder cancer mortality in arsenic-exposed region II of Chile from 1950 to 2000.** *Am J Epidemiol* 2007, **166**(12):1381-1391.
25. Davila-Esqueda ME, Morales JM, Jimenez-Capdeville ME, De la Cruz E, Falcon-Escobedo R, Chi-Ahumada E, Martin-Perez S: **Low-level subchronic arsenic exposure from prenatal developmental stages to adult life results in an impaired glucose homeostasis.** *Experimental and clinical endocrinology & diabetes : official journal, German Society of Endocrinology [and] German Diabetes Association* 2011, **119**(10):613-617.
26. Rossman TG, Klein CB: **Genetic and epigenetic effects of environmental arsenicals.** *Metallomics : integrated biometal science* 2011, **3**(11):1135-1141.
27. Gluckman PD: **Epigenetics and metabolism in 2011: Epigenetics, the life-course and metabolic disease.** *Nature reviews Endocrinology* 2012, **8**(2):74-76.
28. Vickers MH: **Early life nutrition, epigenetics and programming of later life disease.** *Nutrients* 2014, **6**(6):2165-2178.
29. Majumdar S, Chanda S, Ganguli B, Mazumder DN, Lahiri S, Dasgupta UB: **Arsenic exposure induces genomic hypermethylation.** *Environmental toxicology* 2010, **25**(3):315-318.
30. Smeester L, Rager JE, Bailey KA, Guan X, Smith N, Garcia-Vargas G, Del Razo LM, Drobná Z, Kelkar H, Styblo M *et al*: **Epigenetic changes in individuals with arsenicosis.** *Chemical research in toxicology* 2011, **24**(2):165-167.
31. Xie Y, Liu J, Benbrahim-Tallaa L, Ward JM, Logsdon D, Diwan BA, Waalkes MP: **Aberrant DNA methylation and gene expression in livers of newborn mice transplacentally exposed to a hepatocarcinogenic dose of inorganic arsenic.** *Toxicology* 2007, **236**(1-2):7-15.
32. Koestler DC, Avissar-Whiting M, Houseman EA, Karagas MR, Marsit CJ: **Differential DNA methylation in umbilical cord blood of infants exposed to low levels of arsenic in utero.** *Environmental health perspectives* 2013, **121**(8):971-977.
33. Kile ML, Houseman EA, Baccarelli AA, Quamruzzaman Q, Rahman M, Mostofa G, Cardenas A, Wright RO, Christiani DC: **Effect of prenatal arsenic exposure on DNA methylation and leukocyte subpopulations in cord blood.** *Epigenetics* 2014, **9**(5):774-782.
34. Maksimovic J, Gordon L, Oshlack A: **SWAN: Subset-quantile within array normalization for illumina infinium HumanMethylation450 BeadChips.** *Genome biology* 2012, **13**(6):R44.
35. Reinius LE, Acevedo N, Joerink M, Pershagen G, Dahlen SE, Greco D, Soderhall C, Scheynius A, Kere J: **Differential DNA methylation in purified human blood cells: implications for cell lineage and studies on disease susceptibility.** *PloS one* 2012, **7**(7):e41361.

36. Koestler DC, Christensen B, Karagas MR, Marsit CJ, Langevin SM, Kelsey KT, Wiencke JK, Houseman EA: **Blood-based profiles of DNA methylation predict the underlying distribution of cell types: a validation analysis.** *Epigenetics* 2013, **8**(8):816-826.
37. Jaffe AE, Irizarry RA: **Accounting for cellular heterogeneity is critical in epigenome-wide association studies.** *Genome biology* 2014, **15**(2):R31.
38. Houseman EA, Accomando WP, Koestler DC, Christensen BC, Marsit CJ, Nelson HH, Wiencke JK, Kelsey KT: **DNA methylation arrays as surrogate measures of cell mixture distribution.** *BMC bioinformatics* 2012, **13**:86.
39. Huang DW, Sherman BT, Tan Q, Collins JR, Alvord WG, Roayaei J, Stephens R, Baseler MW, Lane HC, Lempicki RA: **The DAVID Gene Functional Classification Tool: a novel biological module-centric algorithm to functionally analyze large gene lists.** *Genome biology* 2007, **8**(9):R183.
40. Kanehisa M, Goto S: **KEGG: kyoto encyclopedia of genes and genomes.** *Nucleic acids research* 2000, **28**(1):27-30.
41. Warde-Farley D, Donaldson SL, Comes O, Zuberi K, Badrawi R, Chao P, Franz M, Grouios C, Kazi F, Lopes CT *et al*: **The GeneMANIA prediction server: biological network integration for gene prioritization and predicting gene function.** *Nucleic acids research* 2010, **38**(Web Server issue):W214-220.
42. Gentleman R, Carey VJ, Huber W, Irizarry RA, Dudoit S: **Bioinformatics and computational biology solutions using R and Bioconductor**, vol. 746718470: Springer; 2005.
43. Smyth GK, Yang YH, Speed T: **Statistical issues in cDNA microarray data analysis.** *Methods Mol Biol* 2003, **224**:111-136.
44. Benjamini Y, Hochberg Y: **Controlling the false discovery rate: a practical and powerful approach to multiple testing.** *Journal of the Royal Statistical Society Series B (Methodological)* 1995:289-300.
45. Huang da W, Sherman BT, Lempicki RA: **Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources.** *Nature protocols* 2009, **4**(1):44-57.
46. Huang da W, Sherman BT, Lempicki RA: **Bioinformatics enrichment tools: paths toward the comprehensive functional analysis of large gene lists.** *Nucleic acids research* 2009, **37**(1):1-13.
47. Chan KH, Huang YT, Meng Q, Wu C, Reiner A, Sobel EM, Tinker L, Lusi AJ, Yang X, Liu S: **Shared molecular pathways and gene networks for cardiovascular disease and type 2 diabetes mellitus in women across diverse ethnicities.** *Circ Cardiovasc Genet* 2014, **7**(6):911-919.
48. Gowd V, Gurukur A, Chilkunda ND: **Glycosaminoglycan remodeling during diabetes and the role of dietary factors in their modulation.** *World J Diabetes* 2016, **7**(4):67-73.
49. Wang SL, Chiou JM, Chen CJ, Tseng CH, Chou WL, Wang CC, Wu TN, Chang LW: **Prevalence of non-insulin-dependent diabetes mellitus and related vascular diseases in southwestern arseniasis-endemic and nonendemic areas in Taiwan.** *Environmental health perspectives* 2003, **111**(2):155-159.
50. Gribble MO, Howard BV, Umans JG, Shara NM, Francesconi KA, Goessler W, Crainiceanu CM, Silbergeld EK, Guallar E, Navas-Acien A: **Arsenic exposure, diabetes prevalence, and diabetes control in the Strong Heart Study.** *American journal of epidemiology* 2012, **176**(10):865-874.
51. Relton CL, Gaunt T, McArdle W, Ho K, Duggirala A, Shihab H, Woodward G, Lyttleton O, Evans DM, Reik W *et al*: **Data Resource Profile: Accessible Resource for Integrated Epigenomic Studies (ARIES).** *Int J Epidemiol* 2015, **44**(4):1181-1190.
52. Luo J, Shu W: **Arsenic-Induced Developmental Neurotoxicity.** *Handbook of Arsenic Toxicology* 2014:363.
53. Gong G, O'Bryant SE: **The arsenic exposure hypothesis for Alzheimer disease.** *Alzheimer disease and associated disorders* 2010, **24**(4):311-316.

54. Vahidnia A, Romijn F, van der Voet GB, de Wolff FA: **Arsenic-induced neurotoxicity in relation to toxicokinetics: effects on sciatic nerve proteins.** *Chemico-biological interactions* 2008, **176**(2-3):188-195.
55. Lemarie A, Morzadec C, Bourdonnay E, Fardel O, Vernhet L: **Human macrophages constitute targets for immunotoxic inorganic arsenic.** *Journal of immunology* 2006, **177**(5):3019-3027.
56. Hsu WL, Tsai MH, Lin MW, Chiu YC, Lu JH, Chang CH, Yu HS, Yoshioka T: **Differential effects of arsenic on calcium signaling in primary keratinocytes and malignant (HSC-1) cells.** *Cell calcium* 2012, **52**(2):161-169.
57. Perera F, Herbstman J: **Prenatal environmental exposures, epigenetics, and disease.** *Reproductive toxicology* 2011, **31**(3):363-373.
58. Skinner MK: **Role of epigenetics in developmental biology and transgenerational inheritance.** *Birth Defects Research Part C: Embryo Today: Reviews* 2011, **93**(1):51-55.
59. Skinner MK: **Environmental epigenetic transgenerational inheritance and somatic epigenetic mitotic stability.** *Epigenetics* 2011, **6**(7):838-842.
60. Lee TW, Kwon H, Zong H, Yamada E, Vatish M, Pessin JE, Bastie CC: **Fyn deficiency promotes a preferential increase in subcutaneous adipose tissue mass and decreased visceral adipose tissue inflammation.** *Diabetes* 2013, **62**(5):1537-1546.
61. Kajimoto Y, Miyagawa J, Ishihara K, Okuyama Y, Fujitani Y, Itoh M, Yoshida H, Kaisho T, Matsuoka T, Watada H *et al*: **Pancreatic islet cells express BST-1, a CD38-like surface molecule having ADP-ribosyl cyclase activity.** *Biochemical and biophysical research communications* 1996, **219**(3):941-946.
62. Paulick MG, Bertozzi CR: **The glycosylphosphatidylinositol anchor: a complex membrane-anchoring structure for proteins.** *Biochemistry* 2008, **47**(27):6991-7000.
63. Pedersen LC, Tsuchida K, Kitagawa H, Sugahara K, Darden TA, Negishi M: **Heparan/chondroitin sulfate biosynthesis. Structure and mechanism of human glucuronyltransferase I.** *The Journal of biological chemistry* 2000, **275**(44):34580-34585.
64. Kreuger J, Kjellen L: **Heparan sulfate biosynthesis: regulation and variability.** *The journal of histochemistry and cytochemistry : official journal of the Histochemistry Society* 2012, **60**(12):898-907.
65. Grande-Allen KJ, Osman N, Ballinger ML, Dadlani H, Marasco S, Little PJ: **Glycosaminoglycan synthesis and structure as targets for the prevention of calcific aortic valve disease.** *Cardiovascular research* 2007, **76**(1):19-28.
66. Ballinger ML, Nigro J, Frontanilla KV, Dart AM, Little PJ: **Regulation of glycosaminoglycan structure and atherogenesis.** *Cellular and molecular life sciences : CMLS* 2004, **61**(11):1296-1306.
67. Schmidli RS, Colman PG, Cui L, Yu WP, Kewming K, Jankulovski C, Harrison LC, Pallen CJ, DeAizpurua HJ: **Antibodies to the protein tyrosine phosphatases IAR and IA-2 are associated with progression to insulin-dependent diabetes (IDDM) in first-degree relatives at-risk for IDDM.** *Autoimmunity* 1998, **28**(1):15-23.
68. Below JE, Gamazon ER, Morrison JV, Konkashbaev A, Pluzhnikov A, McKeigue PM, Parra EJ, Elbein SC, Hallman DM, Nicolae DL *et al*: **Genome-wide association and meta-analysis in populations from Starr County, Texas, and Mexico City identify type 2 diabetes susceptibility loci and enrichment for expression quantitative trait loci in top signals.** *Diabetologia* 2011, **54**(8):2047-2055.
69. Weber M, Hellmann I, Stadler MB, Ramos L, Paabo S, Rebhan M, Schubeler D: **Distribution, silencing potential and evolutionary impact of promoter DNA methylation in the human genome.** *Nature genetics* 2007, **39**(4):457-466.

70. Pai AA, Bell JT, Marioni JC, Pritchard JK, Gilad Y: **A genome-wide study of DNA methylation patterns and gene expression levels in multiple human and chimpanzee tissues.** *PLoS genetics* 2011, **7**(2):e1001316.
71. Pawankar R, Canonica GW, Holgate ST, Lockey RF, Organization WH: **White book on allergy 2011-2012 executive summary.** *World Allergy Organization* 2011.
72. Tezza G, Mazzei F, Boner A: **Epigenetics of allergy.** *Early human development* 2013, **89**:S20-S21.
73. Greer JM, McCombe PA: **The role of epigenetic mechanisms and processes in autoimmune disorders.** *Biologics : targets & therapy* 2012, **6**:307-327.
74. Javierre BM, Hernando H, Ballestar E: **Environmental triggers and epigenetic deregulation in autoimmune disease.** *Discovery medicine* 2011, **12**(67):535-545.
75. Liang L, Willis-Owen SAG, Laprise C, Wong KCC, Davies GA, Hudson TJ: **An epigenome-wide association study of total serum immunoglobulin E concentration.** *Nature* 2015, **520**.
76. Everson TM, Lyons G, Zhang H, Soto-Ramirez N, Lockett GA, Patil VK, Merid SK, Soderhall C, Melen E, Holloway JW *et al*: **DNA methylation loci associated with atopy and high serum IgE: a genome-wide application of recursive Random Forest feature selection.** *Genome Med* 2015, **7**:89.
77. Wang D, Yan L, Hu Q, Sucheston LE, Higgins MJ, Ambrosone CB, Johnson CS, Smiraglia DJ, Liu S: **IMA: an R package for high-throughput analysis of Illumina's 450K Infinium methylation data.** *Bioinformatics (Oxford, England)* 2012, **28**(5):729-730.
78. Lin L-C, Wang S-L, Chang Y-C, Huang P-C, Cheng J-T, Su P-H, Liao P-C: **Associations between maternal phthalate exposure and cord sex hormones in human infants.** *Chemosphere* 2011, **83**(8):1192-1199.
79. Wang S-L, Su P-H, Jong S-B, Guo YL, Chou W-L, Päpke O: **In utero exposure to dioxins and polychlorinated biphenyls and its relations to thyroid function and growth hormone in newborns.** *Environmental health perspectives* 2005:1645-1650.
80. Zhang TT, Okkenhaug K, Nashed BF, Puri KD, Knight ZA, Shokat KM, Vanhaesebroeck B, Marshall AJ: **Genetic or pharmaceutical blockade of p110delta phosphoinositide 3-kinase enhances IgE production.** *J Allergy Clin Immunol* 2008, **122**(4):811-819 e812.
81. Fu SL, Pierre J, Smith-Norowitz TA, Hagler M, Bowne W, Pincus MR, Mueller CM, Zenilman ME, Bluth MH: **Immunoglobulin E antibodies from pancreatic cancer patients mediate antibody-dependent cell-mediated cytotoxicity against pancreatic cancer cells.** *Clin Exp Immunol* 2008, **153**(3):401-409.
82. Jensen-Jarolim E, Achatz G, Turner MC, Karagiannis S, Legrand F, Capron M, Penichet ML, Rodriguez JA, Siccardi AG, Vangelista L *et al*: **AllergoOncology: the role of IgE-mediated allergy in cancer.** *Allergy* 2008, **63**(10):1255-1266.
83. Matta GM, Battaglio S, Dibello C, Napoli P, Baldi C, Ciccone G, Coscia M, Boccadoro M, Massaia M: **Polyclonal immunoglobulin E levels are correlated with hemoglobin values and overall survival in patients with multiple myeloma.** *Clin Cancer Res* 2007, **13**(18 Pt 1):5348-5354.
84. Nabih ES, Kamel HF, Kamel TB: **Association Between CD14 Polymorphism (-1145G/A) and Childhood Bronchial Asthma.** *Biochem Genet* 2016, **54**(1):50-60.
85. March ME, Sleiman PM, Hakonarson H: **Genetic polymorphisms and associated susceptibility to asthma.** *Int J Gen Med* 2013, **6**:253-265.
86. DeWan AT, Triche EW, Xu X, Hsu LI, Zhao C, Belanger K, Hellenbrand K, Willis-Owen SA, Moffatt M, Cookson WO *et al*: **PDE11A associations with asthma: results of a genome-wide association scan.** *J Allergy Clin Immunol* 2010, **126**(4):871-873 e879.
87. Zhang J, Noguchi E, Migita O, Yokouchi Y, Nakayama J, Shibasaki M, Arinami T: **Association of a haplotype block spanning SDAD1 gene and CXC chemokine genes with allergic rhinitis.** *J Allergy Clin Immunol* 2005, **115**(3):548-554.

88. Albrandt K, Orida NK, Liu FT: **An IgE-binding protein with a distinctive repetitive sequence and homology with an IgG receptor.** *Proc Natl Acad Sci U S A* 1987, **84**(19):6859-6863.
89. Konno S, Takahashi D, Hizawa N, Hattori T, Takahashi A, Isada A, Maeda Y, Huang SK, Nishimura M: **Genetic impact of a butyrophilin-like 2 (BTNL2) gene variation on specific IgE responsiveness to Dermatophagoides farinae (Der f) in Japanese.** *Allergol Int* 2009, **58**(1):29-35.
90. Adalsteinsson BT, Gudnason H, Aspelund T, Harris TB, Launer LJ, Eiriksdottir G, Smith AV, Gudnason V: **Heterogeneity in white blood cells has potential to confound DNA methylation measurements.** *PLoS one* 2012, **7**(10):e46705.
91. Talens RP, Boomsma DI, Tobi EW, Kremer D, Jukema JW, Willemsen G, Putter H, Slagboom PE, Heijmans BT: **Variation, patterns, and temporal stability of DNA methylation: considerations for epigenetic epidemiology.** *FASEB journal : official publication of the Federation of American Societies for Experimental Biology* 2010, **24**(9):3135-3144.
92. Bibikova M, Le J, Barnes B, Saedinia-Melnyk S, Zhou L, Shen R, Gunderson KL: **Genome-wide DNA methylation profiling using Infinium(R) assay.** *Epigenomics* 2009, **1**(1):177-200.
93. Gentleman RC, Carey VJ, Bates DM, Bolstad B, Dettling M, Dudoit S, Ellis B, Gautier L, Ge Y, Gentry J *et al*: **Bioconductor: open software development for computational biology and bioinformatics.** *Genome Biol* 2004, **5**(10):R80.
94. Aryee MJ, Jaffe AE, Corrada-Bravo H, Ladd-Acosta C, Feinberg AP, Hansen KD, Irizarry RA: **Minfi: a flexible and comprehensive Bioconductor package for the analysis of Infinium DNA methylation microarrays.** *Bioinformatics (Oxford, England)* 2014, **30**(10):1363-1369.
95. Houseman EA, Molitor J, Marsit CJ: **Reference-free cell mixture adjustments in analysis of DNA methylation data.** *Bioinformatics (Oxford, England)* 2014, **30**(10):1431-1439.
96. Houseman EA, Kile ML, Christiani DC, Ince TA, Kelsey KT, Marsit CJ: **Reference-free deconvolution of DNA methylation data and mediation by cell composition effects.** *BMC Bioinformatics* 2016, **17**:259.
97. Leek JT, Storey JD: **Capturing heterogeneity in gene expression studies by surrogate variable analysis.** *PLoS genetics* 2007, **3**(9):1724-1735.
98. Rahmani E, Zaitlen N, Baran Y, Eng C, Hu D, Galanter J, Oh S, Burchard EG, Eskin E, Zou J *et al*: **Sparse PCA corrects for cell type heterogeneity in epigenome-wide association studies.** *Nat Methods* 2016, **13**(5):443-445.
99. Carroll RJ, Ruppert D, Stefanski LA, Crainiceanu CM: **Measurement error in nonlinear models: a modern perspective:** CRC press; 2012.
100. Gagnon-Bartsch JA, Speed TP: **Using control genes to correct for unwanted variation in microarray data.** *Biostatistics* 2012, **13**(3):539-552.
101. Zou J, Lippert C, Heckerman D, Aryee M, Listgarten J: **Epigenome-wide association studies without the need for cell-type composition.** *Nature methods* 2014, **11**(3):309-311.
102. Lippert C, Listgarten J, Liu Y, Kadie CM, Davidson RI, Heckerman D: **FaST linear mixed models for genome-wide association studies.** *Nature methods* 2011, **8**(10):833-835.
103. Leek JT, Storey JD: **Capturing heterogeneity in gene expression studies by surrogate variable analysis.** *PLoS genetics* 2007, **3**(9):e161.
104. Zou JY: **Correcting for Sample Heterogeneity in Methylome-Wide Association Studies.** *Methods in molecular biology* 2015.
105. Smith AK, Conneely KN, Pace TW, Mister D, Felger JC, Kilaru V, Akel MJ, Vertino PM, Miller AH, Torres MA: **Epigenetic changes associated with inflammation in breast cancer patients treated with chemotherapy.** *Brain, behavior, and immunity* 2014, **38**:227-236.
106. Jaccard P: **Distribution de la Flore Alpine: dans le Bassin des dranses et dans quelques régions voisines:** Rouge; 1901.

107. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, Dolinski K, Dwight SS, Eppig JT *et al*: **Gene ontology: tool for the unification of biology. The Gene Ontology Consortium.** *Nature genetics* 2000, **25**(1):25-29.
108. Kanehisa M, Goto S, Sato Y, Kawashima M, Furumichi M, Tanabe M: **Data, information, knowledge and principle: back to metabolism in KEGG.** *Nucleic acids research* 2014, **42**(Database issue):D199-205.
109. Liu S, Guo X, Wu B, Yu H, Zhang X, Li M: **Arsenic induces diabetic effects through beta-cell dysfunction and increased gluconeogenesis in mice.** *Scientific reports* 2014, **4**:6894.
110. Guo L, Xiao Y, Wang Y: **Monomethylarsonous acid inhibited endogenous cholesterol biosynthesis in human skin fibroblasts.** *Toxicology and applied pharmacology* 2014, **277**(1):21-29.
111. Small EM, Frost RJ, Olson EN: **MicroRNAs add a new dimension to cardiovascular disease.** *Circulation* 2010, **121**(8):1022-1032.
112. Elinder F, Mannikko R, Pandey S, Larsson HP: **Mode shifts in the voltage gating of the mouse and human HCN2 and HCN4 channels.** *The Journal of physiology* 2006, **575**(Pt 2):417-431.
113. Mumford JL, Wu K, Xia Y, Kwok R, Yang Z, Foster J, Sanders WE: **Chronic arsenic exposure and cardiac repolarization abnormalities with QT interval prolongation in a population-based study.** *Environmental health perspectives* 2007, **115**(5):690-694.
114. Hernandez-Castro B, Doniz-Padilla LM, Salgado-Bustamante M, Rocha D, Ortiz-Perez MD, Jimenez-Capdeville ME, Portales-Perez DP, Quintanar-Stephano A, Gonzalez-Amaro R: **Effect of arsenic on regulatory T cells.** *Journal of clinical immunology* 2009, **29**(4):461-469.
115. Biswas D, Banerjee M, Sen G, Das JK, Banerjee A, Sau TJ, Pandit S, Giri AK, Biswas T: **Mechanism of erythrocyte death in human population exposed to arsenic through drinking water.** *Toxicology and applied pharmacology* 2008, **230**(1):57-66.
116. Andrew AS, Jewell DA, Mason RA, Whitfield ML, Moore JH, Karagas MR: **Drinking-water arsenic exposure modulates gene expression in human lymphocytes from a U.S. population.** *Environ Health Perspect* 2008, **116**(4):524-531.
117. Soto-Pena GA, Luna AL, Acosta-Saavedra L, Conde P, Lopez-Carrillo L, Cebrian ME, Bastida M, Calderon-Aranda ES, Vega L: **Assessment of lymphocyte subpopulations and cytokine secretion in children exposed to arsenic.** *FASEB journal : official publication of the Federation of American Societies for Experimental Biology* 2006, **20**(6):779-781.
118. Bethge N, Lothe RA, Honne H, Andresen K, Troen G, Eknaes M, Liestol K, Holte H, Delabie J, Smeland EB *et al*: **Colorectal cancer DNA methylation marker panel validated with high performance in Non-Hodgkin lymphoma.** *Epigenetics : official journal of the DNA Methylation Society* 2014, **9**(3):428-436.
119. Sjoblom T, Jones S, Wood LD, Parsons DW, Lin J, Barber TD, Mandelker D, Leary RJ, Ptak J, Silliman N *et al*: **The consensus coding sequences of human breast and colorectal cancers.** *Science* 2006, **314**(5797):268-274.
120. Blot-Chabaud M, Wanstok F, Bonvalet JP, Farman N: **Cell sodium-induced recruitment of Na(+)-K(+)-ATPase pumps in rabbit cortical collecting tubules is aldosterone-dependent.** *The Journal of biological chemistry* 1990, **265**(20):11676-11681.
121. Li D, Morimoto K, Takeshita T, Lu Y: **Arsenic induces DNA damage via reactive oxygen species in human cells.** *Environmental health and preventive medicine* 2001, **6**(1):27-32.
122. Martindale JL, Holbrook NJ: **Cellular response to oxidative stress: signaling for suicide and survival.** *Journal of cellular physiology* 2002, **192**(1):1-15.
123. Son Y, Kim S, Chung HT, Pae HO: **Reactive oxygen species in the activation of MAP kinases.** *Methods in enzymology* 2013, **528**:27-48.

124. El-Maarri O, Becker T, Junen J, Manzoor SS, Diaz-Lacava A, Schwaab R, Wienker T, Oldenburg J: **Gender specific differences in levels of DNA methylation at selected loci from human total blood: a tendency toward higher methylation levels in males.** *Human genetics* 2007, **122**(5):505-514.
125. Boks MP, Derks EM, Weisenberger DJ, Strengman E, Janson E, Sommer IE, Kahn RS, Ophoff RA: **The relationship of DNA methylation with age, gender and genotype in twins and healthy controls.** *PloS one* 2009, **4**(8):e6767.
126. Teschendorff AE, West J, Beck S: **Age-associated epigenetic drift: implications, and a case of epigenetic thrift?** *Human molecular genetics* 2013, **22**(R1):R7-R15.
127. Barfield RT, Almli LM, Kilaru V, Smith AK, Mercer KB, Duncan R, Klengel T, Mehta D, Binder EB, Epstein MP *et al*: **Accounting for population stratification in DNA methylation studies.** *Genetic epidemiology* 2014, **38**(3):231-241.
128. Fraser HB, Lam LL, Neumann SM, Kober MS: **Population-specificity of human DNA methylation.** *Genome biology* 2012, **13**(2):R8.

Name	TAIWAN COHORT			NHBCS			CpG Information			
	Coef (CI)	P-value	FDR P-value	Coef (CI)	P-value	FDR P-value	Gene	MAPINFO CHR	Relation	CpG Islands Name
cg08198265*	0.57 (0.31, 0.83)	6.77E-05	0.049	0.01(-0.09,0.15)	8.75E-01	0.992	BST1	15708451	4 S_Shelf	chr4:15704640-15705000
cg19850333	-0.36 (-0.53, -0.2)	6.77E-05	0.049	0.04(-0.06,0.15)	4.28E-01	0.977	CCRL2;CCRL2	46448579	3	
cg00049047	-0.37 (-0.54, -0.2)	6.78E-05	0.049	0.06(-0.06,0.19)	3.34E-01	0.967	GDNF	37838425	5 Island	chr5:37836747-37840726
cg25588480	0.41 (0.22, 0.6)	6.78E-05	0.049	-0.05(-0.19,0.09)	4.82E-01	0.988	MINK1;MINK1;MINK1;	4763240	17	
cg18132007	-0.21 (-0.3, -0.11)	6.79E-05	0.049	0.11(0.01,0.2)	2.86E-02	0.719	TP53111;TP53111	44972684	11 Island	chr11:44971048-44972685
cg16364152	-0.43 (-0.63, -0.23)	6.80E-05	0.049	-0.01(-0.13,0.1)	8.27E-01	0.992	RP5-1022P6.2	5591818	20 Island	chr20:5591490-5591875
cg14794023	0.27 (0.15, 0.4)	6.81E-05	0.049	0.1(-0.03,0.23)	1.18E-01	0.967	POM121L12	53103576	7 Island	chr7:53103275-53103801
cg21678813	-0.32 (-0.47, -0.17)	6.81E-05	0.049	0.13(0.02,0.24)	2.64E-02	0.719		20229879	11	
cg23319696	0.45 (0.24, 0.65)	6.82E-05	0.049	0.09(-0.1,0.27)	3.46E-01	0.968		2111870	5 Island	chr5:2111836-2112484
cg06935052	0.29 (0.16, 0.42)	6.84E-05	0.049	-0.04(-0.15,0.07)	4.74E-01	0.988	SMARCB1;SMARCB1	24176449	22	
cg19254118	0.44 (0.24, 0.64)	6.85E-05	0.049	-0.01(-0.15,0.14)	9.26E-01	0.992	ADAT3;SCAMP4	1907972	19 N_Shore	chr19:1908118-1908509
cg19724043	0.27 (0.14, 0.39)	6.86E-05	0.049	0.02(-0.12,0.17)	7.39E-01	0.992		15380395	21 N_Shelf	chr21:15383530-15383813
cg27640763	0.46 (0.25, 0.68)	6.89E-05	0.049	0.17(-0.02,0.37)	7.92E-02	0.893	LUM	91503109	12	
cg00217795	0.3 (0.16, 0.44)	6.90E-05	0.049	-0.03(-0.18,0.11)	6.47E-01	0.992	DIO2;DIO2;DIO2	80677688	14	
cg09568216	-0.52 (-0.76, -0.28)	6.90E-05	0.049	-0.04(-0.13,0.06)	4.48E-01	0.988	NPR3	32779925	5	
cg12523924	-0.33 (-0.49, -0.18)	6.90E-05	0.049	-0.13(0.01,0.25)	3.39E-02	0.719		63255359	5 Island	chr5:63255044-63255407
cg12241963	0.27 (0.14, 0.39)	6.91E-05	0.049	0.02(-0.14,0.17)	8.19E-01	0.992		33807279	6	
cg16275882	-0.26 (-0.38, -0.14)	6.91E-05	0.049	0.13(0.02,0.24)	1.69E-02	0.719		1288184	7 S_Shore	chr7:1286022-1287658
cg16703934	0.3 (0.16, 0.44)	6.91E-05	0.049	-0.07(-0.2,0.06)	2.71E-01	0.967	TSPAN33	1.29E+08	7 S_Shore	chr7:128784695-128785096
cg08622675	0.43 (0.23, 0.63)	6.92E-05	0.049	0.02(-0.11,0.16)	7.36E-01	0.992	KDELRL2;KDELRL2	6524998	7 S_Shore	chr7:6523064-6523897
cg07990873	0.27 (0.15, 0.4)	6.93E-05	0.049	-0.03(-0.14,0.08)	6.11E-01	0.992	ZNF671	58235167	19 N_Shelf	chr19:58238585-58239028
cg02816732	0.38 (0.2, 0.55)	6.94E-05	0.049	-0.12(-0.27,0.03)	1.13E-01	0.967	TNS3	47384337	7	
cg02333281	0.46 (0.25, 0.67)	6.96E-05	0.049	-0.06(-0.27,0.15)	5.93E-01	0.992		6636970	10	
cg26581228	0.51 (0.27, 0.74)	6.98E-05	0.049	-0.03(-0.13,0.06)	5.19E-01	0.989	TRERF1	42326264	6	
cg05310249	-0.43 (-0.62, -0.23)	6.99E-05	0.049	-0.02(-0.12,0.09)	7.34E-01	0.992	NKX2-6	23560590	8 Island	chr8:23559838-23560591
cg07158747	-0.4 (-0.58, -0.21)	6.99E-05	0.049					91196488	1 S_Shelf	chr1:91190489-91192804
cg02487202	0.31 (0.17, 0.45)	7.02E-05	0.049	-0.04(-0.12,0.05)	4.02E-01	0.968	ANKRD11	89358232	16	
cg06896857	0.34 (0.18, 0.5)	7.03E-05	0.049	0.04(-0.07,0.15)	4.74E-01	0.988	PPP1R15A;PPP1R15A	49375797	19 Island	chr19:49375484-49375928
cg00079219	-0.29 (-0.43, -0.16)	7.04E-05	0.049	0.08(-0.05,0.22)	2.03E-01	0.967	HOTAIR	54360131	12 N_Shore	chr12:54360374-54360660
cg08391419	0.31 (0.17, 0.46)	7.04E-05	0.049	0.01(-0.09,0.1)	9.14E-01	0.992	STK32C	1.34E+08	10 S_Shelf	chr10:134071971-134072193
cg05819837	0.35 (0.19, 0.5)	7.07E-05	0.049	0.07(-0.01,0.14)	8.29E-02	0.917	CUX1;CUX1;CUX1	1.02E+08	7	
cg09187936	0.23 (0.12, 0.33)	7.07E-05	0.049	-0.02(-0.12,0.09)	7.54E-01	0.992	SETD1B	1.22E+08	12 S_Shelf	chr12:122265374-122265954
cg04406115	0.32 (0.17, 0.46)	7.09E-05	0.049	-0.02(-0.16,0.13)	8.28E-01	0.992	KDM4B	5065640	19 Island	chr19:5065639-5065919
cg11175091	0.45 (0.24, 0.65)	7.09E-05	0.049	-0.02(-0.18,0.15)	8.47E-01	0.992	MIR1243;ANK2;ANK2;	1.14E+08	4	
cg15266969	0.29 (0.16, 0.42)	7.10E-05	0.049	0.07(-0.05,0.2)	2.23E-01	0.967	SLC22A12;SLC22A12	64369352	11	
cg00704664	1.15 (0.62, 1.68)	7.13E-05	0.049	1.21(0.02,2.4)	4.55E-02	0.776	CDH4	60500578	20 N_Shore	chr20:60501966-60502173
cg00269725	0.73 (0.39, 1.07)	7.23E-05	0.049	-0.03(-0.22,0.16)	7.31E-01	0.992		1.57E+08	6	
cg20918219	-0.22 (-0.32, -0.12)	7.23E-05	0.049	0.06(-0.02,0.15)	1.55E-01	0.967	SCARA3;SCARA3	27493854	8 S_Shelf	chr8:27490959-27491775
cg08557393	0.33 (0.18, 0.48)	7.24E-05	0.049	-0.06(-0.19,0.07)	3.31E-01	0.967	DOK4	57521521	16	
cg25338036	0.47 (0.26, 0.69)	7.24E-05	0.049	-0.02(-0.15,0.1)	6.93E-01	0.992	CSMD1	3047536	8	
cg03131730	0.27 (0.14, 0.39)	7.29E-05	0.049	0.06(-0.03,0.15)	2.01E-01	0.967	CCDC42;CCDC42	8638810	17	
cg05129295	0.39 (0.21, 0.56)	7.29E-05	0.049	0.02(-0.1,0.13)	7.87E-01	0.992		1316294	8	
cg01592801	-0.34 (-0.5, -0.18)	7.30E-05	0.049	-0.14(-0.37,0.1)	2.55E-01	0.967	KCNS2	99438942	8 Island	chr8:99438692-99440425
cg10130718	0.54 (0.29, 0.79)	7.35E-05	0.049	0.08(-0.34,0.5)	7.00E-01	0.992	DMRTB1	53925368	1 Island	chr1:53925191-53926228
cg03897139	0.31 (0.17, 0.45)	7.39E-05	0.05	0.03(-0.08,0.14)	5.50E-01	0.999	DPYSL4	1.34E+08	10 N_Shore	chr10:134019500-134019776
cg05965745	-0.41 (-0.59, -0.22)	7.39E-05	0.05	-0.05(-0.3,0.21)	7.18E-01	0.992	PRDM16;PRDM16	3077798	1 N_Shelf	chr1:3080934-3081292
cg13205528	-0.49 (-0.72, -0.26)	7.44E-05	0.05	-0.03(-0.25,0.19)	7.70E-01	0.992		2705849	1 N_Shore	chr1:2706025-2706961
cg10221365	0.28 (0.15, 0.4)	7.45E-05	0.05	0.06(-0.11,0.22)	5.02E-01	0.989	JMJD5;JMJD5	27214422	16 N_Shore	chr16:27214772-27215678
cg02493798	1.11 (0.6, 1.63)	7.50E-05	0.05	0.19(-0.37,0.76)	5.01E-01	0.989	ALOX12	6899577	17 Island	chr17:6898820-6900427
cg12401679	0.27 (0.14, 0.39)	7.51E-05	0.05	0.01(-0.15,0.17)	9.08E-01	0.992	LOC619207	1.35E+08	10 N_Shore	chr10:135270783-135271061
cg17943647	-0.55 (-0.81, -0.3)	7.51E-05	0.05	0.05(-0.09,0.18)	5.04E-01	0.989	TRIM2;TRIM2	1.54E+08	4	

Footnote: This results were obtained by fitting Robust regression with the M-values of DNA-methylation as response and log10 of maternal urinary arsenic level adjusting for Urinary creatinine, child's gender, batch effect, mother's age, mother's pre-pregnancy body mass index (BMI), mother's education level, and cell proportions of 6 cell types. M-values are logit of DNA methylation, defined as $\log_2 \left[\frac{\beta}{1-\beta} \right]$. CpG sites in bold showed consistent association with Total urinary arsenic level of mother in NHBCS, and CpGs marked with "*" are the CpGs such that their corresponding genes are in the identified KEGG pathways.

Table A1.2

Subject	Lymphocytes				Myeloid Cells	
	CD8T	CD4T	NK	Bcell	Mono	Gran
1	0.086339	0.161789	0.03705	0.125422	0.108936	0.535165
2	0.071657	0.128531	0.025776	0.120749	0.130483	0.579767
3	0.055667	0.117509	0.112656	0.106542	0.120631	0.523271
4	0.094488	0.105196	0.035431	0.098882	0.114146	0.585546
5	0.127137	0.144023	0.0091	0.077574	0.080921	0.568545
6	0.114757	0.171282	0	0.243679	0.127776	0.405114
7	0.221126	0.353146	0	0.191755	0.044692	0.219486
8	0.106903	0.159796	0.021962	0.125591	0.098902	0.526875
9	0.141326	0.176845	0.044873	0.108733	0.126921	0.440781
10	0.110878	0.195382	0	0.23698	0.128623	0.343273
11	0.113305	0.185754	0	0.207857	0.103421	0.4232
12	0.103108	0.306985	1.07E-19	0.046686	0	0.558618
13	0.056519	0.128894	0.064005	0.071728	0.109834	0.58592
14	0.039974	0.109659	0.130263	0.117196	0.080217	0.583316
15	0.131113	0.11647	0.020356	0.058528	0.07092	0.625321
16	0.093583	0.165071	0	0.093822	0.087372	0.581508
17	0.143037	0.104781	0	0.243772	0.1317	0.4711
18	0.091523	0.278061	0.116745	0.146964	0.100536	0.276378
19	0.10395	0.219228	0.014414	0.130644	0.104624	0.467834
20	0.042257	0.123753	0.082375	0.133762	0.142382	0.499199
21	0.124195	0.157698	0	0.131555	0.086831	0.531584
22	0.090689	0.155113	0.12881	0.177518	0.087908	0.419279
23	0.034937	0.096715	0.043709	0.124229	0.1107	0.63069
24	0.070149	0.180137	0.077627	0.126873	0.139379	0.410679
25	0.034482	0.135348	0.060564	0.141314	0.123158	0.560727
26	0.064729	0.038857	0.0107	0.061138	0.110737	0.73377
27	0.030939	0.145214	0.042139	0.103386	0.075217	0.63486
28	0.069561	0.176162	0.012068	0.261309	0.129053	0.39821
29	0.077367	0.200833	0	0.153033	0.084993	0.516621
30	0.046674	0.074382	0.060829	0.086505	0.098737	0.669492
31	0.061626	0.061423	0.024778	0.126057	0.113859	0.648363
32	0.124063	0.137529	0.070476	0.151346	0.218414	0.361695
33	0.115761	0.206974	0	0.178757	0.111402	0.453671
34	0.138564	0.128998	0.031599	0.140466	0.092372	0.519693
35	0.100378	0.128192	0.015141	0.175882	0.154002	0.461204
36	0.083233	0.119729	0.021163	0.19733	0.104917	0.531799
37	0.092963	0.07327	0.058152	0.058765	0.081584	0.686868
38	0.166116	0.326477	-1.39E-17	0.143829	0.073455	0.324551
39	0.11892	0.153317	0	0.254617	0.08274	0.421659
40	0.02767	0.085303	0.082031	0.08386	0.155102	0.5797
41	0.070268	0.05107	0.026707	0.08257	0.106381	0.70062
42	0.072235	0.153685	0.063122	0.143498	0.130354	0.495626
43	0.020794	0.118555	0.159978	0.148758	0.17757	0.401548
44	0.07095	0.161602	0	0.086025	0.115794	0.590266
45	0.016839	0.112772	0.045162	0.222292	0.222783	0.414852

Subject	Lymphocytes				Myeloid Cells	
	CD8T	CD4T	NK	Bcell	Mono	Gran
46	0.131383	0.22382	-6.94E-18	0.154252	0.121875	0.404513
47	0.02962	0.021543	0.020384	0.069937	0.104432	0.771136
48	0.155763	0.175066	0.117977	0.173808	0.076379	0.394373
49	0.090426	0.190173	0.057454	0.152837	0.073335	0.508091
50	0.048113	0.115364	0.029918	0.118548	0.17001	0.538238
51	0.120098	0.269878	0.018139	0.141287	0.126483	0.349225
52	0.057611	0.18277	0.101527	0.152387	0.11435	0.423623
53	0.120721	0.165444	0.042014	0.135022	0.11225	0.471869
54	0.103501	0.131154	0.053606	0.115437	0.14038	0.519567
55	0.075224	0.153331	3.47E-18	0.095834	0.111252	0.587624
56	0.068445	0.11561	0.123447	0.144862	0.208628	0.377559
57	0.086244	0.326938	0	0.168179	0.051681	0.393962
58	0.042673	0.091076	0.033382	0.097413	0.115793	0.641253
59	0.045027	0.060593	0.032034	0.087649	0.096683	0.711222
60	0.034927	0.03991	0.022319	0.105901	0.128001	0.696177
61	0.066205	0.124148	0	0.129001	0.169022	0.53968
62	0.091891	0.139212	0.011755	0.147737	0.129403	0.536163
63	0.087555	0.099459	0.099699	0.18809	0.108163	0.507678
64	0	0.121129	0.212002	0.121968	0.083253	0.507372

Footnote: Cell proportions for 6 cell types, calculated using *estimatecellcounts* in R-package *minfi* using the DNA methylation dataset for 64 subjects in this study

CpG	Estimate	StdErr	RAW_P	ahoc_p	Estimate	StdErr	RAW_P	ahoc_p	UCSC_Re UCSC_Re CHR	UCSC_CpG_Islands_Name	Relation_to_UCSC_CpG_Island
cgl13099261	-2.2995	0.4683	9.72E-05	0.031203							
cgl13229782	-1.2782	0.2449	4.88E-05	0.015653					17 chr17:21178819-21179690	S_Shore	
cgl13332474	1.094	0.1895	1.46E-05	0.00469							
cgl13487156	-1.9427	0.4111	0.000147	0.047272					RFTN1 Body		
cgl13578229	-0.7263	0.1316	2.53E-05	0.008116					ZNF510 TSS200		Island
cgl13612524	-1.1102	0.2239	8.74E-05	0.02806					ERLIN2;ER TSS1500;T		Island
cgl13656556	-2.0559	0.4136	8.49E-05	0.027268					GFPT2 Body		
cgl13689204	-1.6799	0.3358	7.91E-05	0.025376							
cgl13710556	-1.503	0.244	6.40E-06	0.002054						7 chr7:25891956-25892615	S_Shore
cgl13884344	-2.0018	0.3386	1.08E-05	0.003475					RNF38;RN 3'UTR;3'U		
cgl13918937	-1.1164	0.2328	0.000126	0.040306					C19orf10 1stExon		Island
cgl14151682	-0.9587	0.1887	6.65E-05	0.021333					PINK1 TSS200		N_Shore
cgl14354820	-2.0339	0.4079	8.19E-05	0.0263					ANK1 Body		
cgl14398659	-1.1672	0.243	0.000124	0.039658					INTS4 TSS200		Island
cgl14611830	-1.1771	0.2458	0.000128	0.04095					TTL4 TSS200		Island
cgl14631690	-2.4709	0.4646	3.92E-05	0.012594					CLPTM1 Body		Island
cgl14647515	-1.3148	0.2759	0.000134	0.043173					C10orf12 1stExon		
cgl14753070	-2.2551	0.4392	5.89E-05	0.018918					IL27RA 3'UTR		
cgl14894216	-1.619	0.3406	0.000138	0.044354					CATSPER1 TSS1500		
cgl15147833	-1.8938	0.3627	4.86E-05	0.015612					OR1C1 1stExon		
cgl15313617	-1.0822	0.2269	0.000134	0.042874					MAK16 Body		S_Shore
cgl15400591	-0.8711	0.1832	0.000138	0.044175					PROCR Body		S_Shore
cgl15582138	-0.8252	0.1511	2.86E-05	0.009185							
cgl15633603	-1.3646	0.2833	0.00012	0.038478							
cgl15665081	-0.5232	0.09288	1.97E-05	0.006334					CASP4;CAS 5'UTR;Bod		
cgl15712821	-0.9285	0.1827	6.62E-05	0.021252					PTK2B;PTK 5'UTR;5'U		S_Shelf
cgl15897970	0.9015	0.1889	0.000133	0.042548					FLJ43390 TSS200		Island
cgl16033376	-0.7046	0.1485	0.00014	0.0451					PRR24 TSS1500		Island
cgl16078269	-1.6532	0.3306	7.94E-05	0.025497					ZBTB7C TSS1500		
cgl16248798	-1.4764	0.2923	7.11E-05	0.022821						1 chr1:247274585-247275757	N_Shelf
cgl16328462	-2.6089	0.55	0.000141	0.05338						8	
cgl16415931	-1.4967	0.295	6.74E-05	0.021633					ZBED5;ZBE 5'UTR;1stE		Island
cgl16519192	1.0978	0.2116	5.24E-05	0.016813							
cgl16575125	-0.5879	0.1193	9.33E-05	0.029938					PLEKHG4B Body		S_Shelf
cgl16848072	-2.927	0.6205	0.00015	0.048139						3	
cgl16914277	-2.3503	0.4367	3.41E-05	0.010955					HLA-DOA TSS1500		S_Shelf
cgl17253785	-0.8776	0.1833	0.000128	0.041125					C8orf44 TSS1500		
cgl17287974	-1.6002	0.3312	0.000116	0.037256						2	
cgl17448109	-1.3295	0.2826	0.000154	0.049488					RNF115 Body		
cgl17454086	-1.2814	0.2452	4.81E-05	0.01545					KIAA0907 Body		N_Shore
cgl17470103	-2.4262	0.4784	6.78E-05	0.021772						1 chr1:155904072-155904307	N_Shore
cgl17476389	-1.9708	0.365	3.29E-05	0.010545					TTC15 Body		
cgl17529235	-2.3537	0.4614	6.34E-05	0.02036					C17orf51 Body		N_Shelf
cgl17563032	-1.096	0.2318	0.000146	0.047005					SRRT;SRRT Body;Body		
cgl17606881	-1.2973	0.2735	0.000142	0.045445						3 chr3:172468372-172468845	S_Shore
cgl17851795	-2.5593	0.5179	9.05E-05	0.02905					PBX2;GPS1 TSS1500;3		
cgl17856830	-0.8911	0.1784	8.04E-05	0.025808						6	
cgl17879648	-2.783	0.5249	4.07E-05	0.013051						18 chr18:14474587-14475008	N_Shore
cgl17890283	-2.115	0.3903	3.15E-05	0.010113						8 chr8:80942117-80942593	S_Shelf
cgl17993073	-1.3552	0.2657	6.35E-05	0.020381					IGF2R 3'UTR		
cgl18002076	-2.3792	0.4362	2.91E-05	0.009349					KIAA0947 Body		
cgl18038339	-1.3389	0.2233	9.06E-06	0.002907						5	
cgl18477664	-2.2401	0.4535	9.10E-05	0.029219						7 chr7:155898140-155898386	S_Shelf
cgl18481642	-2.6992	0.4812	2.08E-05	0.006668					STAT4 Body		
cgl18821281	-0.8715	0.145	8.80E-06	0.002824					VGLL4;VGI Body;1stEx		S_Shore
cgl18983417	-2.1852	0.3709	1.13E-05	0.003631						15	
cgl19194448	-1.8337	0.3611	6.68E-05	0.021439					RBM45 TSS1500		N_Shore
cgl19236703	-2.155	0.3059	1.05E-06	0.000336					KCNQ2;KC Body;Body		Island
cgl19385365	-0.7872	0.1658	0.00014	0.045055					DGKI Body		
cgl19457770	-1.6858	0.3382	8.23E-05	0.026432					IL1F9 3'UTR		
cgl19752083	-1.4714	0.2951	8.19E-05	0.0263					BAT2 Body		N_Shore
cgl19891951	-1.3769	0.2624	4.58E-05	0.014698					DDN TSS1500		Island
cgl19913448	-1.0509	0.1806	1.32E-05	0.004239					GSR TSS200		Island
cgl19998289	-2.2705	0.4062	2.17E-05	0.006954					RNF34;RN 3'UTR;3'U		
cgl20238308	-0.9452	0.1876	7.30E-05	0.023444					PTPRQ TSS1500		
cgl20268279	-1.2463	0.3384	4.79E-05	0.01536						14	
cgl20278840	-1.7354	0.3449	7.41E-05	0.023798					TMTC1 Body		
cgl20405508	-0.7329	0.1436	6.29E-05	0.020194					PDEFD TSS1500		S_Shore
cgl20848785	-1.4501	0.2728	3.95E-05	0.012668					VIPR2 Body		Island
cgl20943039	0.9471	0.174	2.98E-05	0.009564						15	
cgl21343777	-0.9616	0.188	6.15E-05	0.019747					RGS12;RG; Body;Body		N_Shore
cgl21480464	-1.2647	0.2687	0.000153	0.049253					PEMT;PEV Body;Body		
cgl21700440	-0.4798	0.09635	8.31E-05	0.026677					HS3ST3B1; Body;TSS1		Island
cgl21825397	-2.1304	0.4451	0.000129	0.041254						20	
cgl22047910	-0.8689	0.1618	3.51E-05	0.011279					SLC39A10; 5'UTR;5'U		Island
cgl22067069	-0.7032	0.1475	0.000134	0.042915					HSD17B4 TSS1500		N_Shore
cgl22157525	-1.5517	0.306	6.78E-05	0.021749					NDST1 Body		
cgl22287064	-1.3177	0.2662	8.89E-05	0.028553					MYO15B TSS200		Island
cgl22331108	-1.6337	0.2249	6.80E-07	0.000218						2 chr2:232478359-232479925	Island
cgl22697572	0.6922	0.1431	0.000114	0.03662					LOC65434 TSS200		
cgl22949256	-1.4867	0.2957	7.47E-05	0.023981					SGIP1 Body		
cgl22966316	-1.6837	0.3442	0.000101	0.032548					TTC1 5'UTR		
cgl23374847	-3.0824	0.6247	9.21E-05	0.02957						2	
cgl23671600	1.4123	0.2618	3.32E-05	0.010668					FBXO6 5'UTR		Island
cgl23761017	1.3371	0.2831	0.000148	0.047462						4 chr4:185937242-185937750	N_Shore
cgl23978322	-1.9565	0.3777	5.32E-05	0.017072					FHL5;FHL5 1stExon;5'		

Taiwanese

IoW

CpG Information

CpG	Estimate	StdErr	RAW_P	ahoc_p	Estimate	StdErr	RAW_P	ahoc_p	UCSC_Re	UCSC_Re	CHR	UCSC_CpG	Islands_Name	Relation_to_UCSC_CpG_Island
cg24054668	-1.646	0.2802	1.17E-05	0.003759					IQCE;IQCE	Body;Body	7			
cg24128590	1.0482	0.2054	6.32E-05	0.020289					C14orf43;(3'UTR;3'U'	14	chr14:74185482-74185994	N_Shore	
cg24391991	-0.9326	0.1933	0.000118	0.037822					MAP3K8;N	1stExon;5'	10	chr10:30722378-30723707	Island	
cg24823751	0.6716	0.1389	0.000115	0.036956					MIR1915;(TSS1500;T	10	chr10:21783198-21786420	Island	
cg24843246	-0.906	0.1903	0.000136	0.043565					PNPLA6;P'	Body;Body	19	chr19:7621639-7622262	Island	
cg24868248	-0.8812	0.1619	2.99E-05	0.009604					SFRS1;SFR	TSS1500;T	17	chr17:56083750-56085373	Island	
cg24883586	1.3478	0.2856	0.00015	0.04799					RAB11FIP1	TSS200;TS	8	chr8:37756454-37757339	Island	
cg25008393	-0.9218	0.1888	0.000103	0.033153					PFDN2;NIT	Body;TSS1	1	chr1:161087722-161088112	N_Shore	
cg25023095	1.6941	0.3569	0.00014	0.045061					RRP12;RR'	Body;Body	10			
cg25130710	-1.0804	0.2175	8.56E-05	0.027473					FBRS1	Body	12	chr12:133159225-133160576	Island	
cg25257018	-2.6533	0.5329	8.33E-05	0.02675					KDM4DL	TSS1500	11			
cg25306087	-1.1236	0.2013	2.21E-05	0.007086					OTOP2;US	TSS200;TS	17	chr17:72918995-72921019	Island	
cg25312876	-0.9213	0.1928	0.000131	0.041959					KDM4B	3'UTR	19	chr19:5151332-5151730	Island	
cg25321332	-3.1363	0.6635	0.000147	0.047125							18			
cg25404410	-1.6001	0.3256	9.65E-05	0.030963					KIF3A	TSS1500	5	chr5:132072695-132073429	S_Shore	
cg25453957	-0.9757	0.2014	0.000113	0.036121							17			
cg25584787	0.8155	0.1365	9.50E-06	0.003048					C5orf36	Body	5			
cg25679743	-2.8829	0.5776	8.11E-05	0.026047							19	chr19:42412375-42412584	Island	
cg25871713	-2.3077	0.4592	7.51E-05	0.024117					ITGA9	Body	3			
cg26122004	-1.2898	0.249	5.32E-05	0.017091					C5orf45;C'	Body;Body	5	chr5:179276846-179277068	Island	
cg26606257	-1.3614	0.2829	0.000121	0.03886					MYT1L	Body	2			
cg27040463	-0.6793	0.1228	2.46E-05	0.007896					ABHD13;LI	5'UTR;TSS	13	chr13:108870502-108871328	Island	
cg27055313	-0.5635	0.1188	0.000141	0.045259							14	chr14:106025533-106026386	Island	
cg27324576	-1.1763	0.2473	0.000137	0.044019					VEPH1;VEI	Body;Body	3			
cg27417659	-1.5031	0.2939	6.16E-05	0.019789					TRRAP	Body	7			
cg27420264	-0.8188	0.1662	9.33E-05	0.029965					HSPA8;HSI	TSS200;TS	11	chr11:122932173-122933803	Island	
ch_10_1024	-1.101	0.2143	5.86E-05	0.018807										
ch_4_13366	-0.8257	0.1534	3.40E-05	0.010912										
ch_5_19060	-1.2736	0.2484	6.00E-05	0.019255										

Table A2.2

CpG	ICC	P value
cg11029358	0.979056	4.45E-07
cg12349837	0.961259	1.60E-06
cg26015416	0.944679	1.39E-05
cg17207545	0.937823	2.46E-05
cg04385523	0.931517	0.000117
cg02299937	0.908664	1.29E-05
cg18469813	0.901578	9.13E-05
cg22590761	0.894415	2.65E-05
cg01750170	0.889368	0.000299
cg25288034	0.866309	0.000148
cg08482979	0.840104	0.00124
cg10083824	0.836534	0.001212
cg24699146	0.824922	0.000366
cg07708818	0.814555	0.000392
cg25637226	0.812584	0.002014
cg10002668	0.774514	0.00418
cg05132568	0.771446	0.002125
cg01769501	0.766981	0.004836
cg20458779	0.761597	0.004949
cg15549502	0.74647	0.001206
cg24600706	0.744432	0.001687
cg13984351	0.73477	0.00267
cg18773129	0.734442	0.001436
cg16383005	0.731928	0.006281
cg25499537	0.725816	0.008505
cg04990210	0.724765	0.003261
cg06829645	0.715809	0.002084
cg02448934	0.71478	0.001951
cg07805777	0.687972	0.008096
cg25471923	0.686291	0.002853
cg20075700	0.677139	0.007397
cg16019434	0.670026	0.016181
cg02346442	0.666076	0.008775
cg20431441	0.665892	0.009213
cg24677222	0.654635	0.019408
cg06998210	0.651392	0.005746
cg26624744	0.640641	0.008537
cg11413778	0.637632	0.00811
cg04614823	0.625649	0.009757
cg15367212	0.619501	0.007928
cg11823603	0.606962	0.027105
cg10548968	0.606855	0.007192
cg10343024	0.606723	0.02734
cg16998490	0.598492	0.031624
cg12143439	0.595929	0.01459
cg04723493	0.595728	0.011669

CpG	ICC	P value
cg06382167	0.593588	0.008852
cg25488567	0.593115	0.011583
cg12765716	0.588721	0.010329
cg11735008	0.579966	0.013543
cg00931181	0.574356	0.009914
cg08392484	0.569442	0.036974
cg05590053	0.559366	0.012538
cg00573504	0.539561	0.020963
cg02751327	0.535858	0.015777
cg02030454	0.524532	0.01608
cg13346967	0.519054	0.017335
cg11380128	0.513535	0.043542

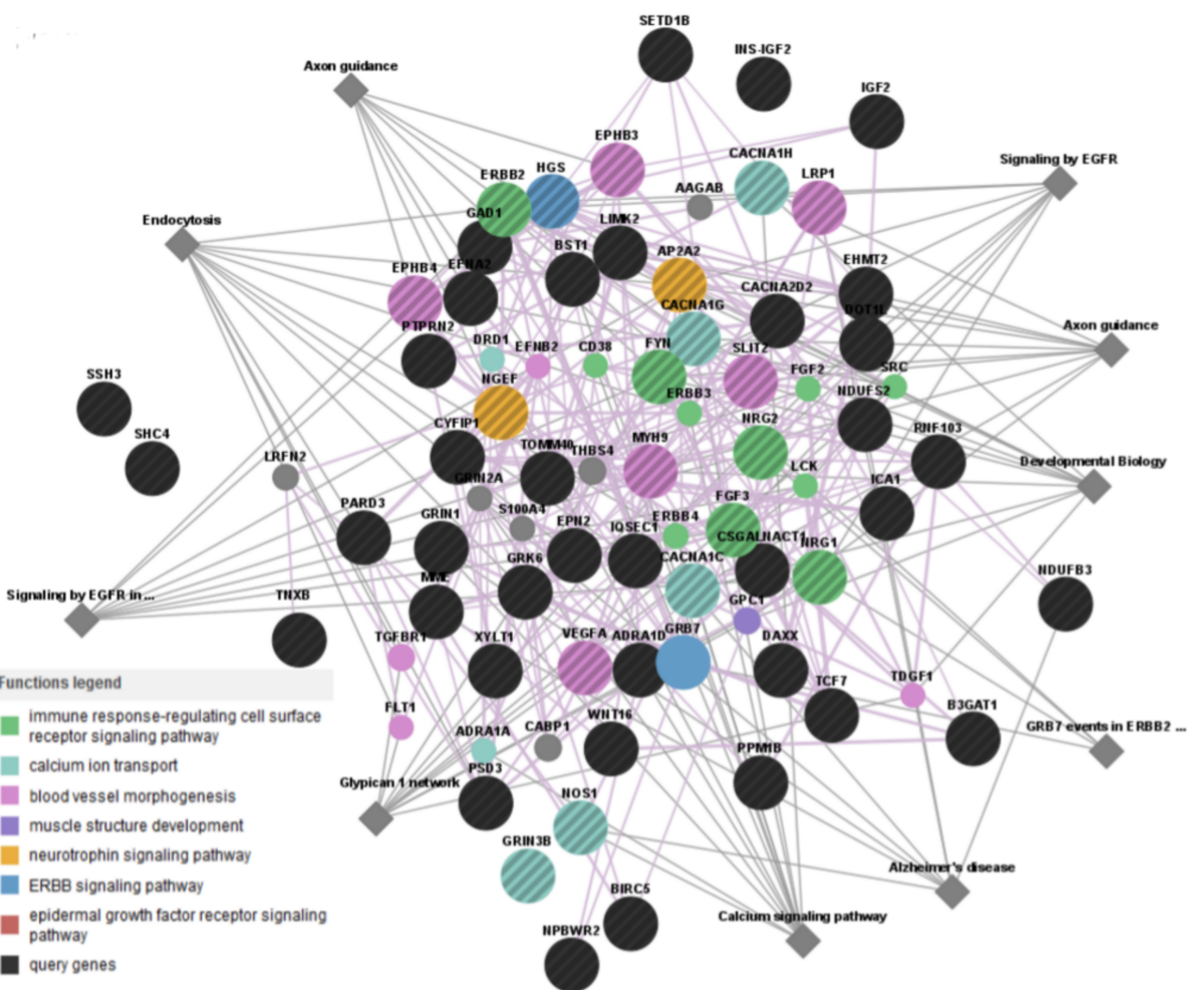


Figure A1.1

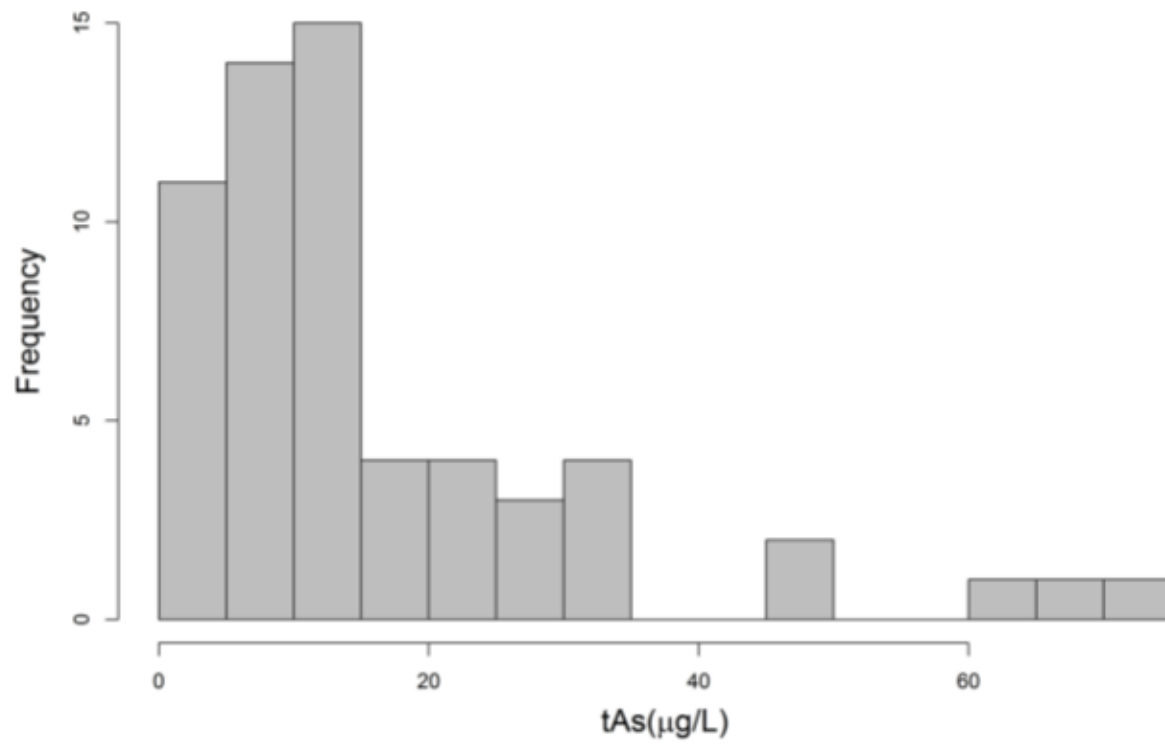


Figure A1.2

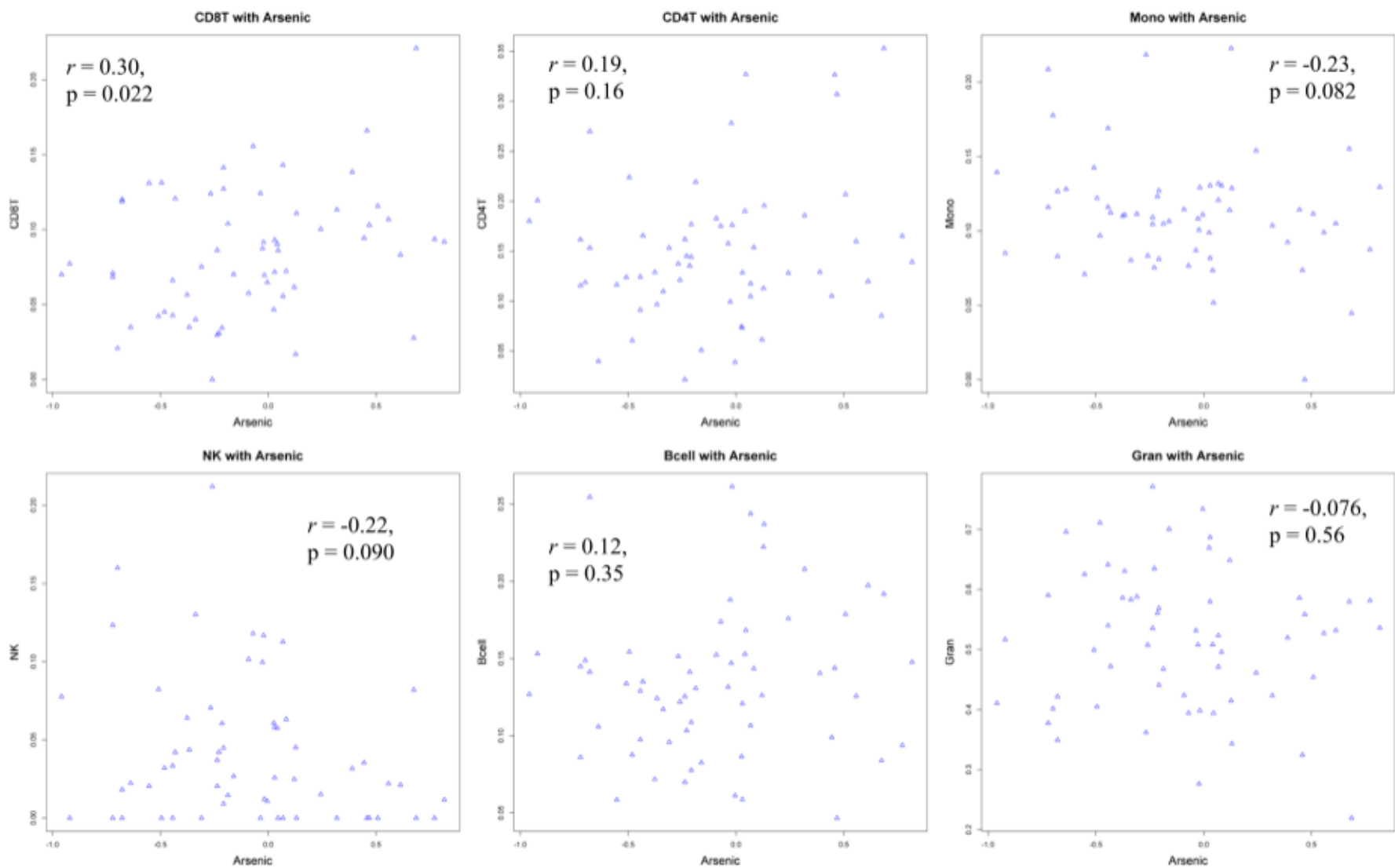


Figure A3.1. Pearson correlations between **inorganic arsenic levels** (in log10 scale) and cell type proportions.

Figure A3.1

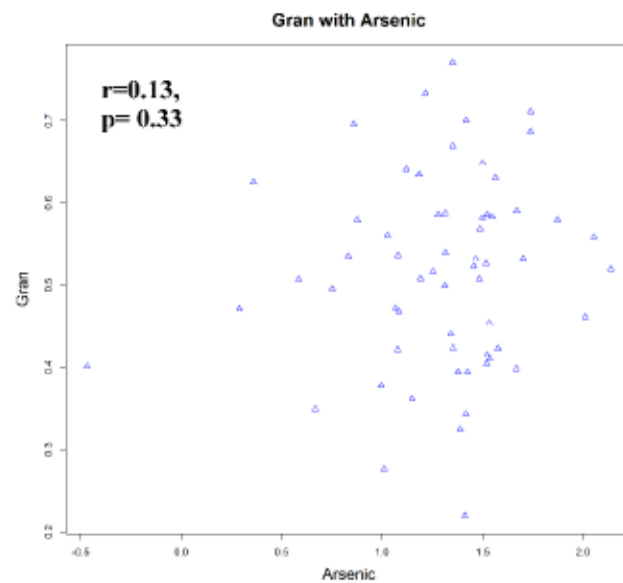
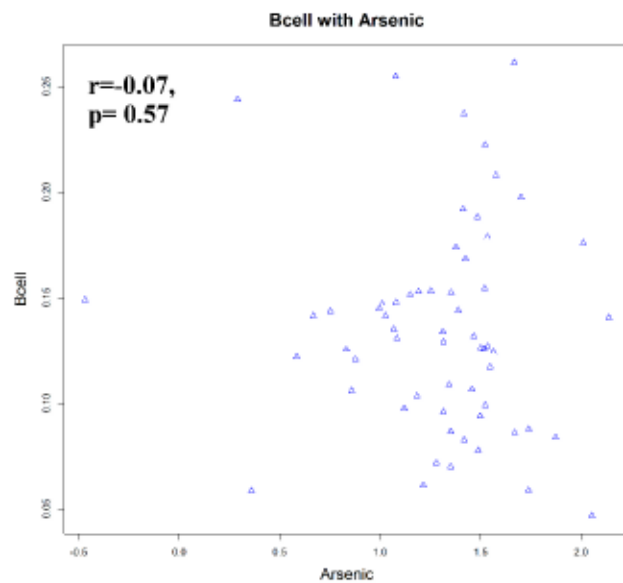
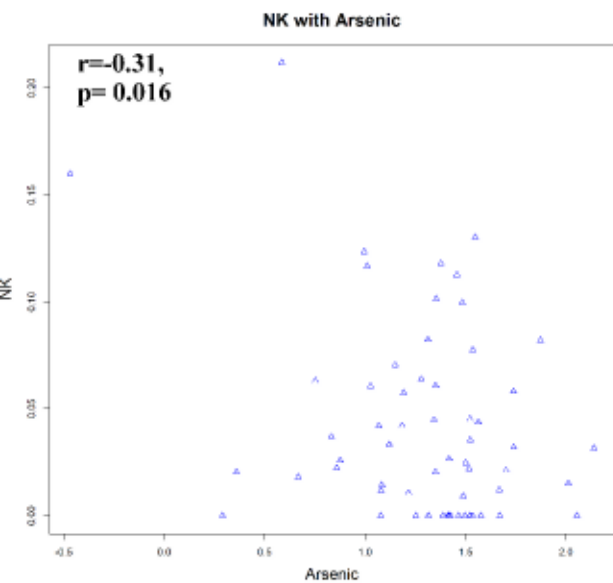
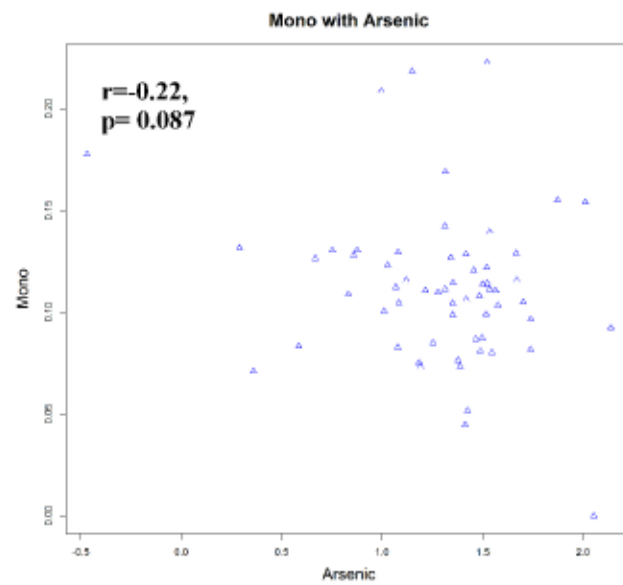
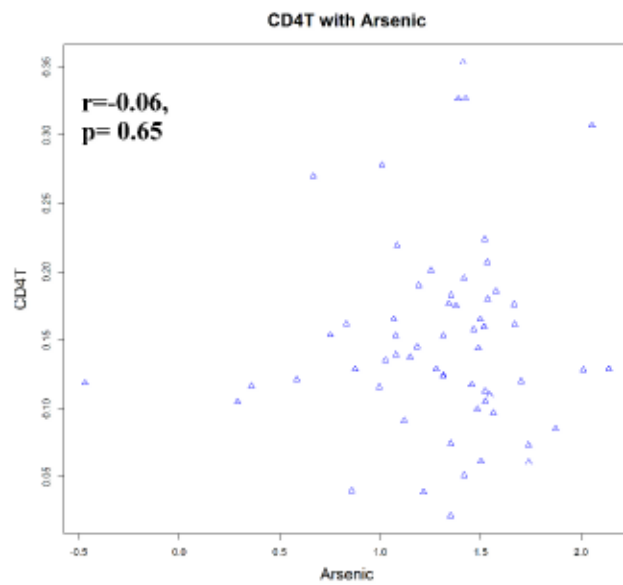
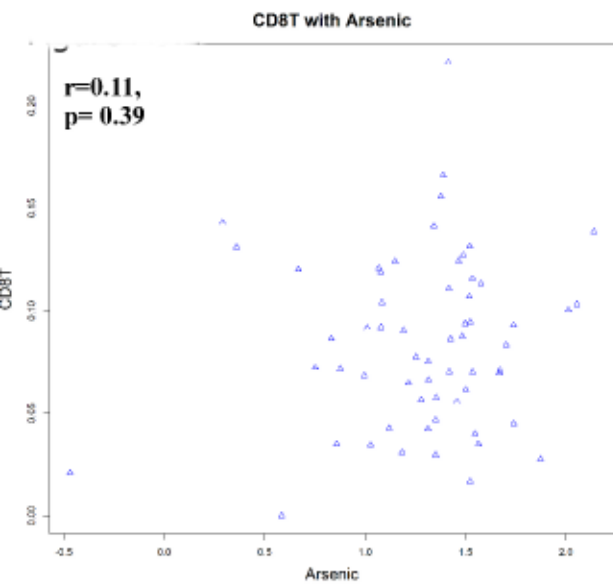


Figure A3.2

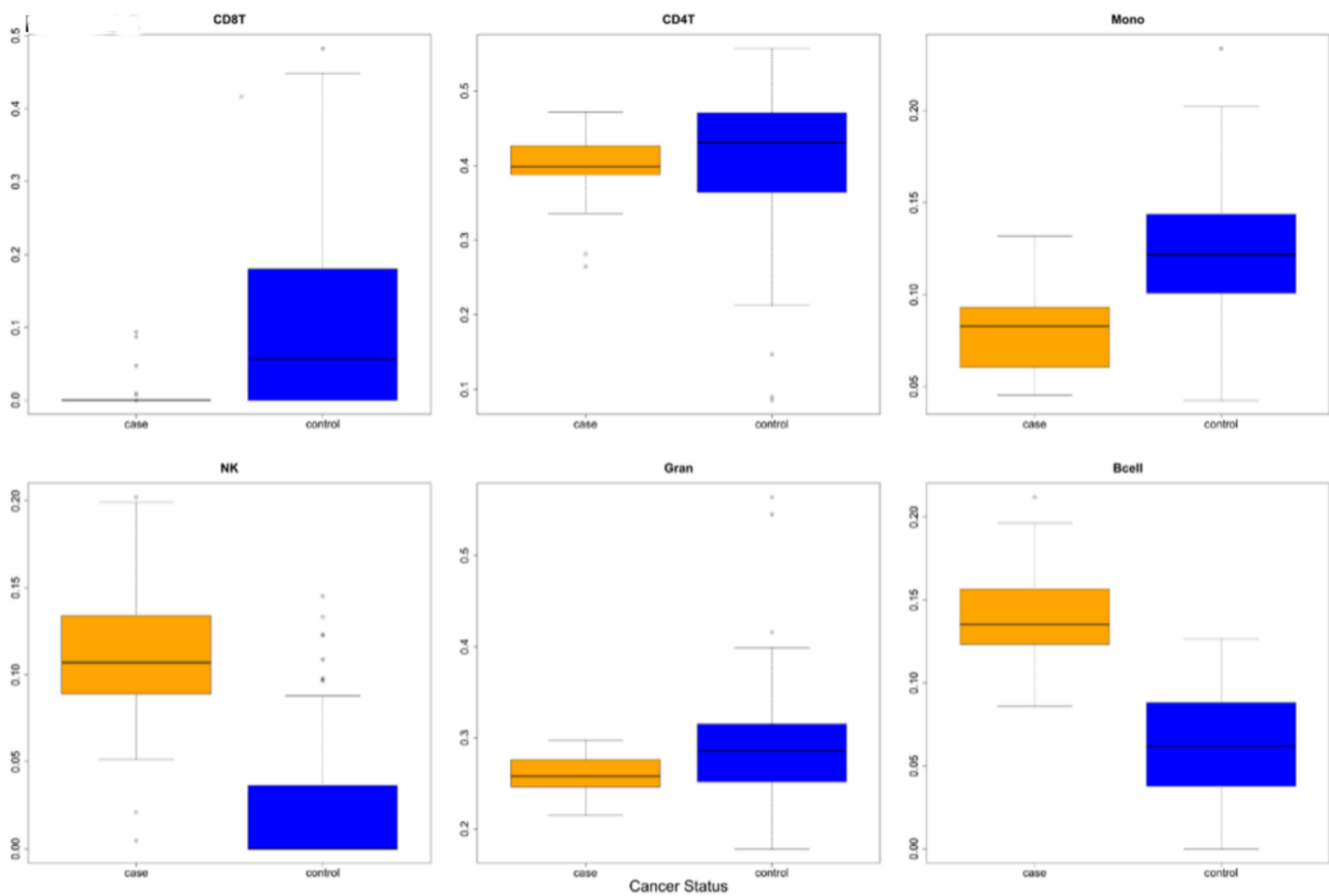


Figure A3.3

R_code_thesis

```
#####
### A) Installing and loading required packages
#####

if (!require("sva")) {
  install.packages("sva", dependencies = TRUE)
  library(sva)
}
if (!require("limma")) {
  install.packages("limma", dependencies = TRUE)
  library(limma)
}
if (!require("MASS")) {
  install.packages("MASS", dependencies = TRUE)
  library(MASS)
}
if (!require("gplots")) {
  install.packages("gplots", dependencies = TRUE)
  library(gplots)
}
if (!require("RColorBrewer")) {
  install.packages("RColorBrewer", dependencies = TRUE)
  library(RColorBrewer)
}
if (!require("LPS")) {
  install.packages("LPS", dependencies = TRUE)
  library(LPS)
}
if (!require("dendextend")) {
  install.packages("dendextend", dependencies = TRUE)
  library(dendextend)
}
if (!require("dendextendRcpp")) {
  install.packages("dendextendRcpp", dependencies = TRUE)
  library(dendextendRcpp)
}
if (!require("colorspace")) {
  install.packages("colorspace", dependencies = TRUE)
  library(colorspace)
}
if (!require("VennDiagram")) {
  install.packages("VennDiagram", dependencies = TRUE)
  library(VennDiagram)
}
#####
### B) Reading in data and transforming DNA-m beta values to M-values
#####
setwd("~/User Directory")
beta<-read.csv("DNA-m filename.csv", header=T) ##read in the DNA-m beta values
covar<-read.csv("covariate filename.csv", header=T)

beta1=beta[, -1]
rownames(beta1)=beta[, 1]

miss = which(is.na(beta), arr.ind=TRUE)[, 1];
beta1 = beta1[-miss, ];
M_val = as.matrix(log2(beta1/(1-beta1))); ##Logit transforming Beta value to M-values

##Read in the covariates from covar dataframe into specific variables

covar1=covar[, 2]; covar2=covar[, 3]; covar1=covar[, 4]; etc.
mod=model.matrix(~covar1+covar2+covar3)
```

R_code_thesis

```
#####  
### C) Performing Robust Regression  
#####  
fit<-lmFit(M_val, mod, method="robust")  
fit1<-eBayes(fit)  
tab <- topTable(fit1, coef = "covar1", number=length(beta1[, 1]), p.val=0.05, adjust =  
"fdr")  
write.table(tab, file='Robust_regression_results.csv', sep=",", row.names=TRUE)  
  
#####  
### D) Estimating Surrogate Variables (SVA) and using it in regression  
#####  
mod0 = model.matrix(~covar2+covar3)  
svobj1 = sva(edata, mod, mod0, n.sv=NULL, method="two-step")  
modSv = cbind(mod, svobj1$sv)  
fit2<-lmFit(M_val, modSv, method="robust")  
fit21<-eBayes(fit2)  
tab21 <- topTable(fit21, coef = "covar1", number=length(beta1[, 1]), p.val=0.05, adjust  
= "fdr")  
  
#####  
### D) Estimating Cell proportions using mi nfi  
#####  
source("https://bioconductor.org/biocLite.R")  
biocLite(c("minfi", "quadprog", "FlowSorted.Blood.450k",  
"IlluminaHumanMethylation450kmanifest",  
"IlluminaHumanMethylation450kanno.ilmn12.hg19"));  
  
lib = c("minfi", "quadprog", "FlowSorted.Blood.450k",  
"IlluminaHumanMethylation450kmanifest",  
"IlluminaHumanMethylation450kanno.ilmn12.hg19");  
lapply(lib, require, character.only = TRUE);  
  
grSet1=read.table("input_data.txt", header=T); grSet=grSet1[, -1];  
rownames(grSet)=grSet1[, 1]; grSet=data.matrix(grSet);  
referenceMset = get('FlowSorted.Blood.450k.compTable');  
  
cell = c("CD8T", "CD4T", "NK", "Bcell", "Mono", "Gran", "Eos");  
compData = minfi::pickCompProbes(referenceMset, cellTypes=cell);  
coefs = compData$coefEsts;  
coefs = coefs[intersect(rownames(grSet), rownames(coefs)), ];  
rm(referenceMset);  
  
counts = minfi::projectCellType(grSet[rownames(coefs), ], coefs);  
rownames(counts) = colnames(grSet);  
write.table(counts, file="Ewasher_data_minfi_cell.csv", sep=",");  
  
#####  
### D) Estimating Cell proportions using Houseman  
#####  
##Visit this website  
https://drive.google.com/forward?id=0B6IN3-9RV-LBeUt5bENMSXI sY2s&usp=sharing  
##to download all the relevant files  
  
beta=t(as.matrix(read.table("../input_data.txt", sep="\t", header=T)))  
source("Rcodes_Cell_mixture.R")  
  
#####  
### E) Using Refactor to adjust for Cell types  
#####  
##Download the Refactor source code from following link:
```

<http://www.cs.tau.ac.il/~heran/cozygene/software/refactor.html>

```

sd1=apply(beta1, 1, sd, na.rm = F)
include = which(sd1>=quantile(sd1, 0.05))
O = beta1[include, ]
M_val1 = M_val[include, ]

cpgnames <- rownames(O)
for (site in 1:nrow(O))
{
  model <- lm(O[site, ] ~ covar1+covar2+covar3)
  O_adj[site, ] = residuals(model)
}
O = O_adj
print('Running a standard PCA...')
pcs = prcomp(scale(t(O)));

coeff = pcs$rotation
score = pcs$x

print('Compute a low rank approximation of input data and rank sites...')
x = score[, 1:7]%%t(coeff[, 1:7]);
An = scale(t(O), center=T, scale=F)
Bn = scale(x, center=T, scale=F)
An = t(t(An)*(1/sqrt(apply(An^2, 2, sum))))
Bn = t(t(Bn)*(1/sqrt(apply(Bn^2, 2, sum))))

# Find the distance of each site from its low rank approximation.
distances = apply((An-Bn)^2, 2, sum)^0.5 ;
dsort = sort(distances, index.return=T);
ranked_list = dsort$ix

print('Compute ReFACTor components...')
sites = ranked_list[1:500];
pcs = prcomp(scale(t(O[sites, ])));
first_score <- score[, 1:7];
score = pcs$x

mod=model.matrix(~covar1+covar2+covar3+first_score)
fit1<-lmFit(M_val1, mod, method="ls")
fit1<-eBayes(fit1)
tab1 <- topTable(fit1, coef = "covar1", number=length(beta1[, 1]), p.val =0.05, adjust = "fdr")

write.table(tab1, file="Refactor_results.csv", sep=" ", row.names=T)

#####
### F) Using RefFreeCellMix to adjust for Cell types
#####

library(RefFreeEWAS)
cell <-RefFreeCellMix(beta1, mu0=NULL, K=7, iters=5, Yficial=NULL, verbose=TRUE)
mod1<-model.matrix(~covar1+covar2+covar3+cell$Omega)

fit1<-lmFit(M_val, mod1, method="robust")
fit1<-eBayes(fit1)
tab1 <- topTable(fit1, coef = "covar1", number=length(beta1[, 1]), p.val =0.05, adjust = "fdr")

write.table(tab1, file="RefFreeCellMix_results.csv", sep=" ", row.names=T)

#####

```

```

R_code_thesis
### F) Using RUV to adjust for Cell types
#####

library(data.table)
library(psych)
library(pracma)

##Specify beta values of top 600 CpG sites based upon the reference dataset

beta<-read.csv("Top_600_Mnfi_27k.csv", header=T) ##Specify Top 600
beta1=beta[, -1]
rownames(beta1)=beta[, 1]
beta1_t<-t(beta1)

pca_mval <-prcomp(beta1_t, center=T, scale.=T)
plot(pca_mval, type = "l")

###Based upon scree plot we can choose number of PC's##
rawLoadings <- pca_mval$rotation[, 1:7] %% diag(pca_mval$sdev, 7, 7)
rotatedLoadings <- varimax(rawLoadings, normalize = TRUE, eps = 1e-5)$loadings
invLoadings <- t(pracma::pinv(rotatedLoadings))
scores <- scale(beta1_t) %% invLoadings

###x1, x2, x3, x4, x5 and x6 are primary and secondary covariates,
###PC1, PC2, PC3 and PC4 are principal components
#####
#####
mod<-model.matrix(~x1+x2+x3+x4+x5+x6+PC1+PC2+PC3+PC4)
fit<-lmFit(edata, mod, method="robust")
fite<-eBayes(fit)
tab <- topTable(fite, coef = "x1", number=length(edata[, 1]), p.val=0.05, adjust =
"fdr")

#####
#####
### F) Bar plots
#####
#####

data<-read.csv("coef.csv", header=T)
data1=data[, -1]
rownames(data1)=data[, 1]

category<-factor(colnames(data1))

col<-c("orange1", "blue")[category]

barplot(t(data1), beside=T, col = cols, width = 0.82, space = NULL, names.arg = NULL,
legend.text = NULL, horiz = FALSE, density = NULL, las=2, border =
par("fg"), main = NULL, sub = NULL, xlab = NULL, ylab = NULL, xlim =
NULL, ylim = c(-0.5, 0.5), xpd = TRUE, log = "", axes = TRUE, axisnames = TRUE,
plot = TRUE, cex.names = 0.85, axis.lty = 1, offset = 0, add = FALSE,
args.legend = NULL)
abline(h=0)
title("Bar plot", "", "", "Methylation")
text(seq(1, 29, by=1), par("usr")[3] - 0.2, labels = rownames(data1), srt = 90, pos
= 2, offset = 0, vfont = NULL, col = NULL, font = NULL)

axis(1, at=seq(1, 29, by=1), tick = F, line = NA, pos = NA, outer = FALSE, font = NA,
lty = "solid", labels = FALSE, lwd = 1.5, col = NULL, col.ticks = NULL, hadj = NA,
padj = NA, lwd.ticks=1)
legend(1, 0.5, c("CAT", "HDM"), lty=c(1, 1), lwd=c(3, 3), col=c("orange1", "blue"))

```


R_code_thesis

```
box(bty="l")
```

```
#####  
#####  
### F) Manhattan plots  
#####  
#####
```

```
##Please refer to this website  
http://genome.sph.umich.edu/wiki/Code\_Sample:\_Generating\_Manhattan\_Plots\_in\_R
```

```
#####  
#####  
### G) Heatmap  
#####  
#####
```

```
png("Plotname.png",      # create PNG for the heat map  
     width = 5*300,      # 5 x 300 pixels  
     height = 7*300,  
     res = 300,          # 300 pixels per inch  
     pointsize = 5)      # smaller font size
```

```
Rowv<-mat_data %>% dist %>% hclust %>% as.dendrogram %>%  
  set("branches_k_color", k = 5) %>% set("branches_lwd", 1.5) %>%  
  ladderize
```

```
heatmap.2(mat_data,      ###mat_data contains the correlation values or beta values  
          cex.mai n=5.0,  
          key=T,  
          keysize=0.85,  
          symkey=F,  
          main = "Correlation", # heat map title  
          notecol="black",      # change font color of cell labels to black  
          density.info="density", # turns off density plot inside color legend  
          trace="none",        # turns off trace lines inside the heat map  
          margins =c(13, 11),  # widens margins around plot  
          col=my_palette,      # use on color palette defined earlier  
          breaks=col_breaks,  
          dendrogram="row",    # only draw a row dendrogram  
          scale="none",  
          Rowv=Rowv,  
          Colv="NA",          # turn off column clustering  
          cexCol=2.0,  
          cexRow=2.0  
          )
```

```
dev.off()
```

```
#####  
#####  
### G) Venn Diagrams  
#####  
#####  
ars<-read.csv("DNA-m data.csv", header=T)  
venn.plot<-venn.diagram(  
x = list(  
SVA = ars$SVA,  
Houseman = ars$Houseman,  
mifi = ars$mifi ,
```

```
RefFreeEWAS = ars$RefFreeEWAS,  
RefFreeCellMix = ars$RefFreeCellMix  
)  
filename = 'Venn_Ewasher1.png',  
output = TRUE,  
height = 35,  
width = 65,  
resolution = 300,  
units = 'in', cat.pos = c(0, 310, 215, 145, 50),  
rotation.degree=15,  
fill = gray.colors(5, start = 0.2, end = 0.9, gamma = 1.5, alpha = 0.7),  
lty = "solid",  
cex=5,  
cat.cex=7,  
Scaled=TRUE  
);
```

```

SAS code for Epi genome wide mixed modeling
ods graphics off; ods html close; ods listing close;

/*DISPLAY P-VALUES TO A HIGHER NUMBER OF DECIMAL PLACES*/
ods path sasuser.templat(update) sashelp.templmst(read);
proc template;
  edit Common.PValue;
  notes "Default p-value column";
  just = r;
  format = pvalue15.9;
end;
run;
/*
/ Unless you want to keep this edited template, delete it.
/ this only deletes the version of the template in your SASUSER
/ library, returning you to the regularly scheduled official one in
/ SASHELP.
/-----*/

/* The first 5 rows of data below are ID and covariates and next rows are
DNA-methylation data*/
/* I have two rows for ID by the names ID_POS and ID*/

PROC IMPORT OUT= WORK.IgEt
  DATAFILE= "C:\Users\akaushl1\Documents\IGE_MIXED\iow11.csv"
  DBMS=CSV REPLACE;

RUN;

Data Work.IgEt1;
  SET Work.IgEt (obs=5);
run;
Data Work.IgEt2;
  SET Work.IgEt (firstobs=6);
run;

Data work.IgEt;
  Set Work.IgEt1 Work.IgEt2;
run;
proc transpose data= Work.IgEt1 et
  out= Ige_new (where=(upcase(_name_) ne 'ID_POS'));
  var _all_;
  id id_pos;
run;

Data Ige_new1;
Set Ige_new(firstobs=2);
run;

DATA Ige_new_Long ;
  SET Ige_new1;
  Ige = serum_ige_10; time = 10; OUTPUT;
  Ige = serum_ige_18; time = 18; OUTPUT;

  DROP serum_ige_10 serum_ige_18;
RUN;

Data Ige_new_Long;
  Set Ige_new_Long;
  if Ige<0 then Ige=. ;
run;
data Ige_new_Long1;
  set Ige_new_Long;

```

SAS code for Epi genome wide mixed modeling

```

tcat=time;
newI gE = i nput(I gE, best32. );
cb_i ge=i nput(cb_i ge, best32. );
bw=i nput(BI RTHWT, best32. );
I I gE=l og10(newI gE);
drop I gE newI gE cb_i ge;
rename I I gE=I gE;

run;
proc print data= I gE_new_l ong1; run;

data I gE1_new;
set I gE_new1(drop= serum_i ge_10 serum_i ge_18 _NAME_ ID cb_i ge BI RTHWT);
run;

proc contents data = I gE1_new
out = vars(keep = varnum name)
nopri nt;
run;
proc print data= I gE1_new; run;

data _null_;
mVars +1;
len = 0;
length str $10;
call execute( cats( '%str( %%l et name)', mVars , '=' ) );
do WHILE( NOT eof );
set work.vars(keep= name) end= eof;
have+1;

str = trim( name );
len + ( 3 + length(str) );
if EOF OR len gT 20000 then LEAVE;
end;
call execute( trim(str) !!' ; ' );
putlog 'info: ' have= 'in ' mVars=;
if eof then
do;
call symputx( 'n_name_vars', mVars );
stop;
end;

run;

```

SAS code for Epi genome wide mixed modeling

```

%macro split;
%do i = 1 %to &n_name_vars;
data IgE1_new&i.;
set IgE1_new (keep=name&i);
run;
%end; ;
%mend split;
%split;
/*
data IgE1_new1;
  set IgE1_new (keep=&x1);
  run;
*/

/* Macro to split the large dataset into bunch of small dataset */

%macro split1;
%do i = 1 %to &n_name_vars;

proc contents data = IgE1_new&i.
out = vars&i. (keep = name type)
noprint;
run;
%end; ;
%mend split1;
%split1;

proc print data=vars1; run;

%macro split2;

%do k = 1 %to &n_name_vars;
data vars&k.;
set vars&k.;
if type=2 and name ne 'id';
newname=trim(left(name))||"_n";

proc sql noprint;
select trim(left(name)), trim(left(newname)),
       trim(left(newname))||'|'='||trim(left(name))
into :c_list separated by '|', :n_list separated by '|',
      :renam_list separated by '|';
from vars&k.;

data IgE_new_long2&k.;
set IgE_new_long1 (keep=lgE ID cbi ge TIME TCAT bw &&name&k.);
array ch(*) $ &c_list;
array nu(*) &n_list;
do i = 1 to dim(ch);
  nu(i)=input(ch(i), 5.);
end;
drop i &c_list;
rename &renam_list;

run;

options nosymbolgen;
%end;

```

SAS code for Epi genome wide mixed modeling

```

%mend split2;
%split2;

proc print data=IgE_new_long21; run;

%macro mylogita(indata, indvars, dep, myout =_out );
  %let k=1;
  %let ind = %scan(&indvars, &k);
  %do %while(&ind NE);
    title "The dependent variable is &dep";
    title2 "The independent variables are &ind";

ODS output SolutionF = est1&k Tests3=est2&i;

proc mixed data=&indata method=ML ;
  class ID;
  model &dep= &in cbi ge time bw/solution ddfm=bw;

  random int time/Subject=ID type=AR(1); *R=1,2 RCORR;
run;

ODS OUTPUT CLOSE;

%let k = %eval (&k + 1);
  %let ind = %scan(&indvars, &k);
  %end;

  data &myout;
  set
  %do i = 1 %to &k - 1;
    est1&i est2&i
  %end;
  ;
run;

%mend;

*run the program;

%macro split3;

%do j = 1 %to &n_name_var;
%mylogita(work.IgE_new_long2&j., &&name&j., IgE, myout = myparms&j.)
%end;
%mend split3;
%split3;

*ods trace off;

%macro split4;

%do k = 1 %to &n_name_vars;

Data myparm1 (Keep = Effect Estimate StdErr df);
  Set myparms&k. (where=(probt NE . and Effect not in
(' time', ' Intercept', ' tcat', ' cbi ge', ' bw')));
  run;
title;
Data myparm2 (Keep = Effect probf);
  Set myparms&k. (where=(probf NE . and Effect not in
(' time', ' Intercept', ' tcat', ' cbi ge', ' bw')));

```

SAS code for Epi genome wide mixed modeling

```

run;

* 1. Sort myparms1 by "Effect" & save sorted file as mayparms11 ;
PROC SORT DATA=myparm1 OUT=myparms11;
  BY Effect;
RUN;

* 2. Sort myparms2 by "Effect" & save sorted file as mayparms21 ;
PROC SORT DATA=myparm2 OUT=myparms21;
  BY Effect;
RUN;

* 3. Merge myparms11 and myparms21 by Effect in a data step ;
DATA myparms12&k. ;
  MERGE myparms11 myparms21;
  BY Effect;
RUN;
%end;
%mend split4;
%split4;

%macro combine;

data big;
  set
    %do i = 1 %to &n_name_vars;

      myparms12&i.
    %end;
  ;
run;

%mend;

%combine;

Data big1;
set big(rename=(probf=RAW_P));
run;
proc multtest pdata=big1 ADAPTI VEHOCHBERG out=big3p;
run;

PROC EXPORT DATA= WORK. big3p
  OUTFILE=
"C:\Users\akaushl1\Documents\I gE_MI XED\Resul ts_mi xed_i ow1_common_wi th_tai _j une1.csv"
  DBMS=CSV REPLACE;
  PUTNAMES=YES;
RUN;

```

IRB Approval letter

Hello,

The University of Memphis Institutional Review Board, FWA00006815, has reviewed and approved your submission in accordance with all applicable statuses and regulations as well as ethical principles.

PI NAME: Akhilesh Kaushal

CO-PI:

PROJECT TITLE: Comparison of different cell type correction methods for genome-wide epigenetics studies

FACULTY ADVISOR NAME (if applicable): Hongmei Zhang

IRB ID: #4075

APPROVAL DATE: 6/62016

EXPIRATION DATE:

LEVEL OF REVIEW: Exempt

Please Note: Modifications do not extend the expiration of the original approval

Approval of this project is given with the following obligations:

- 1. If this IRB approval has an expiration date, an approved renewal must be in effect to continue the project prior to that date. If approval is not obtained, the human consent form(s) and recruiting material(s) are no longer valid and any research activities involving human subjects must stop.**
- 2. When the project is finished or terminated, a completion form must be completed and sent to the board.**
- 3. No change may be made in the approved protocol without prior board approval, whether the approved protocol was reviewed at the Exempt, Exedited or Full Board level.**
- 4. Exempt approval are considered to have no expiration date and no further review is necessary unless the protocol needs modification.**

Approval of this project is given with the following special obligations:

Thank you,

James P. Whelan, Ph.D.

Institutional Review Board Chair

The University of Memphis.

Note: Review outcomes will be communicated to the email address on file. This email should be considered an official communication from the UM IRB.