

Supplemental Information for “Human methylome variation across Infinium 450K data on the Gene Expression Omnibus”

Sean K. Maden Reid F. Thompson Kasper D. Hansen Abhinav Nellore

Contents

Gene Expression Omnibus (GEO) data	2
Sample identification	2
Sample data acquisition	2
DNAm assay data	2
Sample metadata	2
Preprocessing – GSE-wise annotations	2
Postprocessing – Learning formatted labels with regular expressions	2
Sample type predictions	3
DNAm model-based predictions	3
Metadata column descriptions	3
Quality and summary metrics	4
BeadArray metrics	4
Minfi quality metrics	4
Detection p-values	4
Methylated and unmethylated log2 median signals	4
Genetic relatedness	4
Statistical analyses	5
Overview	5
Scripts and code	5
Principal component analyses	5
GEO year-wise sample and study availability	5
Metadata concordance analyses and plot	5
Analyzing control metric signal	5
Analyzing PCA across samples and probes	6
DNAm variability analyses in 7 tissues	6
Works cited	6

This document contains supplemental information for the manuscript “Human methylome variation across Infinium 450K data on the Gene Expression Omnibus”. It contains information about data access, aggregation, and analyses. For additional information, consult documentation in the `recountmethylation` R/Bioconductor package. For details about data files and running scripts, consult the `inputs-and-outputs.Rmd` vignette.

Gene Expression Omnibus (GEO) data

We identified DNAm array samples published to GEO and available in the GEO Data Sets database as of March 31, 2019.

Sample identification

We used the Entrez Utilities software (v.10.9) to quantify DNAm array samples available by year (Figures 1 and S1, see script `edirect_query.py`).

Sample data acquisition

We developed server software in Python to automate and version downloads of GSM IDATs and GSE SOFT files from the GEO Data Sets database. The software, `recount-methylation-server`, is freely available as a GitHub repo.

DNAm assay data

We extracted DNAm assay data. Signals were read from sample IDATs into an R session using the `minfi` package (v.1.29.3, Aryee et al. (2014)). Descriptions of the principal data types are as follows:

- Red signal: Red color channel signal(s) for beads, read from sample IDATs.
- Green signal: Green color channel signal(s) for beads, read from sample IDATs.
- Methylated signal: Methylated signal amount calculated from red and green signals.
- Unmethylated signal: Unmethylated signal amount, calculated from red and green signals.
- Beta-value: Fraction of DNAm ($M/[U + M + e]$) calculated from red and green signals.

We calculated quality metrics from raw/unnormlized signals, using red and green signals for BeadArray metrics and the log2 array medians of methylated and unmethylated signals. DNAm model-based predictions used the noob-normalized Beta-values (see below). These DNAm array data were stored as databases and are accessible using the companion `recountmethylation` R/Bioc package.

Sample metadata

Metadata preprocessing and postprocessing was performed for downloaded GSE SOFT files, which resulted in the metadata table used for manuscript analyses (Table S1, see below).

Preprocessing – GSE-wise annotations

Preprocessing was performed on free text tags and labels extracted from sample characteristics columns in the GSE SOFT files. This was performed on a GSE-wise basis, where the most frequent and informative labels were retained. Care was taken to retain only sample-specific information.

Postprocessing – Learning formatted labels with regular expressions

Preprocessed metadata were postprocessed by string pattern matching with regular expressions (see script `md_postprocessing.R`). Labels were applied hierarchically, wherein specific labels have general labels applied automatically and separated by semicolons. Resulting labels are all lowercase, and underscores replace spaces for applicable labels. For example, a label of “peripheral blood” is converted into “peripheral_blood;blood”.

Sample type predictions

We scraped GSM metadata from downloaded SOFT files into JSON-formatted files. We removed GSE-specific data like experiment title and methods text. For example, the filtered JSON file for sample GSM937258 was:

```
[
{
  "!Sample_characteristics_ch1": "tissue site: Bone",
  "!Sample_characteristics_ch1.1": "tissue type: Tumor",
  "!Sample_characteristics_ch1.2": "sample id (short form for plotting): 24-6b",
  "!Sample_characteristics_ch1.3": "description: Xiphoid Met",
  "!Sample_source_name_ch1": "Prostate cancer metastasis",
  "!Sample_title": "Prostate cancer metastasis 16032"
}
]
```

JSON files were run using a fork of the MetaSRA-pipeline software repo (Bernstein et al. (2017)). This produces mapped ontology term labels and predictions for a series of sample types. The most likely sample types were retained under the “sampletype” column in the postprocessed metadata. For applicable samples, cell line annotations were appended from the Cellosaurus resource (see below).

DNAm model-based predictions

We performed model-based predictions of sex, age, and cell type fractions using noob-normalized DNAm (Beta-values, see below) compiled from available GSM IDATs. Sex and cell fractions were predicted using the `getSex` and `estimateCellCounts` functions from the `minfi` package, and age predictions were performed using the `agep` function from the `watermelon` (v.1.28.0) package. The `getSex` and `estimateCellCounts` functions perform light preprocessing that may be influenced by input sample order. Repeated calculations with randomized sample ordering showed > 95% concordance for both predicted variables across repetitions.

Metadata column descriptions

Descriptions of the sample metadata columns (Table S1) are as follows:

1. `gsm` : Sample record
2. `gsm_title` : Sample record title
3. `gseid` : Study Record ID
4. `disease` : Learned disease or study group term
5. `tissue` : Learned tissue type term
6. `sampletype` : Most likely sample type from MetaSRA-pipeline (`msraptype` is type label, `msrapconf` is prediction confidence, and `ccid`, `ccacc`, and `cccat` are cellosaurus id, accession, and attribute info where applicable)
7. `arrayid_full` : Full array ID (format: `Sentrix ID_chip coordinate`)
8. `basename` : Basename for sample IDATs (e.g. common name for paired IDAT files)
9. `age` : Age from mined sample metadata (`valm` is age value, `unitm` is age units)
10. `predage` : Age (years) predicted from norm. Beta-values using `watermelon::agep()`
11. `sex` : Sex from mined sample metadata (M = male, F = female)
12. `predsex` : Sex predicted from unnorm. signal using `minfi::getSex()` (M = male, F = female)
13. `predcell.CD8T` : Predicted CD8T cell fraction from `minfi::estimateCellCounts()`
14. `predcell.CD4T` : Predicted CD4T cell fraction from `minfi::estimateCellCounts()`
15. `predcell.NK` : Predicted Natural Killer cell fraction from `minfi::estimateCellCounts()`
16. `predcell.Bcell` : Predicted Bcell fraction from `minfi::estimateCellCounts()`
17. `predcell.Mono` : Predicted Monocyte fraction from `minfi::estimateCellCounts()`
18. `predcell.Gran` : Predicted Granulocyte fraction from `minfi::estimateCellCounts()`

19. storage : Sample storage procedure, annotated from metadata (FFPE or F)

Quality and summary metrics

Several types of quality and summary metrics were generated from the DNAm assay data (Table S2). These were used to determine metric performances, performance difference across preparations, and study-wise sample sub-threshold frequencies (Figure 2).

BeadArray metrics

Seventeen BeadArray metrics (Table S2, columns 3-19) were calculated from red and green signals (Table S3, see script `beadarray_cgctrlmetrics.R`). Metric formulae and thresholds were determined by consulting the platform documentation ((“Illumina Genome Studio Methylation Module V1.8” 2010), (“BeadArray Controls Reporter Software Guide” 2015)) and related functions in the `ewastools` package (Heiss and Just (2018)), with the following notes on calculations:

- Use extension Grn A/T probes for system background.
- For metrics where denominators would be 0 for some samples, use a uniform denominator offset of 1.
- Use C and U 1-3 and 4-6 for Bisulfite Conversion I.
- Use probe address “34648333” and “43603326” (DNP, Biotin subtype) for Biotin Staining Background.
- For Specificity I, use PM, MM 1-3 for green signal, 4-6 for red signal.
- For Specificity II, use just probes S1-3, as probe S4 unavailable in control probe annotation.

Minfi quality metrics

Additional quality metrics were calculated with the function `.buildControlMatrix450k` from the `minfi` package (Table S2, columns 20-61, Fortin et al. (2014)).

Detection p-values

The quantities of probes below 3 detection p-value cutoffs (0.01, 0.05, and 0.1) were determined (Table S2, columns 69-71). Detection p-values were calculated using the `detectionP` function from the `minfi` package.

Methylated and unmethylated log2 median signals

Low log2 medians for M and U signal may indicate poor sample quality (Aryee et al. (2014)). We calculated methylated (M) and unmethylated (U) signal log2 medians (Table S2, columns 72-73) from the raw/unnormalized signals.

Genetic relatedness

Sample genetic identities were calculated for GSE records using the `call_genotypes` function from the `ewastools` package (Heiss and Just (2018)). This method uses array probes mapping to high-frequency SNPs for the HM450K platform in a probabilistic model to determine whether samples share the same genetic identity. GSM IDs and number of GSM records sharing the same genetic identity from the same GSE record were determined (Table S2, columns 74-75).

Statistical analyses

Overview

In summary, we generated summary statistics using the `scikit-learn` Python package or base R functions, conducted statistical tests using base R, and visualized data using base R and the `ggplot2` R package. Unless noted otherwise, p-value adjustments used the Benjamini-Hochberg method. Group label enrichments were calculated with Binomial tests, and correlations used Spearman tests.

Scripts and code

Code to reproduce analyses are provided in the R and `inst/scripts` directories. Descriptions of script inputs and outputs are summarized in the `inputs-and-outputs` vignette. For additional analysis examples, see the “`data_analyses`” vignette in the `recountmethylation` package.

Principal component analyses

Approximate array-wide principal component analysis was performed using noob-normalized Beta-values from autosomal probes. Feature hashing, or the “hashing trick” (Weinberger et al. (2010)), was used as an intermediate dimensionality reduction step (see data files in `inst/extdata/pca_fh1k_all_gsm35k/` directory). This was implemented with the following Python function:

```
def feature_hash(arr, target_dim=1000):
    low_d_rep = [0 for _ in range(target_dim)]
    for i, el in enumerate(arr):
        hashed = mmh3.hash(str(i))
        if hashed > 0:
            low_d_rep[hashed % target_dim] += arr[i]
        else:
            low_d_rep[hashed % target_dim] -= arr[i]
    return low_d_rep
```

GEO year-wise sample and study availability

Yearly GSM and GSE record quantities were obtained by platform using the Entrez Utilities software (see script). Data were plotted using the `ggplot2` package.

Metadata concordance analyses and plot

For mined and predicted sex, we calculated the fraction concordance. For mined and predicted age, we used ANOVAs to calculate variance percentages and p-values for multivariate models including GSE ID, cancer status, and sex. Predicted age and chronological age were regressed to calculate R-squared values and correlated using Spearman’s test. These were plotted with the regression model using `ggplot2` (Figure 1b). See example 1 in the `data_analyses` vignette of the `recountmethylation` package for details.

Analyzing control metric signal

We assessed control metric signal (BeadArray metrics and methylated or unmethylated log2 medians) using base R and certain packages. We calculated metadata term enrichment for the `tissue` and `disease` annotations using the Bimodal test with multiple testing adjustment using the Benjamini-Hochberg method

(p -adjusted $< 1e-3$ significance). Significance of signal differences between storage types (frozen and FFPE) were calculated using T-tests. Stacked barplots, confidence intervals, and violin plots in Figure 2 were generated using `ggplot2`. Study-wise sub-threshold frequency heatmaps were generated using `ComplexHeatmap`.

Analyzing PCA across samples and probes

We performed PCA using the feature-hashed probe datasets and the `prcomp` R function with default arguments (see above). Plots were generated using the `factoextra` package and base R.

DNAm variability analyses in 7 tissues

Samples were categorized as being from 1 of 7 tissues by manual review of studies (Table S6). Likely cancers or samples from cancer patients were removed from consideration. We removed samples occurring in the bottom 5th quantile for *both* methylated and unmethylated log2 median signal, samples failing at least 1 BeadArray quality metric, and likely replicates or samples with shared genetic identity (see above, Table S2). After metadata and quality filters, analysis was restricted to the 7 remaining tissues with at least 100 remaining samples from at least 2 studies (Figure S5b).

Because few targeted studies sampled across tissues and because tissue groups explained less variation than study/GSE ID, we performed preprocessing within tissue groups in parallel. Linear adjustment on study ID was performed using the `removeBatchEffect` function from the `limma` (v.3.39.12) package. Probes were then removed if they showed significant (p -adj. < 0.01) and substantial ($> 10\%$) variance contributions in ANOVAs. ANOVA models consisted of GSE ID and DNAm model predictions for age, sex, and blood cell fractions (see above, Table S1). This filter removed 8 - 40% of autosomal probes across the 7 tissues, respectively (Figure S5c).

Probes with low variance across all 7 tissues were identified by taking the union of probes in the lowest 5th quantile variance for each tissue (Table S7, see script). The top 2,000 probes with tissue-specific high variances (14,000 probes total, Table S8) were selected from an absolute quantile filter and a binned quantile filter, where probes showed the highest 5th quantile variance and occurred in just 1 of the 7 tissue groups for each method (see script).

Works cited

- Aryee, Martin J., Andrew E. Jaffe, Hector Corrada-Bravo, Christine Ladd-Acosta, Andrew P. Feinberg, Kasper D. Hansen, and Rafael A. Irizarry. 2014. "Minfi: A Flexible and Comprehensive Bioconductor Package for the Analysis of Infinium DNA Methylation Microarrays." *Bioinformatics* 30 (10): 1363–9. <https://doi.org/10.1093/bioinformatics/btu049>.
- "BeadArray Controls Reporter Software Guide." 2015, October. https://support.illumina.com/content/dam/illumina-support/documents/documentation/chemistry_documentation/infinium_assays/infinium_hd_methylation/beadarray-controls-reporter-user-guide-1000000004009-00.pdf.
- Bernstein, Matthew N., AnHai Doan, Colin N. Dewey, and Jonathan Wren. 2017. "MetaSRA: Normalized Human Sample-Specific Metadata for the Sequence Read Archive." *Bioinformatics* 33 (18): 2914–23. <https://doi.org/10.1093/bioinformatics/btx334>.
- Fortin, Jean-Philippe, Aurélie Labbe, Mathieu Lemire, Brent W Zanke, Thomas J Hudson, Elana J Fertig, Celia MT Greenwood, and Kasper D Hansen. 2014. "Functional Normalization of 450k Methylation Array Data Improves Replication in Large Cancer Studies." *Genome Biology* 15 (11): 503.
- Heiss, Jonathan A., and Allan C. Just. 2018. "Identifying Mislabeled and Contaminated DNA Methylation Microarray Data: An Extended Quality Control Toolset with Examples from GEO." *Clinical Epigenetics* 10 (June). <https://doi.org/10.1186/s13148-018-0504-1>.

“Illumina Genome Studio Methylation Module V1.8.” 2010. Illumina. https://support.illumina.com/content/dam/illumina-support/documents/documentation/software_documentation/genomestudio/genomestudio-2011-1/genomestudio-methylation-v1-8-user-guide-11319130-b.pdf.

Weinberger, Kilian, Anirban Dasgupta, Josh Attenberg, John Langford, and Alex Smola. 2010. “Feature Hashing for Large Scale Multitask Learning.” *arXiv:0902.2206 [Cs]*, February. <http://arxiv.org/abs/0902.2206>.