

Statistical Methods for the Analysis of Microarray and Short-Read Sequencing Data

R Project

Student: Damianos Melidis, Msc in Computational Biology and Bioinformatics
Student id: 12-945-416, dmelidis@student.ethz.ch

December 28, 2013

Abstract

The main goal of this project is to reproduce experimental results using *R* from a chosen publication [Lovisa E. Reinus et al.].

1. Introduction

The authors are motivated to study epigenetic marks for blood samples because of the limitations of extracting large number of (affected and non-affected) tissue samples. In this study they extract 10 different blood cell types from 6 male donors. The cell types are **whole blood**, peripheral blood mononuclear cells (**PBMC**) *CD4 T*, *CD8 T*, *CD56 NK*, *CD19 B*, *CD14 monocytes* and **granulocytes** *Neutrophils*, *Eosinophils*. The Illumina® Human Methylation 450k array used to measure the methylation status for each cell type in the majority of known genes (RefSeq® database). The researchers would like to investigate the difference in methylation between each cell type. Besides trying to examine if only the whole blood methylation status is a sufficient epigenetic mark, the methylation of CpG sites comparing cell types with whole blood cell type. Interestingly authors after extracting the probes related to candidate genes for inflammatory diseases, they assess the similarity of differential methylation of the corresponding CpG sites between all cell types. In this project Figure 2 and 3a as well as table 1 and 2 from [Lovisa E. Reinus et al.] are reproduced using *R*.

3. Reproduced Results

In this section we will show the procedure that was followed to reproduce some results from [Lovisa E. Reinus et al.]. Following the bioinformatics analysis of *Materials and Methods* section the raw data files (.idat) were partitioned to distinct folder for *each* blood cell type, following the given sample_sheet.IDAT.xlsx, creating 10 different folders. Afterwards normalization and background correction was performed. Then the median for probe's methylation signal of the 6 donors is calculated and the principal component analysis and hierarchical clustering was done resulting to Figure 1 and 2, respectively. Comparing these figures with figure 2 of the publication we can see some variation in the PCA (but the median signal and not all individual signals are used (insufficient workspace memory)) and the left clade of reproduced clustering tree is the same as the right clade of Figure 2a. Using *limma* we perform linear fitting for differential methylation between each pair of cells and tested by Bayesian moderated t-test on our normalized methylation signals. This gives us the table 1.

The authors would like to investigate if the whole blood is a sufficient cell type to study epigenetic marks. In this direction they compared the differentially methylated sites between each cell population *and whole blood*. In order to assign the methylation status to *unmethylated*, *marginal* or *methylated* they suggest to fit a gamma model which represents one *normal* distribution for

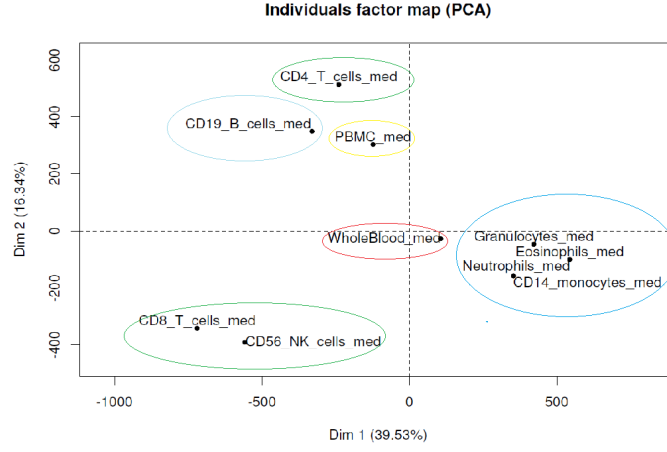


Figure 1: PCA for median signal of the 6 donors

marginal status (centered on 0) one gamma distribution (shifted normal to the positive values) representing *methyated* and one gamma shifted to the negative values *unmethyated*. Afterwards they extract the most significant differentially methylated sites depending on the Benjamini and hochberg corrected p-value. The same analysis is conducted leading to Figure 3 (Figure 3a in [Lovisa E. Reinus et al.]) and table 2. The *R* code used for the whole analysis is given in the next section.

Table 1: Number of differentially methylated CpG sites comparing each pair of cell types.

	Whole blood	PBMC	CD4	CD8	CD56	CD19	CD14	Granulocytes	Neutrophils	Eosinophils
Whole blood	0	126409	128806	115707	119812	124910	128369	130564	128346	128346
PBMC	0	0	130640	117666	121543	126428	128753	130764	128710	128710
CD4	0	0	0	120955	124171	129337	131042	132817	130931	130931
CD8	0	0	0	0	112702	116520	117051	119306	117212	117212
CD56	0	0	0	0	0	120199	121768	123798	121713	121713
CD19	0	0	0	0	0	0	127843	129459	127685	127685
CD14	0	0	0	0	0	0	0	135082	133381	133381
Granulocytes	0	0	0	0	0	0	0	0	135082	135082
Neutrophils	0	0	0	0	0	0	0	0	0	133800
Eosinophils	0	0	0	0	0	0	0	0	0	0

Table 2: Differential methylated probes for each cell type in comparison to whole blood

	M-value	Comparison on calls	Unmethyated	Marginal	Methyated
PBMC	126409	30000	0.76	0.011	0.23
CD4	128806	30000	0.74	0.02	0.26
CD8	115707	26706	0.73	0.02	0.25
CD56	119812	28745	0.74	0.02	0.24
CD19	124910	28630	0.72	0.02	0.26
CD14	128369	29774	0.72	0.04	0.24
Granulocytes	130564	30000	0.74	0.03	0.24
Neutrophils	128346	30000	0.74	0.03	0.24
Eosinophils	128346	30000	0.74	0.03	0.23

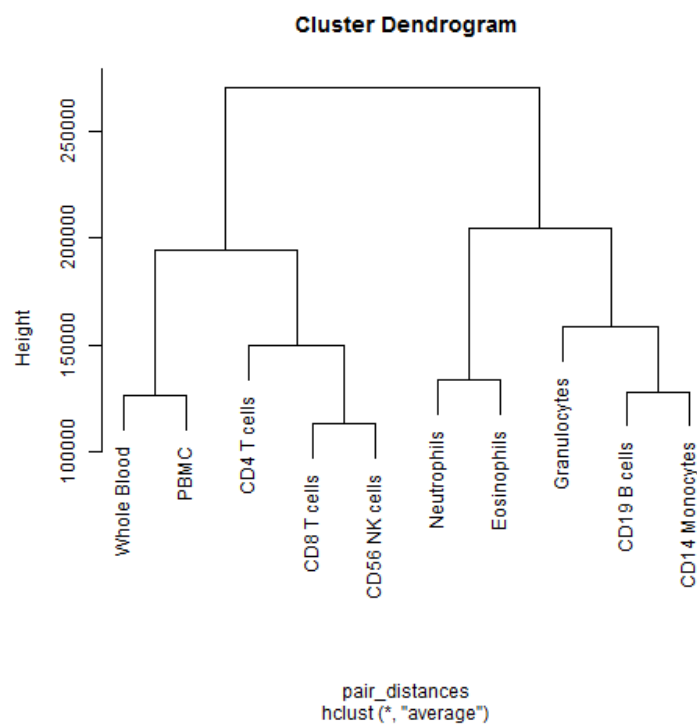


Figure 2: Hierarchical clustering for median signal of the 6 donors

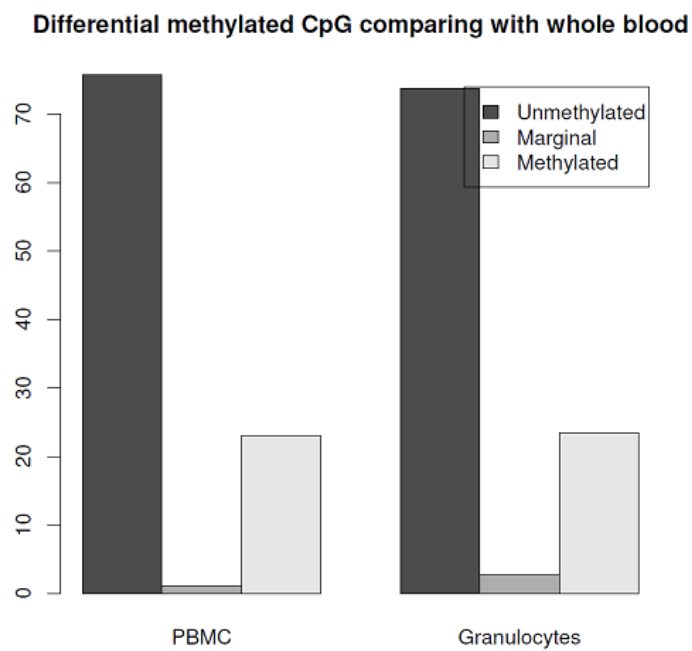


Figure 3: Distribution of methylation status for PBMC and Granulocytes comparing to whole blood.

4. Appendix

In the following the actual code is given:

```
1 setwd("C:/Users/damian/Desktop/R_project/idadat_files")
2 library("limma")
3 library("lumi")
4 library("minfi")
5 library("FactoMineR")
6 IlluminaHumanMethylation450k("IlluminaHumanMethylation450k.db")
7 #read idat files for each cell type
8 WholeBlood_Int <- read.450k.exp("WholeBlood")
9 PBMC_Int <- read.450k.exp("PBMC")
10 CD19_Int <- read.450k.exp("CD19")
11 CD4_Int <- read.450k.exp("CD4")
12 CD56_Int <- read.450k.exp("CD56")
13 CD8_Int <- read.450k.exp("CD8")
14 CD14_Int <- read.450k.exp("CD14")
15 Eosinophils_Int <- read.450k.exp("Eosinophils")
16 Granulocytes_Int <- read.450k.exp("Granulocytes")
17 Neutrophils_Int <- read.450k.exp("Neutrophils")
18
19 #normalization of data with background correction ##needs IlluminaHumanMethylation450kmanifest
20 WholeBlood_Meth.norm <- preprocessIllumina(WholeBlood_Int, bg.correct = TRUE, normalize = "controls",
21 reference = 2)
22 PBMC_Meth.norm <- preprocessIllumina(PBMC_Int, bg.correct = TRUE, normalize = "controls", reference = 2)
23 CD19_Meth.norm <- preprocessIllumina(CD19_Int, bg.correct = TRUE, normalize = "controls", reference = 2)
24 CD4_Meth.norm <- preprocessIllumina(CD4_Int, bg.correct = TRUE, normalize = "controls", reference = 2)
25 CD56_Meth.norm <- preprocessIllumina(CD56_Int, bg.correct = TRUE, normalize = "controls", reference = 2)
26 CD8_Meth.norm <- preprocessIllumina(CD8_Int, bg.correct = TRUE, normalize = "controls", reference = 2)
27 CD14_Meth.norm <- preprocessIllumina(CD14_Int, bg.correct = TRUE, normalize = "controls", reference = 2)
28 Eosinophils_Meth.norm <- preprocessIllumina(Eosinophils_Int, bg.correct = TRUE, normalize = "controls",
29 reference = 2)
30 Granulocytes_Meth.norm <- preprocessIllumina(Granulocytes_Int, bg.correct = TRUE, normalize = "controls",
31 reference = 2)
32 Neutrophils_Meth.norm <- preprocessIllumina(Neutrophils_Int, bg.correct = TRUE, normalize = "controls",
33 reference = 2)
34
35 #get the M values
36 WholeBlood_Mvalue <- getM(WholeBlood_Meth.norm)
37 PBMC_Mvalue <- getM(PBMC_Meth.norm)
38 CD19_Mvalue <- getM(CD19_Meth.norm)
39 CD4_Mvalue <- getM(CD4_Meth.norm)
40 CD56_Mvalue <- getM(CD56_Meth.norm)
41 CD8_Mvalue <- getM(CD8_Meth.norm)
42 CD14_Mvalue <- getM(CD14_Meth.norm)
43 Eosinophils_Mvalue <- getM(Eosinophils_Meth.norm)
44 Granulocytes_Mvalue <- getM(Granulocytes_Meth.norm)
45 Neutrophils_Mvalue <- getM(Neutrophils_Meth.norm)
46
47 #create a matrix with rows the probes and columns the all cell types for each donor
48 mat_Mvalue <- rbind(WholeBlood_Mvalue, PBMC_Mvalue, CD19_Mvalue, CD4_Mvalue, CD56_Mvalue, CD8_Mvalue, CD14_
49 Mvalue, Eosinophils_Mvalue, Granulocytes_Mvalue, Neutrophils_Mvalue)
50
51 #give name to columns
52 WB_donors <- paste("WB_Donor_", 1:6)
53 PBMC_donors <- paste("PBMC_Donor_", 1:6)
54 CD19_donors <- paste("CD19_Donor_", 1:6)
55 CD4_donors <- paste("CD4_Donor_", 1:6)
56 CD56_donors <- paste("CD56_Donor_", 1:6)
57 CD8_donors <- paste("CD8_Donor_", 1:6)
58 CD14_donors <- paste("CD14_Donor_", 1:6)
59 Eosinophils_donors <- paste("Eosinophils_Donor_", 1:6)
60 Granulocytes_donors <- paste("Granulocytes_Donor_", 1:6)
61 Neutrophils_donors <- paste("Neutrophils_Donor_", 1:6)
62 colnames(mat_Mvalue) <- c(WB_donors, PBMC_donors, CD19_donors, CD4_donors, CD56_donors, CD8_donors, CD14_
63 donors, Eosinophils_donors, Granulocytes_donors, Neutrophils_donors)
64
65 #get the median values as representative for the six individuals
66 GenomicRegions <- rownames(WholeBlood_Mvalue)
67 WholeBlood_med <- apply(WholeBlood_Mvalue, 1, median, na.rm = T)
68 PBMC_med <- apply(PBMC_Mvalue, 1, median, na.rm = T)
69 CD19_B_cells_med <- apply(CD19_Mvalue, 1, median, na.rm = T)
70 CD4_T_cells_med <- apply(CD4_Mvalue, 1, median, na.rm = T)
71 CD56_NK_cells_med <- apply(CD56_Mvalue, 1, median, na.rm = T)
72 CD8_T_cells_med <- apply(CD8_Mvalue, 1, median, na.rm = T)
73 CD14_monocytes_med <- apply(CD14_Mvalue, 1, median, na.rm = T)
74 Eosinophils_med <- apply(Eosinophils_Mvalue, 1, median, na.rm = T)
75 Granulocytes_med <- apply(Granulocytes_Mvalue, 1, median, na.rm = T)
76 Neutrophils_med <- apply(Neutrophils_Mvalue, 1, median, na.rm = T)
77
78 #####Analysis for Figure 2 and Table 1#####
79 #create a matrix n * m, where (n) rows are the number of different cells and (m) columns are the genomic
80 regions that measured for methylation
81 mat_Mvalue_med <- rbind(WholeBlood_med, PBMC_med, CD4_T_cells_med, CD8_T_cells_med, CD56_NK_cells_med,
82 CD19_B_cells_med, CD14_monocytes_med, Granulocytes_med, Neutrophils_med, Eosinophils_med)
83 #replace the -Inf (totally unmethylated) and Inf (totally methylated) with 1.0 and 0.0 respectively
84 mat_Mvalue_med[ mat_Mvalue_med == Inf ] <- max(mat_Mvalue_med[mat_Mvalue_med < Inf])
85 mat_Mvalue_med[ mat_Mvalue_med == -Inf ] <- min(mat_Mvalue_med[mat_Mvalue_med > -Inf])
86
87 #-----Figure 2.B-----#
88 #do PCA to find the most important genomic regions for this dataset (e.g. eigenvectors of previous matrix)
89 PCA(mat_Mvalue_med)
90
91 #-----Figure 2.B-----#
92 #find differential methylation between each pair of cells
93 mat_Mvalue_med <- t(mat_Mvalue_med) #transpose matrix in order to get the cell population in the rows and
94 the CpG sites in the columns
95 diff.exp.probes <- matrix(0,10,10) # create a matrix to save the differential expression between each
96 pair of cells (e.g diagonal matrix)
```

```

85 for (i in 1:9){#loop through cell and compute the differential methylated CpG sites with all other cells
86   c <- i+1
87   for(j in c:10){
88     ExtractedMatrix <- cbind(mat_Mvalue_med[,i],mat_Mvalue_med[,j]) #get the CpG sites for the two
89     #show(ExtractedMatrix)
90     fit <- lmFit(ExtractedMatrix)
91     fit <- eBayes(fit)
92     diff_exp_probes[i,j] = nrow(topTable(fit , coef=1, number=300000, lfc = 4))
93   }
94 }
95 cell_names <- c("Whole Blood", "PBMC", "CD4 T cells", "CD8 T cells", "CD56 NK cells", "CD19 B cells", "
96   CD14 Monocytes", "Granulocytes", "Neutrophils", "Eosinophils")
97 rownames(diff_exp_probes) <- cell_names
98 colnames(diff_exp_probes) <- cell_names
99 #-----Table 1-----#
100 write.table(diff_exp_probes,"clipboard",sep="\t", col.names = T, row.names = T)
101 #-----Table 1-----#
102 #-----Figure 2.A-----#
103 #perform hierarchical clustering for cell populations (using euclidean distance)
104 pair_distances <- dist(diff_exp_probes, method = "euclidean")
105 png("CellPopulation_Clustering_med22.png")
106 plot(hclust(pair_distances, method="average"))
107 dev.off()
108 #-----Figure 2.A-----#
109 #####Analysis for Figure 2 and Table 1#####
110 #####Analysis for Figure 3 and Table 2#####
111 table2 <- matrix(0,9,5)
112 table2[,1] = t(diff_exp_probes[1,2:10]) #get the first column showing differentially methylated sites
113   between cell types and whole blood
114 for (i in 2:10){
115   ExtractedMatrix <- cbind(mat_Mvalue_med[,1],mat_Mvalue_med[,i]) #compare each cell type with whole blood
116   fit <- lmFit(ExtractedMatrix)
117   fit <- eBayes(fit)
118   significant_meth_sites <- topTable(fit , coef=NULL, number=30000, adjust.method="BH", sort.by="B", resort.by=
119     NULL, p.value=0.01, lfc=6)
120   table2[i-1,2] = nrow(significant_meth_sites) #get the number of significant methylation sites
121   methylation_selected_sites <- mat_Mvalue_med[rownames(significant_meth_sites),2] #get the significantly
122     methylated sites
123   fittedGamma <- gammaFitEM(methylation_selected_sites, initialFit=NULL, maxIteration = 150, tol = 1e-04,
124     plotMode = FALSE, verbose = FALSE)
125   #plotGammaFit(methylation_selected_sites, gammaFit=fittedGamma)
126   status <- methylationCall(fittedGamma)
127   methyl_status <- table(status)
128   margin_proportion <- methyl_status[1] / table2[i-1,2]
129   methyl_proportion <- methyl_status[2] / table2[i-1,2]
130   unmethyl_proportion <- methyl_status[3] / table2[i-1,2]
131   table2[i-1,3] <- unmethyl_proportion
132   table2[i-1,4] <- margin_proportion
133   table2[i-1,5] <- methyl_proportion
134 }
135 table2
136 rownames(table2) <- cell_names(2:10)
137 colnames(table2) <- c("M-value comparison", "Comparison on calls", "Unmethylated", "Marginal", "Methylated
138   ")
139 #-----Table 2-----#
140 write.table(table2,"clipboard",sep="\t", col.names = T, row.names = T)
141 #-----Table 2-----#
142 #-----Figure 3.A-----#
143 PBMC_Granulocytes_Meth_status <- 100 * rbind(table2[1,3:5], table2[7,3:5])
144 rownames(PBMC_Granulocytes_Meth_status) <- c("PBMC", "Granulocytes")
145 barplot(PBMC_Granulocytes_Meth_status, legend=T, beside=T, main='Differential methylated CpG comparing with
146   whole blood')
147 #-----Figure 3.A-----#
148 #####Analysis for Figure 3 and Table 2#####
149 #####Incomplete part#####
150 x <- IlluminaHumanMethylation450kENTREZID
151 # Get the probe identifiers that are mapped to an ENTREZ Gene ID
152 mapped_probes <- mappedkeys(x)
153 # Convert to a list
154 xx <- as.list(x[mapped_probes])

```

References

[Lovisa E. Reinius et al.] Lovisa E. Reinius et al. Differential DNA methylation in Purified Human Blood Cells: Implications for Cell Lineage and Studies on Disease Susceptibility Plos One, vol. 7 no. 7 (2012).