

Estimating DNA methylation levels and finding differentially methylated regions using the charm package

Martin Aryee*, Rafael Irizarry

February, 2010

Johns Hopkins School of Medicine / Johns Hopkins School of Public Health
Baltimore, MD, USA

1 Introduction

The Bioconductor package **charm** can be used to analyze DNA methylation data generated using McrBC fractionation and two-color Nimblegen microarrays. It is customized for use with the from the custom CHARM microarray [1], but can also be applied to many other Nimblegen designs.

Functions include:

- Quality control
- Finding suitable control probes for normalization
- Percentage methylation estimates
- Identification of differentially methylated regions

As input we will need raw Nimblegen data (.xys) files and a corresponding annotation package built with pdInfoBuilder. This vignette uses the following packages:

- **charm**: contains the analysis functions
- **charmData**: an example dataset
- **pd.charm.hg18.example**: the annotation package for the example dataset
- **BSgenome.Hsapiens.UCSC.hg18**: A BSgenome object containing genomic sequence used for finding non-CpG control probes

Each sample is represented by two xys files corresponding to the untreated (green) and methyl-depleted (red) channels. The 532.xys and 635.xys suffixes indicate the green and red channels respectively.

*aryee@jhu

2 Analyzing data from the custom CHARM microarray

Load the charm package:

```
R> library(charm)
```

3 Read in raw data

Get the name of your data directory (in this case, the example data):

```
R> dataDir <- system.file("data", package = "charmData")
```

```
R> dataDir
```

```
[1] "/thumper/ctsa/genomicsR/install/R-2.10.0/lib64/R/site-library/charmData/data"
```

First we read in the sample description file:

```
R> pd <- read.delim(file.path(dataDir, "phenodata.txt"))
```

```
R> pd
```

	filename	sampleID	tissue
1	136421_532.xys	441_liver	liver
2	136421_635.xys	441_liver	liver
3	136600_532.xys	449_spleen	spleen
4	136600_635.xys	449_spleen	spleen
5	3788602_532.xys	449_liver	liver
6	3788602_635.xys	449_liver	liver
7	3822402_532.xys	441_spleen	spleen
8	3822402_635.xys	441_spleen	spleen
9	5739902_532.xys	624_colon	colon
10	5739902_635.xys	624_colon	colon
11	5875602_532.xys	441_colon	colon
12	5875602_635.xys	441_colon	colon

A valid sample description file should contain at least the following (arbitrarily named) columns:

- a filename column
- a sample ID column
- a group label column (optional)

The sample ID column is used to pair the methyl-depleted and untreated data files for each sample. The group label column is used when identifying differentially methylated regions between experimental groups.

The `validatePd` function can be used to validate the sample description file. When called with only a sample description data frame and no further options `validatePd` will try to guess the contents of the columns.

```
R> res <- validatePd(pd)
```

```
FileNames:
```

```
  fileNameColumn not specified. Trying to guess.
```

```
  OK - Found in column 1
```

```
Sample names:
```

```
  sampleNameColumn column not specified. Trying to guess.
```

```
  OK - Found in column 2
```

```
Group labels:
```

```
  groupColumn column not specified. Trying to guess.
```

```
  OK - Found in column 3
```

Now we read in the raw data. The `readCharm` command makes the assumption (unless told otherwise) that the two xys files for a sample have the same file name up to the suffixes 532.xys (untreated) and 635.xys (methyl-depleted).

```
R> rawData <- readCharm(files = pd$filename, path = dataDir,  
  sampleKey = pd)
```

```
Checking designs for each XYS file... Done.
```

```
Allocating memory... Done.
```

```
Reading /thumper/ctsa/genomicsR/install/R-2.10.0/lib64/R/site-library/charmData/data/136421_
```

```
Reading /thumper/ctsa/genomicsR/install/R-2.10.0/lib64/R/site-library/charmData/data/136600_
```

```
Reading /thumper/ctsa/genomicsR/install/R-2.10.0/lib64/R/site-library/charmData/data/3788602_
```

```
Reading /thumper/ctsa/genomicsR/install/R-2.10.0/lib64/R/site-library/charmData/data/3822402_
```

```
Reading /thumper/ctsa/genomicsR/install/R-2.10.0/lib64/R/site-library/charmData/data/5739902_
```

```
Reading /thumper/ctsa/genomicsR/install/R-2.10.0/lib64/R/site-library/charmData/data/5875602_
```

```
Checking designs for each XYS file... Done.
```

```
Allocating memory... Done.
```

```
Reading /thumper/ctsa/genomicsR/install/R-2.10.0/lib64/R/site-library/charmData/data/136421_
```

```
Reading /thumper/ctsa/genomicsR/install/R-2.10.0/lib64/R/site-library/charmData/data/136600_
```

```
Reading /thumper/ctsa/genomicsR/install/R-2.10.0/lib64/R/site-library/charmData/data/3788602_
```

```
Reading /thumper/ctsa/genomicsR/install/R-2.10.0/lib64/R/site-library/charmData/data/3822402_
```

```
Reading /thumper/ctsa/genomicsR/install/R-2.10.0/lib64/R/site-library/charmData/data/5739902_
```

```
Reading /thumper/ctsa/genomicsR/install/R-2.10.0/lib64/R/site-library/charmData/data/5875602_
```

```
R> rawData
```

```
TilingFeatureSet2 (storageMode: lockedEnvironment)
```

```
assayData: 243129 features, 6 samples
```

```
  element names: channel1, channel2
```

```
phenoData
```

```
  sampleNames: 136421, 136600, ..., 5875602 (6 total)
```

```
  varLabels and varMetadata description:
```

```
    sampleID: NA
```

```
    tissue: NA
```

```
    ...: ...
```

```

channel2DateTime: date/time from raw files
(6 total)
additional varMetadata: channel
featureData
  featureNames: 1, 2, ..., 243129 (243129 total)
  fvarLabels and fvarMetadata description: none
experimentData: use 'experimentData(object)'
Annotation: pd.charm.hg18.example

```

4 Array quality assessment

We can calculate array quality scores and generate a pdf report with the `qcReport` command.

A useful quick way of assessing data quality is to examine the untreated channel where we expect every probe to have signal. Very low signal intensities on all or part of an array can indicate problems with hybridization or scanning. The CHARM array and many other designs include background probes that do not match any genomic sequence. Any signal at these background probes can be assumed to be the result of optical noise or cross-hybridization. Since the untreated channel contains total DNA a successful hybridization would have strong signal for all untreated channel genomic probes. The array signal quality score (`pmSignal`) is calculated as the average percentile rank of the signal probes among these background probes. A score of 100 means all signal probes rank above all background probes (the ideal scenario).

```

R> qual <- qcReport(rawData, file = "qcReport.pdf")
R> qual

```

	pmSignal	sd1	sd2
136421	78.56437	0.1950274	0.1932112
136600	81.46541	0.1755225	0.1227921
3788602	83.95419	0.1249030	0.2409803
3822402	81.43751	0.1180708	0.1824810
5739902	82.55727	0.1490854	0.2035761
5875602	79.38069	0.3130266	0.3962373

The PDF quality report is shown in Appendix A. Three quality metrics are calculated for each array:

1. Average signal strength: the average percentile rank of untreated channel signal probes among the background (anti-genomic) probes.
2. Untreated channel signal standard deviation. The array is divided into a series of rectangular blocks and the average signal level calculated for each. Since probes are arranged randomly on the array there should be no large differences between blocks. Arrays with spatial artifacts have a large standard deviation between blocks.

3. Methyl-depleted channel signal standard deviation.

5 Percentage methylation estimates and differentially methylated regions (DMRs)

We now calculate probe-level percentage methylation estimates for each sample. As a first step we need to identify a suitable set of unmethylated control probes from CpG-free regions to be used in normalization.

```
R> library(BSgenome.Hsapiens.UCSC.hg18)
R> ctrlIdx <- getControlIndex(rawData, subject = Hsapiens)
```

The minimal code required to estimate methylation would be `p <- methp(rawData, controlIndex=ctrlIdx)`. However, it is often useful to get `methp` to produce a series of diagnostic density plots to help identify non-hybridization quality issues. The `plotDensity` option specifies the name of the output pdf file, and the optional `plotDensityGroups` can be used to give groups different colors.

```
R> grp <- pData(rawData)$tissue
R> p <- methp(rawData, controlIndex = ctrlIdx, plotDensity = "density.pdf",
  plotDensityGroups = grp)
```

Spatial normalization
Background removal
Within sample normalization: loess
Between sample normalization: quantile
Estimating percentage methylation

```
R> head(p)

      136421    136600    3788602    3822402    5739902
1 0.17250259 0.3895986 0.3881209 0.5793288 0.3813294
2 0.84184931 0.6773743 0.3465772 0.8999109 0.5788529
3 0.09046145 0.0641674 0.1614183 0.1426585 0.2339750
4 0.77692120 0.4944354 0.4772154 0.4760221 0.3861171
5 0.69668343 0.5593289 0.4191725 0.4469628 0.4110470
6 0.66978815 0.7949903 0.7856860 0.7403811 0.8982202
      5875602
1 0.2708504
2 0.9183155
3 0.7293811
4 0.4904633
5 0.4008233
6 0.8522786
```

The density plots are shown in Appendix B.
We can now identify differentially methylated regions using `dmrFinder`:

```
R> dmr <- dmrFinder(rawData, p = p, groups = grp,
  compare = c("colon", "liver", "colon", "spleen"))
```

Computing group medians and SDs for 3 groups:

```
1
2
3
```

Done.

Smoothing:

	1%
=	1%
=	2%
=	3%
==	3%
==	4%
==	5%
===	5%
===	6%
===	7%
====	7%
====	8%
====	9%
=====	9%
=====	10%
=====	11%
=====	11%
=====	12%

	=====	13%
	=====	13%
	=====	14%
	=====	15%
	=====	15%
	=====	16%
	=====	17%
	=====	17%
	=====	18%
	=====	19%
	=====	19%
	=====	20%
	=====	21%
	=====	21%
	=====	22%
	=====	23%
	=====	23%
	=====	24%
	=====	25%
	=====	25%
	=====	26%
	=====	27%
	=====	27%

=====		28%
=====		29%
=====		29%
=====		30%
=====		31%
=====		31%
=====		32%
=====		33%
=====		33%
=====		34%
=====		35%
=====		35%
=====		36%
=====		37%
=====		37%
=====		38%
=====		39%
=====		39%
=====		40%
=====		41%
=====		41%
=====		42%
=====		43%

			43%
	=====		
	=====		44%
	=====		
	=====		45%
	=====		
	=====		45%
	=====		
	=====		46%
	=====		
	=====		47%
	=====		
	=====		47%
	=====		
	=====		48%
	=====		
	=====		49%
	=====		
	=====		49%
	=====		
	=====		50%
	=====		
	=====		51%
	=====		
	=====		51%
	=====		
	=====		52%
	=====		
	=====		53%
	=====		
	=====		53%
	=====		
	=====		54%
	=====		
	=====		55%
	=====		
	=====		55%
	=====		
	=====		56%
	=====		
	=====		57%
	=====		
	=====		57%
	=====		
	=====		58%

	=====	59%
	=====	59%
	=====	60%
	=====	61%
	=====	61%
	=====	62%
	=====	63%
	=====	63%
	=====	64%
	=====	65%
	=====	65%
	=====	66%
	=====	67%
	=====	67%
	=====	68%
	=====	69%
	=====	69%
	=====	70%
	=====	71%
	=====	71%
	=====	72%
	=====	73%
	=====	73%

	=====	74%
	=====	75%
	=====	75%
	=====	76%
	=====	77%
	=====	77%
	=====	78%
	=====	79%
	=====	79%
	=====	80%
	=====	81%
	=====	81%
	=====	82%
	=====	83%
	=====	83%
	=====	84%
	=====	85%
	=====	85%
	=====	86%
	=====	87%
	=====	87%
	=====	88%
	=====	89%

```

|
|=====| 89%
|=====| 90%
|=====| 91%
|=====| 91%
|=====| 92%
|=====| 93%
|=====| 93%
|=====| 94%
|=====| 95%
|=====| 95%
|=====| 96%
|=====| 97%
|=====| 97%
|=====| 98%
|=====| 99%
|=====| 99%
|=====| 100%

```

Done.

Finding DMRs for each pairwise comparison.

colon-liver..

colon-spleen..

Done

R> names(dmr)

```

[1] "tabs"      "p"         "l"
[4] "chr"       "pos"       "pns"
[7] "index"     "controlIndex" "gm"
[10] "groups"    "args"      "comps"
[13] "package"

```

```
R> names(dmr$tabs)

[1] "colon-liver" "colon-spleen"

R> head(dmr$tabs[[1]])
```

	chr	start	end	p1	p2
447	chr12	88272817	88273811	0.8804792	0.1408243
1773	chr6	52637747	52638747	0.7584872	0.1313882
492	chr13	27090247	27091263	0.8188426	0.1291102
151	chr11	14620645	14621065	0.8799767	0.3324493
643	chr15	58673117	58673711	0.8562095	0.2512823
795	chr17	198791	199552	0.8167333	0.2229383

	regionName	indexStart	indexEnd	area
447	chr12:88266873-88274292	41215	41238	17.751716
1773	chr6:52635302-52638967	163819	163843	15.677477
492	chr13:27090144-27095500	46022	46041	13.794648
151	chr11:14620645-14623686	28888	28900	7.117856
643	chr15:58669815-58674073	59008	59024	10.283762
795	chr17:198024-209044	74031	74047	10.094516

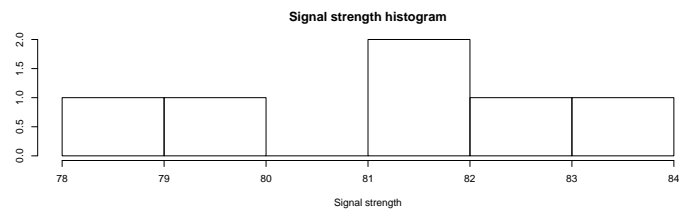
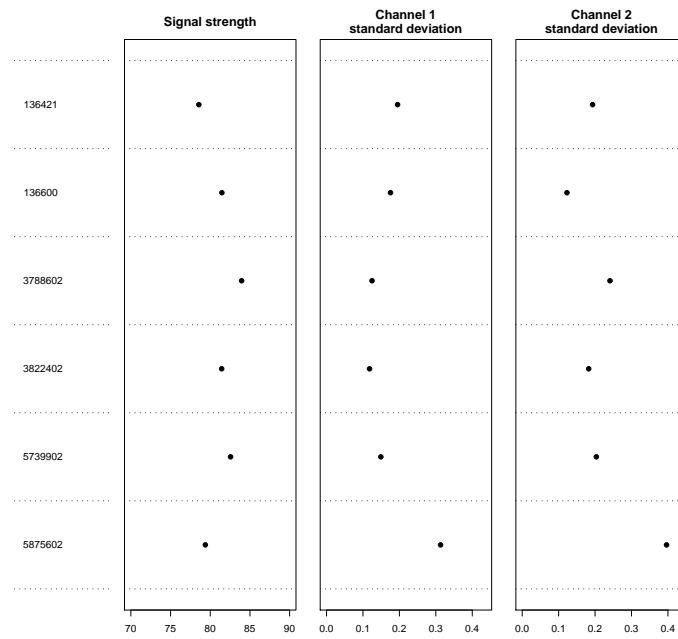
	ttarea
447	818.4631
1773	731.6463
492	711.4216
151	478.1902
643	476.9114
795	425.0093

When called without the `compare` option, `dmrFinder` performs all pairwise comparisons between the groups.

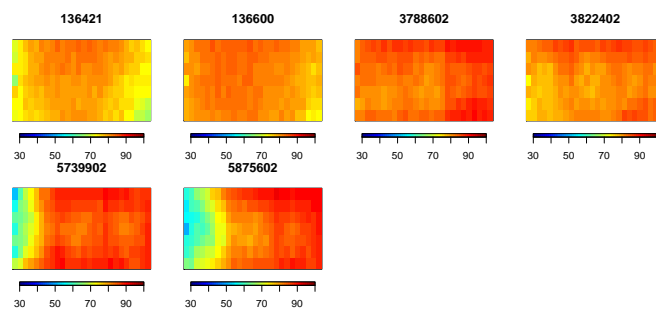
References

- [1] Irizarry et al. Comprehensive high-throughput arrays for relative methylation (charm). *Genome Research*, 18(5):780–790, 2008.

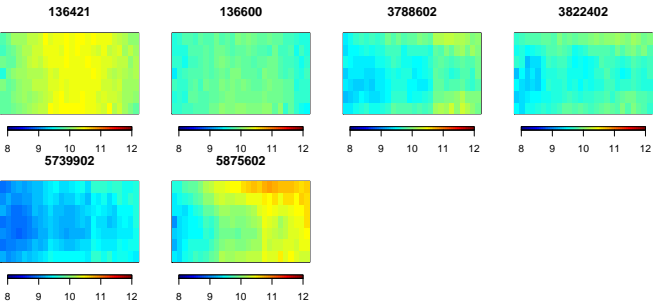
6 Appendix A: Quality report



Untreated Channel: PM probe quality

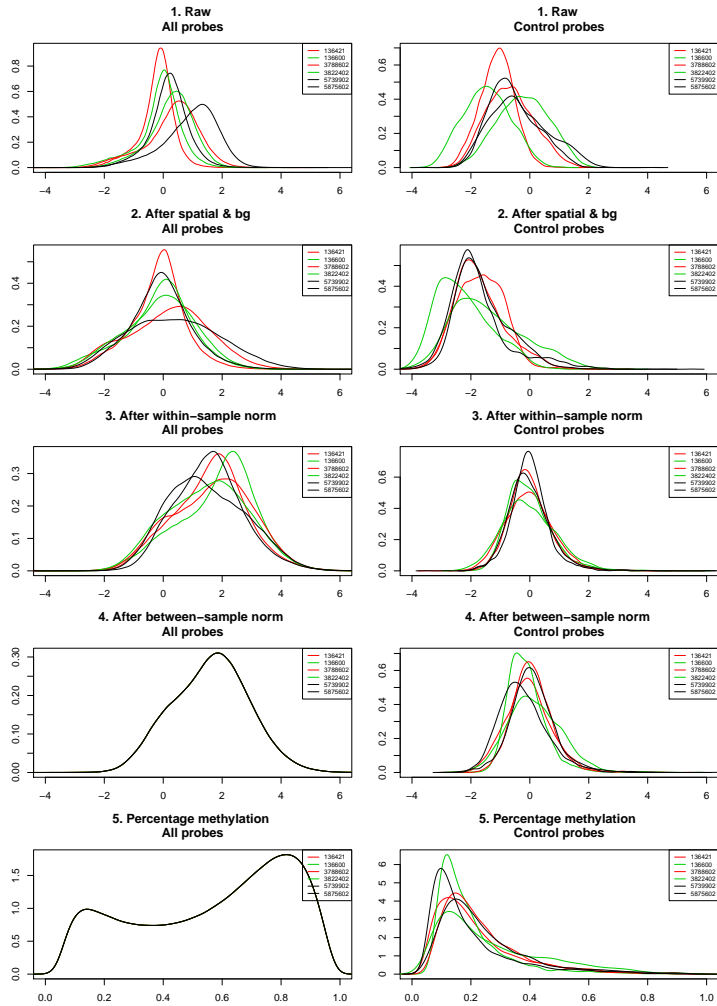


Enriched Channel: PM signal intensity



7 Appendix B: Density plots

Each row corresponds to one stage of the normalization process (Raw data, After spatial and background correction, after within-sample normalization, after between-sample normalization, percentage methylation estimates). The left column shows all probes, while the right column shows control probes.



8 Details

This document was written using:

```
R> sessionInfo()
```

```
R version 2.10.0 (2009-10-26)  
x86_64-unknown-linux-gnu
```

```
locale:
```

```
[1] LC_CTYPE=en_US.iso885915  
[2] LC_NUMERIC=C  
[3] LC_TIME=en_US.iso885915  
[4] LC_COLLATE=en_US.iso885915  
[5] LC_MONETARY=C  
[6] LC_MESSAGES=en_US.iso885915  
[7] LC_PAPER=en_US.iso885915  
[8] LC_NAME=C  
[9] LC_ADDRESS=C  
[10] LC_TELEPHONE=C  
[11] LC_MEASUREMENT=en_US.iso885915  
[12] LC_IDENTIFICATION=C
```

```
attached base packages:
```

```
[1] tools      stats      graphics  grDevices  utils  
[6] datasets  methods   base
```

```
other attached packages:
```

```
[1] BSgenome.Hsapiens.UCSC.hg18_1.3.16  
[2] BSgenome_1.14.2  
[3] Biostrings_2.14.8  
[4] IRanges_1.4.8  
[5] pd.charm.hg18.example_0.9.0  
[6] RSQLite_0.7-3  
[7] DBI_0.2-4  
[8] charm_0.9.34  
[9] SQN_1.0  
[10] nor1mix_1.1-1  
[11] mclust_3.3.2  
[12] fields_6.01  
[13] spam_0.15-5  
[14] oligo_1.10.2  
[15] preprocessCore_1.8.0  
[16] oligoClasses_1.8.0  
[17] Biobase_2.6.0
```

```
loaded via a namespace (and not attached):
```

```
[1] affxparser_1.18.0 affyio_1.14.0 gtools_2.6.1
[4] MASS_7.3-3        multtest_2.2.0 siggenes_1.20.0
[7] spatial_7.3-1      splines_2.10.0 survival_2.35-7
```