

Working notes

Effect of normalization on the prediction of DNA methylation status of MGMT promoter: example with HM-450K Infinium data from TCGA and the R package `mgmtstp27`.
(document in construction!)

BADY P., Hegi M.E.

2 octobre 2014

License : GPL version 2 or newer
Copyright (C) 2000-2014 Pierre Bady
This program/document is free software; you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation; either version 2 of the License, or (at your option) any later version.
This program/document is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY; without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

Résumé

This document contains the description and the use of some functions proposed in the R packages `mgmtstp27`. Additionally, it provides information related to the effect of normalization of HM450K/HM27k Infinium platform on the prediction of the DNA methylation status of the MGMT promoter (Bady et al. 2012).

Table des matières

1	Motivations	2
2	Data	2
3	Preprocessing and Normalization	2

4	Normalisation effect on the prediction	3
5	Comparison of the three datasets (PP-27K, PP-450K, TCGA-450K)	3
6	Normalization effect for the training dataset (M-GBM, [2])	7
7	Conclusion	9
8	Acknowledgments	9
9	Appendix	10
9.1	Import raw HM-27K data (format .IDAT)	10
9.2	Session	10

1 Motivations

In this document, we propose to evaluate the effect of the normalization of the data from Infinium HM-450K platform (DNA methylation) on the prediction of the DNA methylation of the MGMT promoter from the model proposed in [2].

2 Data

Dataset came from TCGA project (The Cancer Genome Atlas Research Network 2008, <http://cancergenome.nih.gov/>). The DNA methylation was evaluated by platform Infinium HM-450K. The first dataset come from the older version dated to 2012-05-25 where the level 1 (see TCGA documentation) directly contained the preprocessed information from 74 samples (e.g. unmethylated and methylated intensities). However, little information was provided to describe the normalization/preprocessing used to prepare this dataset. A second version (2012-07-31) of this data set in raw format (the information of the two colors is separated in two different files) were used to determine the normalization used in the initial dataset and to compare the preprocessing methods. The dataset used in [2] as training dataset (M-GBM), was analyzed in a similar way.

3 Preprocessing and Normalization

For the initial dataset, we have some doubts on the method used to preprocess the data. Concerning the new dataset (updated version), we used two different methods available in Genome Studio :

- "raw" version corresponding to the method initially used to prepare the data from HM-27K platform. Preprocessing means converting the Red and Green channel into unmethylated and methylated signal.
- The second method corresponds to a new method proposed by Illumina to preprocess the HM-450K data. The procedure includes background correction and normalization using a sample as reference (the second by default, see documentation of R package minfi, Kasper and Martin 2012).
- swamn
- Scaling correction for chemistry

The functions from the R package minfi (Kasper and Martin 2012) were used to perform both these normalizations. In this study, we didn't take into account the chemistry effect because the two probes considered in our model came from the chemistry I only. The R package lumi provided additional functions for normalization (Du and Lin 2008, not used here). Analyses and Graphical representations were performed using R-3.1.1 ([4]) and the R package minfi ([1]) and methylumi ([3]).

4 Normalisation effect on the prediction

The importation and preparation of the three datasets were relatively facilitated by the use of the function from R package minfi. The functions preprocessRaw and preprocessIllumina provided the two new datasets from the last update (see R code below). The dataset used in the table S4 ([2], R object called predTCGA450K) was built manually because the old structure of the level 1 data was not compatible with the functions of R packages minfi or methylumi.

```
#-----
# data importation
#-----
library(minfi)
library(IlluminaHumanMethylation450kmanifest)
# data importation
datadir <- paste(getwd(), "/JHU_USC__HumanMethylation450/Level_1/", sep="")
list.files(datadir)
infofile0 <- read.table("file_manifest.txt", h=TRUE, sep="\t")
infofile1 <- infofile0[infofile0$Level==1,]
rgset0 <- read.450k.exp(datadir)
# preprocessing
rawdata0 <- preprocessRaw(rgset0)
normdata0 <- preprocessIllumina(rgset0)
# meylation and unmethylation data
rawunmeth0 <- getUnmeth(rawdata0)
rawmeth0 <- getMeth(rawdata0)
normunmeth0 <- getUnmeth(normdata0)
normmeth0 <- getMeth(normdata0)
# table containing the probes used in the model
load("promoterprobes.rda")
rawunmeth1 <- rawunmeth0[promoterprobes,]
rawmeth1 <- rawmeth0[promoterprobes,]
normunmeth1 <- normunmeth0[promoterprobes,]
normmeth1 <- normmeth0[promoterprobes,]
mvalueraw1 <- log2((rawmeth1+1)/(rawunmeth1+1))
mvaluenorm1 <- log2((normmeth1+1)/(normunmeth1+1))
# initial dataset (74 samples used in the table S4 in Bady et al. 2012)
load("/export/scratch/data/monikaproject/TCGA6/DNAmethylation/450k/predTCGA450.rda")
```

5 Comparison of the three datasets (PP-27K, PP-450K, TCGA-450K)

In this section, we only kept the samples common to the three datasets. Consequently, we had three measures by probes for a given sample :

—
—

The dataset called PP-27K corresponds to the raw dataset from TCGA (update 2012-07-31) after classical preprocessing/normalization (that correspond to the normalization initially used for HM-27K platform). It contained 124 samples before matching step.

The dataset called PP-450K corresponds to the raw dataset from TCGA (update 2012-07-31) after "new" Illumina preprocessing. It contained 124 samples before matching step.

TCGA-450K corresponds to the dataset (update 2012-05-25) used for the prediction in the table S4 (Bady et al. 2012). It contained 74 samples.

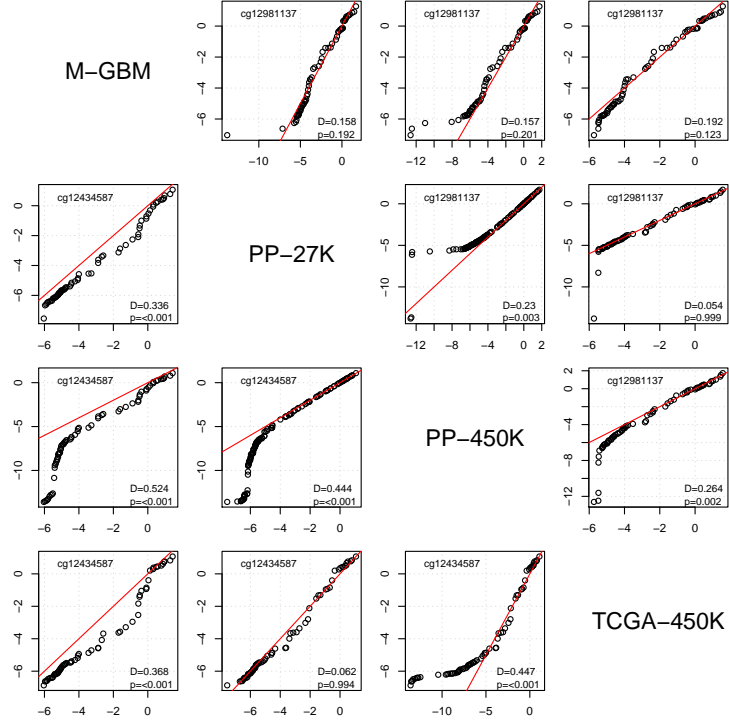


FIGURE 1 – Comparison of M-value distributions between the three "unmatched" datasets and the training dataset (M-GBM). The M-values of the probes cg12434587 and cg12434587 used in MGMT-STP27 were compared by quantile-quantile representation (QQ-plot). The red line corresponds to the line $y=x$. The terms 'D' and 'p' refer to the comparison of the distribution by the Kolmogorov-Smirnov test. The platform Illumina used is indicated for each dataset. When the p-value is inferior to 0.05, the two distributions are considered as significantly different.

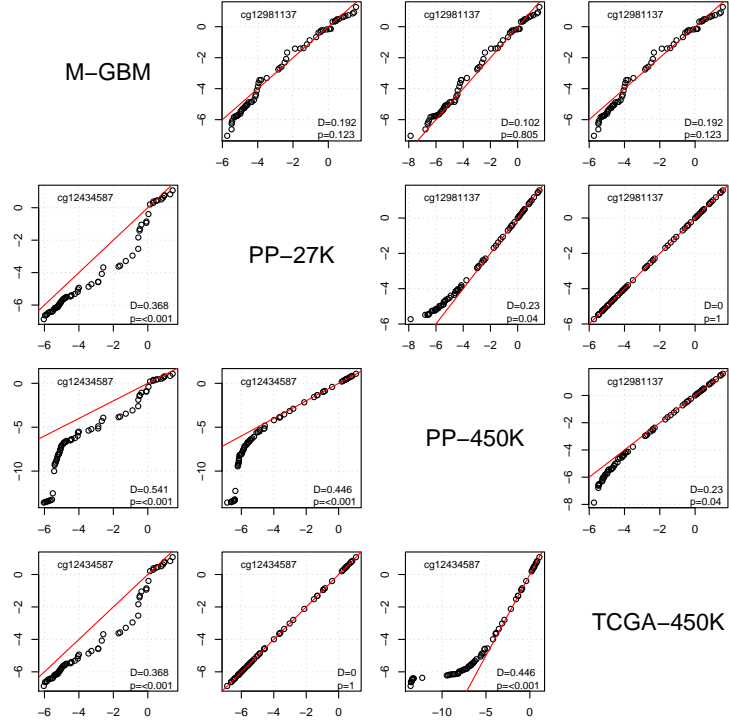


FIGURE 2 – Comparison of M-value distributions between the three "matched" datasets and the training dataset (M-GBM). The M-values of the probes cg12434587 and cg12434587 used in MGMT-STP27 were compared by quantile-quantile representation (QQ-plot). The red line corresponds to the line $y=x$. The terms 'D' and 'p' refer to the comparison of distribution by the Kolmogorov-Smirnov test. The platform Illumina used is indicated for each dataset. When the p-value is inferior to 0.05, the two distributions are considered as significantly different.

After matching based on the sample names, the three datasets contained 74 samples. The analyses in Figure 2 and Figure 3 show that the initial dataset (TCGA-450K) is exactly similar to the dataset normalized by the "raw" preprocessing (PP-27K). The procedure used to preprocess the initial dataset is certainly the same and corresponds to the procedure used to prepare the HM-27K datasets. Consequently, the prediction proposed in the table S4 (Bady et al. 2012) is the same.

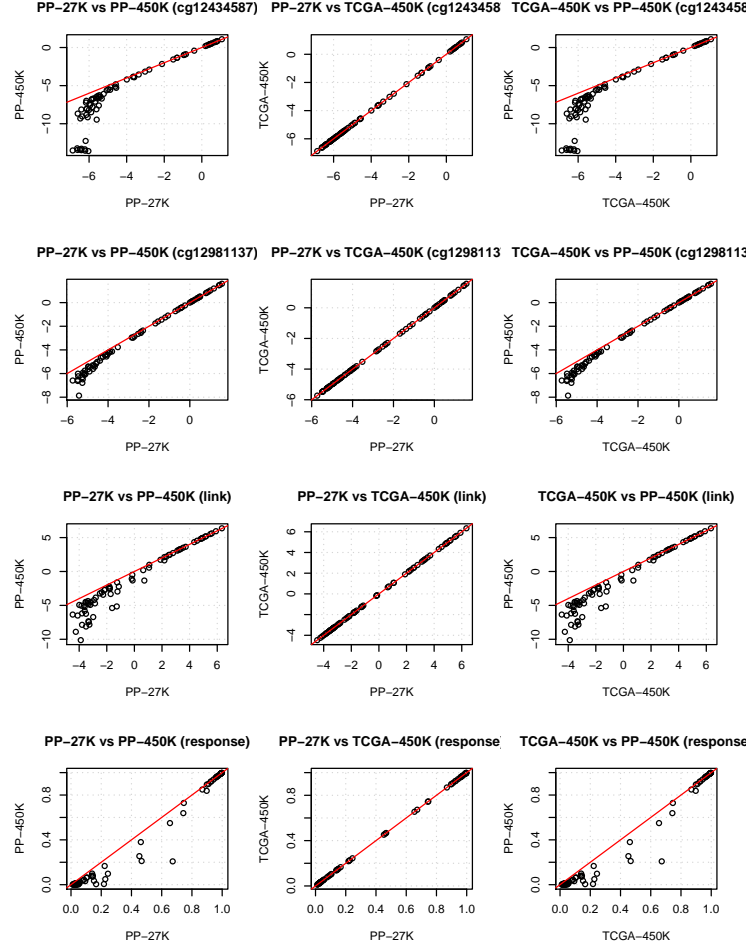


FIGURE 3 – Comparisons of the values of the both probes (cg12434587 and cg12981137) and predictions (link and response values) between the three matched dataset from TCGA.

The correlation between PP-27K and TCGA-450K datasets is perfect. The discrepancies between training (M-GBM) and PP-450K datasets were excessively increased by the normalization proposed in Genome Studio. The highest deviations between PP-27K and PP-450K were observed for the probe cg12434587 and they were mainly observed for the low M-values (Figure 3). The evaluation of the concordance between predicted statuses is provided below :

```

predraw2 <- predict(step27k,dfraw2,type="response")
mgmtraw2 <- ifelse(predraw2>=step27k$perf$cut,"M","U")
prednorm2 <- predict(step27k,dfnorm2,type="response")
mgmtnorm2 <- ifelse(prednorm2>=step27k$perf$cut,"M","U")
predini2 <- predict(step27k,predTCGA2,type="response")
mgmtini2 <- ifelse(predini2>=step27k$perf$cut,"M","U")
table(mgmtraw2,mgmtini2)
table(mgmtraw2,mgmtnorm2)
table(mgmtnorm2,mgmtini2)

```

We observe that the initial dataset was in perfect concordance with the dataset normalized by "raw" preprocessing. When the dataset was normalized by new Illumina procedure, we observe that three samples were not correctly classified.

6 Normalization effect for the training dataset (M-GBM, [2])

As previously, three datasets were considered in these analyses :

- M-PP-27K corresponds to the raw dataset after classical preprocessing/normalization (that correspond to the normalization initially used for HM-27K platform).
- M-PP-450K corresponds to the raw dataset after "new" Illumina preprocessing
- M-GBM-450K corresponds to the exact training dataset used to perform the model proposed in [2]

As expected, the results observed were very similar to the ones presented in the previous section. We observed deviations between PP-27K (identical to M-GBM dataset) and PP-450K for the both probes and they mainly observed the low M-values (Figure 4).

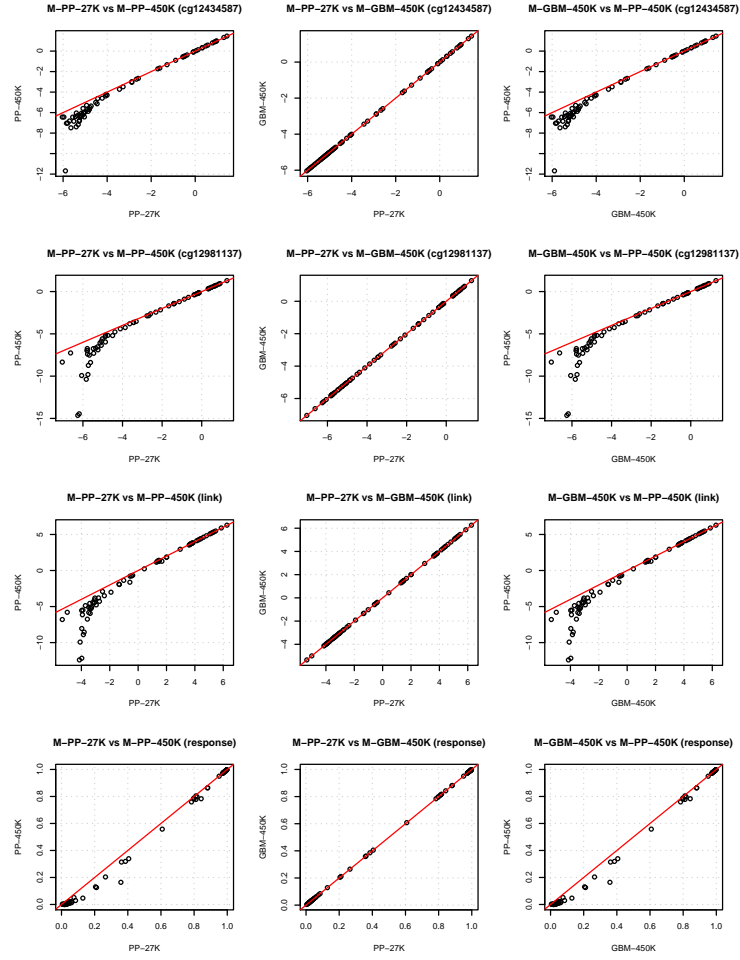


FIGURE 4 – Comparisons of the values of the both probes (cg12434587 and cg12434587) and predictions (link and response values) between the three matched datasets from M-GBM data used as training dataset in [2]

Four samples were not correctly classified. The evaluation of the concordance between the predicted statuses is provided below :

```
# comparison prediction
predraw2 <- predict(step27k,dfraw2,type="response")
mgmtraw2 <- ifelse(predraw2>=step27k$perf$cut,"M","U")
prednorm2 <- predict(step27k,dfnorm2,type="response")
mgmtnorm2 <- ifelse(prednorm2>=step27k$perf$cut,"M","U")
predini2 <- predict(step27k,predGBM2,type="response")
mgmtini2 <- ifelse(predini2>=step27k$perf$cut,"M","U")
table(mgmtraw2,mgmtini2)
table(mgmtraw2,mgmtnorm2)
table(mgmtnorm2,mgmtini2)
```


7 Conclusion

to conclude, comments and instructions (as they come) related to the choice of normalization for using the model predicting the DNA methylation of the MGMT promoter, are given below :

- Original data from TCGA (update 2012-05-25) was preprocessed as HM-27K data (e.g. the function `rawpreprocessRaw` from R package `minfi`).
- The normalization can affect the prediction of the DNA methylation (it's not really a surprise).
- The generalization of the model can be affected by the new normalization proposed by Illumina (`preprocessIllumina`). The reference samples used during the normalization procedure were fixed within each dataset and they were not the same among the datasets.
- The model in Bady et al. (2012) requires to use the initial preprocessing (normalization) proposed initially by Illumina in *GenomeStudio* that corresponds to the function `preprocessRaw` from R package `minfi`.
- The predictions proposed in the table S4 for the dataset based on HM-450K of the paper are consistent with our previous comments/recommendations (see above).
- There are no problem/bias induced by chemistry type, because the two probes used in the model come from the chemistry I as in the HM-27K platform.

8 Acknowledgments

The results published here are in part based upon data generated by The Cancer TCGA Genome Atlas pilot project established by the NCI and NHGRI. Information about TCGA and the investigators and institutions who constitute the TCGA research network can be found at (<http://cancergenome.nih.gov>). The dbGaP accession number to the specific version of the TCGA data set is phs000178.v8.p7.

Références

- [1] Martin J Aryee, Andrew E Jaffe, Hector Corrada-Bravo, Christine Ladd-Acosta, Andrew P. Feinberg, Kasper D. Hansen, and Rafael A. Irizarry. `Minfi` : A flexible and comprehensive Bioconductor package for the analysis of Infinium DNA Methylation microarrays. *Bioinformatics*, 2014.
- [2] Pierre Bady, Davide Sciuscio, Annie-Claire Diserens, Jocelyne Bloch, Martin J. van den Bent, Christine Marosi, Pierre-Yves Dietrich, Michael Weller, Luigi Mariani, Frank L. Heppner, David R. McDonald, Denis Lacombe, Roger Stupp, Mauro Delorenzi, and Monika E. Hegi. Mgmt methylation analysis of glioblastoma on the infinium methylation beadchip identifies two distinct cpg regions associated with gene silencing and outcome, yielding a prediction model for comparisons across datasets, tumor grades, and cimp-status. *Acta Neuropathologica*, 124(4) :547–560, 2012. Times Cited : 6.
- [3] Sean Davis, Pan Du, Sven Bilke, Tim Triche, Jr., and Moiz Bootwalla. *methylumi : Handle Illumina methylation data*, 2014. R package version 2.10.0.

- [4] R Core Team. *R : A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2014.

9 Appendix

9.1 Import raw HM-27K data (format .IDAT)

To import raw data (format **.IDAT**), the function contains in R package `minfi` ([1]) don't work with HM-27K. However, it's possible to import data with functions from R package `methylumi` ([3]).

```
require(methylumi)
rgset0<- methylumIDAT(barcode=as.character(File.Name),idatPath=datadir)
# no normalization for HM-27k,
# see help "For HumanMethylation27 data, the function does nothing"
norm27k <- normalizeMethylumiSet(rgset0)
u27k <- unmethyated(norm27k)
m27k <- methyated(norm27k)
mvalue0 <- log2((m27k+1)/(u27k+1))
```

9.2 Session

```
print(sessionInfo(),locale=FALSE)
R version 3.1.1 (2014-07-10)
Platform: x86_64-w64-mingw32/x64 (64-bit)

attached base packages:
[1] stats      graphics  grDevices  utils      datasets  methods   base

loaded via a namespace (and not attached):
[1] tools_3.1.1
```