



Federated analyses for large epigenetic studies

Josine Min

Senior Research Fellow in Genetic and Epigenetic Epidemiology MRC IEU Population Health Sciences

Association of common genetic variation with diseases and traits

- Genome wide association studies led to the discovery of >70,000 genetic associations to common diseases and traits
- 90% of GWAS variants are located in non-coding regions





Integration with other genomic layers



Zeggini et al. 2020

What does "epigenetics" mean?

Epigenetics is the study of mitotically heritable changes in gene expression that occur without changes in DNA sequence.



What influences DNA methylation?



mQTL as regulatory variant



mQTL: a genetic variant that is associated with DNA methylation at a CpG site.

DNA methylation dataset



- Arrays measure 3-5% of the genome
- On an array we measure methylation in a population of cells
- One CpG can be 0, 0.5 or 1
- Across a population we get a methylation level between 0 and 1

mQTL analysis - analogous to GWAS for continuous traits

Linear regression equation $Y_i = b_0 + b_1 X_i + e_i$ where

 Y_i = continuous trait value for individual i X_i = number of 'A' alleles an individual has



	sample 1	sample 2	sample 3	samp
rs1050745	1	2	1	
rs10507452	0	0	0	
rs10507456	2	1	1	
rs10507457	2	1	2	
rs10507458	0	1	0	
rs10507462	2	2	2	
rs10507463	2	2	2	
rs10507465	0	1	1	
rs10507466	1	1	2	
rs10507467	1	2	1	
rs10507481	1	2	1	
rs10507482	1	0	2	
rs10507483	0	0	0	
rs10507489	1	2	2	
rs10507490	2	2	2	
rs10507491	0	0	0	
rs10507495	0	0	0	
rs10507496	2	1	2	

Association test is whether $b \neq 0$

Whole genome mQTL analysis is a GWAS for each CpG



800,000 CpGs x 10M SNPs 800,000 GWAS!

Why federated analyses?

- Larger sample sizes are needed
 - Combining multiple studies
- Relevance of distant associations
 - To increase power to identify distant mQTLs have smaller effect sizes
- Reproducible and efficient
 - Sharing code improves research quality -> no analysis plans
 - Consistency across large-scale analyses of genetic and DNAm datasets
 - Automated control of confounding and batch effects
- To enable collaboration
 - No sharing of individual data

Genetics of DNA methylation Consortium

>60 cohorts with genetic and epigenetic data on blood samples

EPIC

450k

Total

37

23

60



65186

21022

86208

Inc	lucion	oritoria
	IUSIOII	Cillena

- Genetic data
- Methylation array data
- Blood samples/sorted celltypes
- N>200

20 15

European

- Multi-ancestry
 - Diverse epigenetic epidemiology partnership (DEEP)
 - Fine mapping

Github repository

https://github.com/genetics-of-dna-methylation-consortium/godmc_phase2/wiki



- Input data folder
- Config file
- Processed data
- Results
- Results upload
- 55 cohorts ran first module

'ING'S



GoDMC2 wiki

Home

ejh243 edited this page on Jan 17 \cdot 12 revisions

Welcome to the GoDMC Phase 2 pipeline. This wiki will guide you through the analysis. The pipeline has been designed so that minimal coding is required by analysts, which hopefully ensures that it is quick to execute, avoids duplication of efforts, increases reproducibility and reduces potential for error. The main objectives of the pipeline are:

Perform a mQTL analysis

Perform a full GWAS on every methylation probe available in the sample. Variations of this analysis will be performed, including:

- a standard mQTL analysis using bestguess imputed SNPs covering the full surface using new software
- cell-type interacting mQTL analyses for abundant cell types
- an inversion-mQTL analysis using inversions inferred from SNP data
- a var-mQTL analysis which looks for SNPs that influence the variance of a probe.
- a sex-stratified mQTL analysis for Chromosome X and Y

IMPORTANT: Please note that in GoDMC Phase 2 we are focusing on blood samples of Europeans only that are profiled with EPIC v1 or EPIC v2 arrays. Datasets should have at least 200 individuals and ideally be imputed with the HRC panel.

GWAS of DNAm derived phenotypes

We will be performing a GWAS on:

- biological age proxies derived from epigenetic clocks
- cumulative smoking exposure derived from DNA methylation data
- · blood cell type proportions derived by DNA methylation
- MZ twinning

EWAS of polygenic risk scores

We will be conducting PRS EWASs on:

Edit New page





https://github.com/genetics-of-dna-methylation-consortium/godmc_phase2

Example of module



Example Module 1



Age check

mQTL Age distribution



GoDMC2 mQTL catalog

Software/Package	Language	Approx Speed	Approx Storage
MatrixEQTL	R	<1 day	>500T
TensorQTL	python3	~1-2 days	>500T
HASE	python2	~1-2 days	~100G
OSCA	C++	~1 day	30T

HASE can store complete surface -> partial derivatives

Building a mQTL resource for the community



Research:

. . .

- Open source mQTL catalog
- Comparison with variance mQTLs and interaction mQTLs
- Causal relationship with clinical traits
- Genetic correlations
- Natural selection
- Functional properties

QC-systematic heterogeneity



Meta-regression Number of SNPs Number of CpGs Sample size Methylation array Average MAF Average Info Lambda Average beta value across all probes Average sd across all probes Ancestry (UK vs Non UK, NL vs non NL etc. Average age Fraction of males Relatedness yes/no Cord vs whole blood Height (N=24) BMI (N=24) Normalization method

• Used 427 independent trans SNPs with p<10-14

Diverse DNA methylation data to understand disease discordance



The Diverse Epigenetic Epidemiology Partnership (DEEP) study aims to improve population health by exploring the effects of genomic and environmental diversity on DNA methylation and disease risk across the global human population



Trait-methylation association discovery





DEEP Study



Medical Research Council



DIVERSE EPIGENETIC EPIDEMIOLOGY PARTNERSHIP Pls Josine Min & Hannah Elliott (UoB) Co-Is Prachand Issarapu, Andrew Prentice (LSHTM/MRCG), Giriraj Chandak (CCMB) & 20 partners



www.deep-epigenetics.org

DATASHIELD nondisclosive federated analysis



- Pooled analysis
- No individual level data sharing

Non genetic variation of DNA methylation

Nondisclosive Principal component analysis



- Batch and cohort effects
- Heritable and nonheritable probes
- Remove genetic component
- Correlate PCs to

. . .

- Cell counts
- Inflammation measures (IL6)
- Age at infection
- **Climate variation**

• Inform future analysis of health outcomes by identifying factors that contribute to inter-individual DNAm variation

Non genetic variation of DNA methylation

- Comparison between populations
 - EWAS of ancestry/ethnicity
 - EWAS of environmental exposures
 - Ittle variation within population
 - large variation between populations
 - age at infection (CMV, EBV)





Acknowledgements

IEU

Caroline Relton Gibran Hemani Hannah Elliott Tom Gaunt Haotian Tang

Diverse Epigenetic Epidemiology Partnership

OGY



GoDMC Developer's team GoDMC core team

