# DRD report

*Roberto Polverelli Monti*
*LM Bioinformatics*
*A.Y. 2018/2019*
*Bologna*

## Step 1-4

### Preliminary work

- Start with a clean environment, and load the *minfi* package:

```
1  rm(list=ls())
2  setwd("~/reportDRD")
3  library("minfi")
```

- Import targets, then load the manifest:

```
1  targets <- read.metharray.sheet("inputData")
2  load("Illumina450Manifest.RData")
```

- Subset the manifest as to remove those probes that do not map to any chromosome:

```
1  Illumina450Manifest <-
   droplevels(Illumina450Manifest[Illumina450Manifest$CHR!="",])
```

- Extract probe intensities:

```
1  RGset <- read.metharray.exp(targets=targets)
2  save(RGset,file = "RGset.Rdata")
```

- Store the red and green fluorescence intensities:

```
1  red <- data.frame(getRed(RGset))
2  green <- data.frame(getGreen(RGset))
```

- For convenience's sake, store the address and obtain the probe type:

```
1  my_address <- 52682510
2  probe_type <-
   as.character(Illumina450Manifest[Illumina450Manifest$AddressA_ID==my
   _address,]$Infinium_Design_Type)
```

- Separately obtain the red and green intensities for the address:

```
1  my_address_green <- green[row.names(green)==my_address,]
2  my_address_red <- red[row.names(red)==my_address,]
3  my_intensities <-
   cbind(t(my_address_red),t(my_address_green),probe_type)
4  my_intensities <- cbind(row.names(my_intensities),my_intensities)
5  rownames(my_intensities) <- NULL
6  colnames(my_intensities) <- c("Sample","Red","Green","Type")
7  View(my_intensities)
```

| | Sample | Red | Green | Type |
|---|---|---|---|---|
| 1 | X5775278008_R01C01 | 941 | 3573 | II |
| 2 | X5775278008_R02C01 | 1259 | 3820 | II |
| 3 | X5775278008_R04C01 | 875 | 866 | II |
| 4 | X5775278008_R05C01 | 745 | 716 | II |
| 5 | X5775278035_R01C01 | 675 | 1652 | II |
| 6 | X5775278035_R02C01 | 462 | 393 | II |
| 7 | X5775278035_R04C01 | 1846 | 4182 | II |
| 8 | X5775278035_R05C01 | 1331 | 923 | II |

**Table 1:** the red and green intensities for the given address.

- Convert red/green channel intensities to methylation signal, and save it to a file:

```
1  MSet.raw <- preprocessRaw(RGset)
2  save(MSet.raw,file="MSet.raw.RData")
```
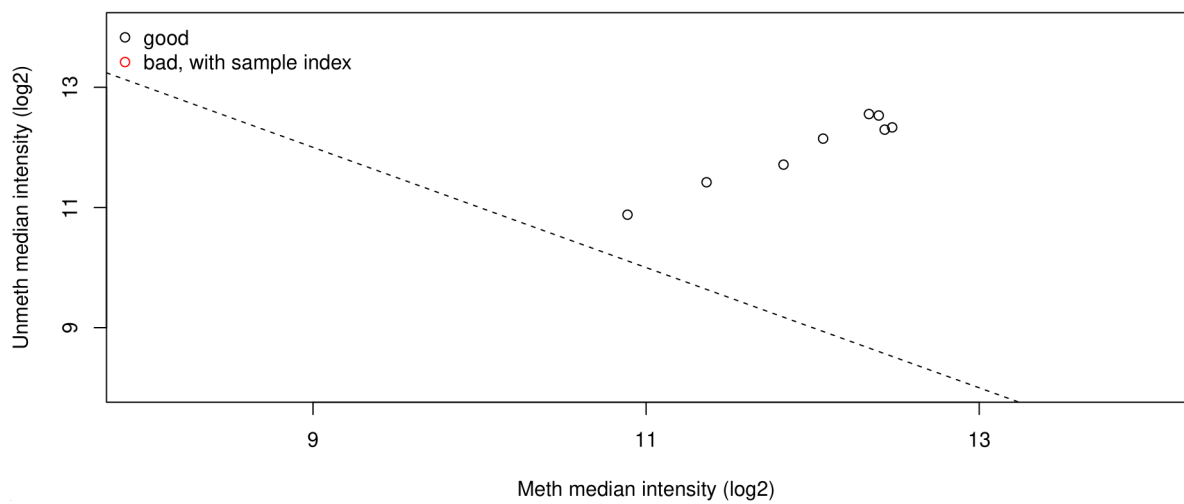
# Step 5

## Quality checks

- Obtain and plot the chipwide medians of the Meth and Unmeth channels:

```
1  pdf(file="step5.1.pdf")
2  qc <- getQC(MSet.raw)
3  plotQC(qc)
4  dev.off()
```
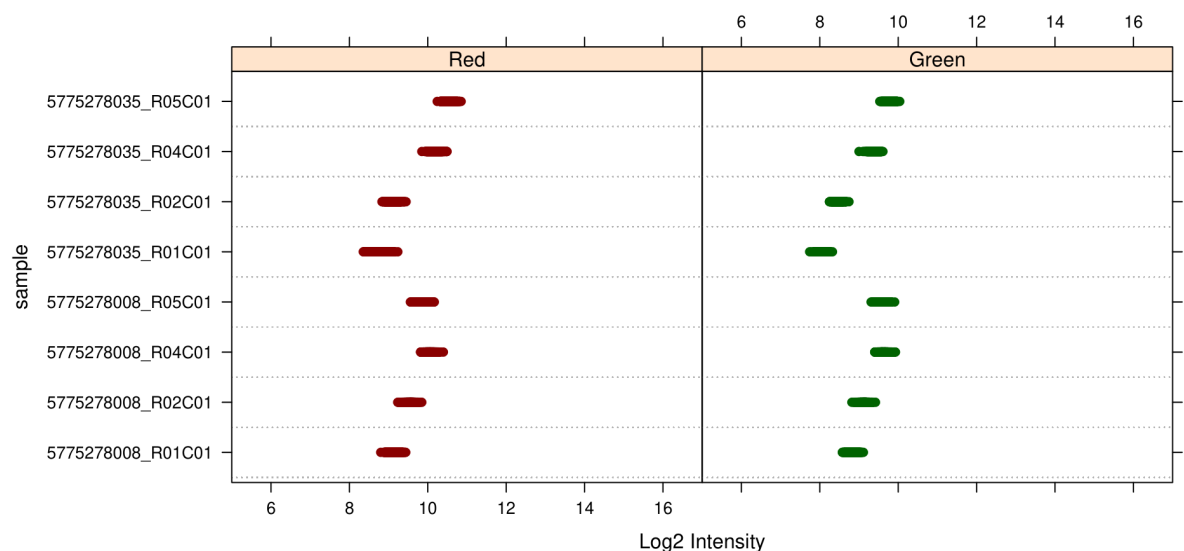
**Figure 1:** no sample reporting bad detection intensity.

- Produce the negative control probe signal plot component of the QC report:

```
1  pdf(file="step5.2.pdf")
2  controlStripPlot(RGset,controls="NEGATIVE")
3  dev.off()
```



**Figure 2:** for all eight samples, Log2 intensities fall within detection range.

- Check the number of bad probes for each sample, with detection P-value 0.01 as threshold:

```
1  det_p <- detectionP(RGset)
2  failed_positions <- det_p > 0.01
3  failed_count <- c()
4  for (i in colnames(failed_positions)){
5    failed_count.newrow <- t(as.integer(summary(failed_positions[,i])[-1]))
6    failed_count.newrow <- cbind(i,failed_count.newrow)
7    failed_count <- rbind(failed_count,failed_count.newrow)
8  }
9  rm(i,failed_count.newrow)
10  colnames(failed_count) <- c("Sample","Passed","Failed")
11  View(failed_count)
```

| | Sample | Passed | Failed |
|---|---|---|---|
| 1 | 5775278008_R01C01 | 484946 | 566 |
| 2 | 5775278008_R02C01 | 485061 | 451 |
| 3 | 5775278008_R04C01 | 484665 | 847 |
| 4 | 5775278008_R05C01 | 484631 | 881 |
| 5 | 5775278035_R01C01 | 485038 | 474 |
| 6 | 5775278035_R02C01 | 484845 | 667 |
| 7 | 5775278035_R04C01 | 485182 | 330 |
| 8 | 5775278035_R05C01 | 484481 | 1031 |

**Table 2:** number of positions reporting a detection P-value lower or higher than the threshold.

- Store in a vector all those probes with a reported detection P-value higher than 0.01 in more than 1% of the samples, and check whether this number (i.e. 2191) is in agreement with the number of good probes and the total number of probes:

```
failed_proportion <- rowMeans(failed_positions)
bad_probes <- names(failed_proportion[failed_proportion > 0.01])
good_probes <- names(failed_proportion[failed_proportion <= 0.01])

length(good_probes) + length(bad_probes) ==
length(failed_proportion)
# [1] TRUE
length(bad_probes)
# [1] 2191
```

# Step 6

## $\beta$-values and M-values

- Obtain the raw $\beta$-values and M-values from the MethylSet object:

```
beta <- getBeta(MSet.raw)
M <- getM(MSet.raw)
```

- Subset the data into *WT* and *DS* subjects:

```
1  wt_names <- targets[targets$Group=="WT","Basename"]
2  wt_names <- sub(".+/","",wt_names)
3  wt_beta <- beta[,wt_names]
4  wt_M <- M[,wt_names]
5
6  ds_names <- targets[targets$Group=="DS","Basename"]
7  ds_names <- sub(".+/","",ds_names)
8  ds_beta <- beta[,ds_names]
9  ds_M <- M[,ds_names]
```
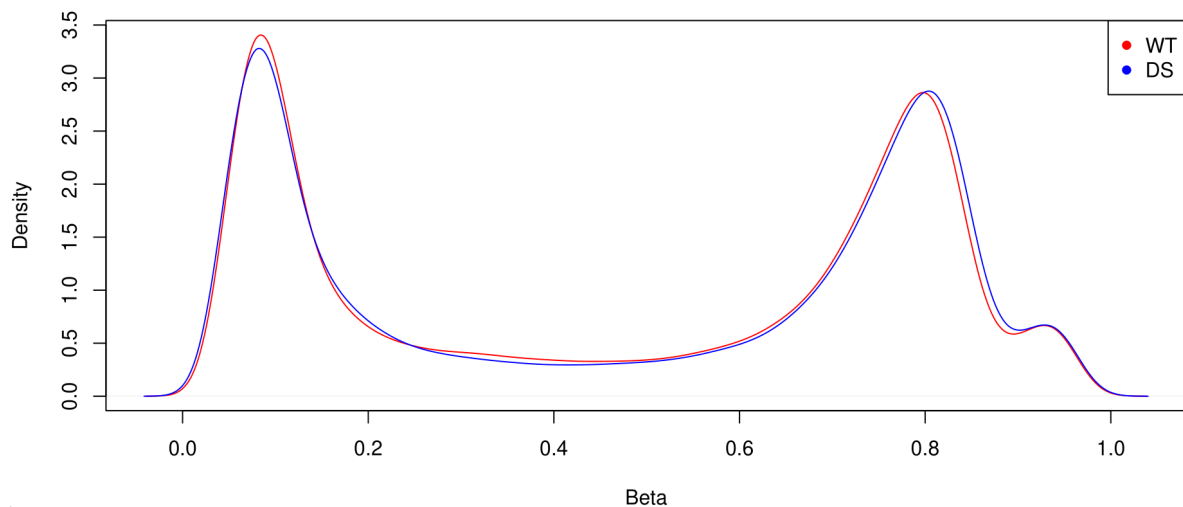
- Obtain the mean methylation values for each subset, and plot their densities:

```
1   mean_wt_beta <- apply(wt_beta,1,mean)
2   mean_ds_beta <- apply(ds_beta,1,mean)
3
4   mean_wt_M <- apply(wt_M,1,mean)
5   mean_ds_M <- apply(ds_M,1,mean)
6
7   pdf(file="step6.1.pdf")
8   plot(density(mean_wt_beta,na.rm=TRUE),col=4,main="Mean Beta-values,
       density",xlab="Beta")
9   lines(density(mean_ds_beta,na.rm=TRUE),col=3)
10  legend('topright',pch=c(16,16),col=c(4,3),legend=c("WT","DS"))
11  dev.off()
```
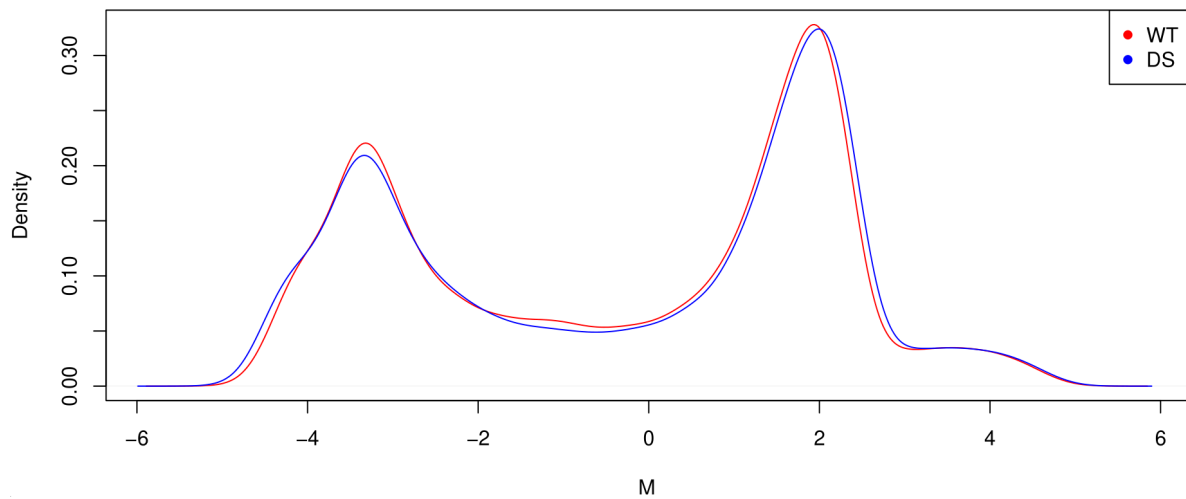


**Figure 3:** density of mean $\beta$-values.

```
1  pdf(file="step6.2.pdf")
2  plot(density(mean_wt_M,na.rm=TRUE),col=4,main="Mean M-values,
      density",xlab="M")
3  lines(density(mean_ds_M,na.rm=TRUE),col=3)
4  legend('topright',pch=c(16,16),col=c(4,3),legend=c("WT","DS"))
5  dev.off()
```

**Figure 4:** density of mean M-values.

# Step 7

## Subsetting

- Subset the manifest according to probe chemistry:

```
1  type_I_manifest <-
   Illumina450Manifest[Illumina450Manifest$Infinium_Design_Type=="I",]
2  type_II_manifest <-
   Illumina450Manifest[Illumina450Manifest$Infinium_Design_Type=="II",]
```

- Subset beta according to type of probe, and compute both the mean and S.D.:

```
1  beta_I <- beta[rownames(beta) %in% type_I_manifest$IlmnID,]
2  mean_beta_I <- apply(beta_I,1,mean)
3  sd_beta_I <- apply(beta_I,1,sd)
4
5  beta_II <- beta[rownames(beta) %in% type_II_manifest$IlmnID,]
6  mean_beta_II <- apply(beta_II,1,mean)
7  sd_beta_II <- apply(beta_II,1,sd)
```

## Normalization

- Normalize fluorescence intensities, subset them according to type of probe, and compute the mean and the S.D.:

```
1   norm_RGset <- preprocessFunnorm(RGset)
2   norm_beta <- getBeta(norm_RGset)
3
4   norm_beta_I <- norm_beta[rownames(norm_beta) %in%
    type_I_manifest$IlmnID,]
5   mean_norm_beta_I <- apply(norm_beta_I,1,mean)
6   sd_norm_beta_I <- apply(norm_beta_I,1,sd)
7
8   norm_beta_II <- norm_beta[rownames(norm_beta) %in%
    type_II_manifest$IlmnID,]
9   mean_norm_beta_II <- apply(norm_beta_II,1,mean)
10  sd_norm_beta_II <- apply(norm_beta_II,1,sd)
```
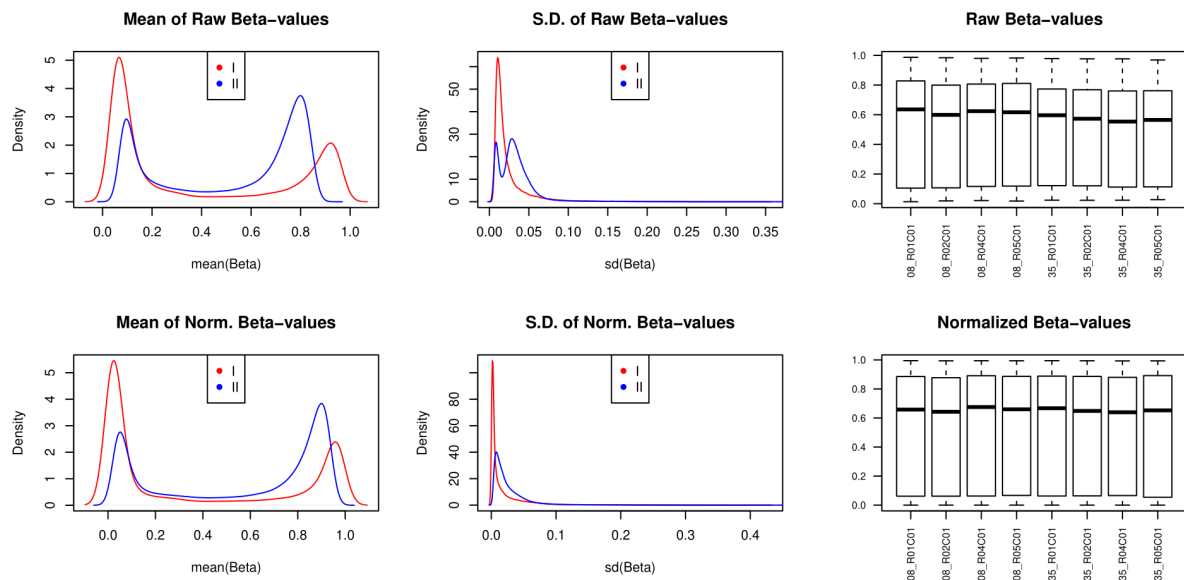
## Comparison

- Produce a six panel plot comparing raw and normalized data:

```
1   pdf(file="step7.pdf",width=10,height=5)
2   par(mfrow=c(2,3))
3
4   plot(density(mean_beta_I,na.rm=T),main="Mean of Raw Beta-
    values",xlab="mean(Beta)",col=4)
5   lines(density(mean_beta_II,na.rm=T),col=3)
6   legend('top',pch=c(16,16),col=c(4,3),legend=c("I","II"))
7
8   plot(density(sd_beta_I,na.rm=T),main="S.D. of Raw Beta-
    values",xlab="sd(Beta)",col=4)
9   lines(density(sd_beta_II,na.rm=T),col=3)
10  legend('top',pch=c(16,16),col=c(4,3),legend=c("I","II"))
11
12  colnames(beta) <- sub("57752780","",colnames(beta))
13  boxplot(beta,main="Raw Beta-values",las=2,cex.axis=0.7)
14  colnames(beta) <- sub("","57752780",colnames(beta))
15
16  plot(density(mean_norm_beta_I,na.rm=T),main="Mean of Norm. Beta-
    values",xlab="mean(Beta)",col=4)
17  lines(density(mean_norm_beta_II,na.rm=T),col=3)
18  legend('top',pch=c(16,16),col=c(4,3),legend=c("I","II"))
19
20  plot(density(sd_norm_beta_I,na.rm=T),main="S.D. of Norm. Beta-
    values",xlab="sd(Beta)",col=4)
21  lines(density(sd_norm_beta_II,na.rm=T),col=3)
22  legend('top',pch=c(16,16),col=c(4,3),legend=c("I","II"))
23
24  colnames(norm_beta) <- sub("57752780","",colnames(norm_beta))
25  boxplot(norm_beta,main="Normalized Beta-values",las=2,cex.axis=0.7)
26  colnames(norm_beta) <- sub("","57752780",colnames(norm_beta))
27  dev.off()
```

**Figure 5:** after normalization, shifing of strongly methylated sites allows for type I and II probes distribution to better coincide; type I probes are much more represented in weakly methylated sited, as would be expected giving that they preferentially target CpG islands. Normalized data also shows shorter variance which makes type I and II more easily comparable. Experimental variation across samples is also reduced in normalized data.

# Step 8

## Filtering

- Extract the M-values from the normalized RGset, then filter the normalized data retaining only the good probes:

```
1  norm_M <- getM(norm_RGset)
2
3  filt_norm_beta <- norm_beta[!row.names(norm_beta) %in% bad_probes,]
4  filt_norm_M <- norm_M[!row.names(norm_M) %in% bad_probes,]
```

- Checking sizes just to be sure:

```
1  nrow(norm_beta) - nrow(filt_norm_beta) == length(bad_probes)
2  # [1] TRUE
3  nrow(norm_M) - nrow(filt_norm_M) == length(bad_probes)
4  # [1] TRUE
5  nrow(filt_norm_beta) == length(good_probes)
6  # [1] TRUE
7  nrow(filt_norm_M) == length(good_probes)
8  # [1] TRUE
```

# Step 9

## Homo/heteroscedasticity

- Compute mean and S.D. for the now filtered data:

```
1  mean_filt_norm_beta <- apply(filt_norm_beta,1,mean)
2  sd_filt_norm_beta <- apply(filt_norm_beta,1,sd)
3
4  mean_filt_norm_M <- apply(filt_norm_M,1,mean)
5  sd_filt_norm_M <- apply(filt_norm_M,1,sd)
```
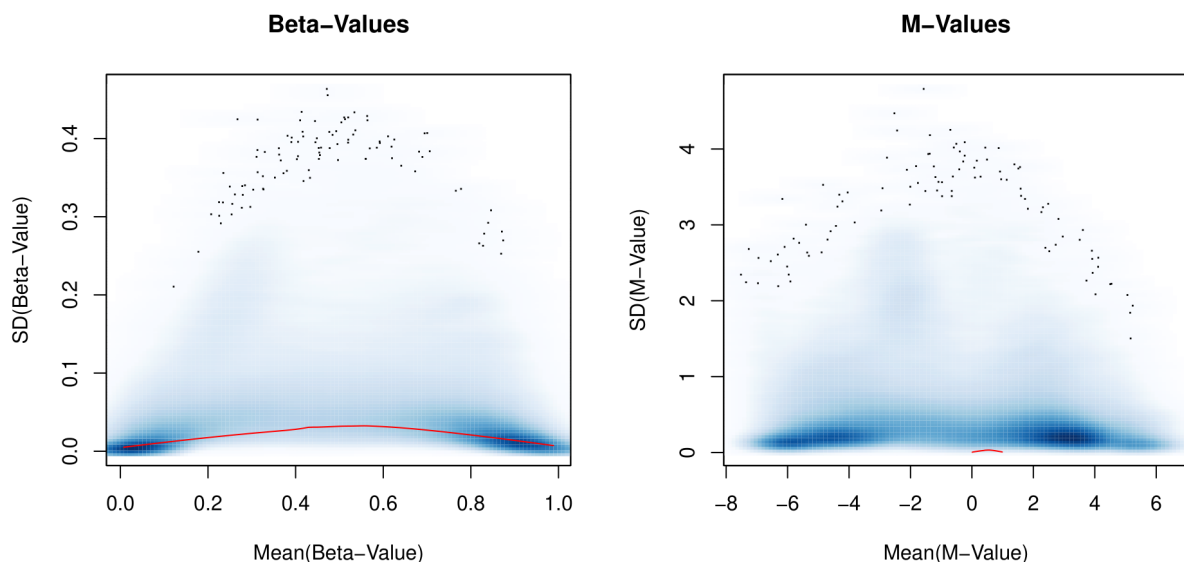
- We can assess the homo/heteroscedasticity of the data by plotting the mean against the S.D.:

```
1  pdf(file="step9.pdf",width=10,height=5)
2  par(mfrow=c(1,2))
3  smoothScatter(mean_filt_norm_beta,sd_filt_norm_beta,main="Beta-
   Values",xlab="Mean(Beta-Value)",ylab="SD(Beta-Value)")
4  lines(lowess(mean_filt_norm_beta,sd_filt_norm_beta),col="red")
5  smoothScatter(mean_filt_norm_M,sd_filt_norm_M,main="M-
   Values",xlab="Mean(M-Value)",ylab="SD(M-Value)")
6  lines(lowess(mean_filt_norm_beta,sd_filt_norm_beta),col="red")
7  dev.off()
```



**Figure 6:** $\beta$-values appear to be more heavily heteroscedastic, showing lower variance within the two groups coinciding with near-zero and near-one methylation signal; M-values are still heteroscedastic in nature, only less prominently.

# Step 10

## Principal Component Analysis

- Perform PCA on the normalized $\beta$-values:

```
1  pca <- prcomp(t(filt_norm_beta),scale=TRUE)
2  summary(pca)
3
4  # Importance of components:
5  #                PC1       PC2       PC3       PC4       PC5       PC6
      PC7         PC8
6  # (1)     322.4629 287.7958 276.2468 256.3168 238.4123 224.9492
   216.93110 4.614e-12
7  # (2)       0.2151   0.1714   0.1579   0.1359   0.1176   0.1047
   0.09737 0.000e+00
8  # (3)       0.2151   0.3865   0.5444   0.6803   0.7979   0.9026
   1.00000 1.000e+00
9  #
10 ## (1):  Standard deviation
11 ## (2):  Proportion of Variance
12 ## (3):  Cumulative Proportion
```
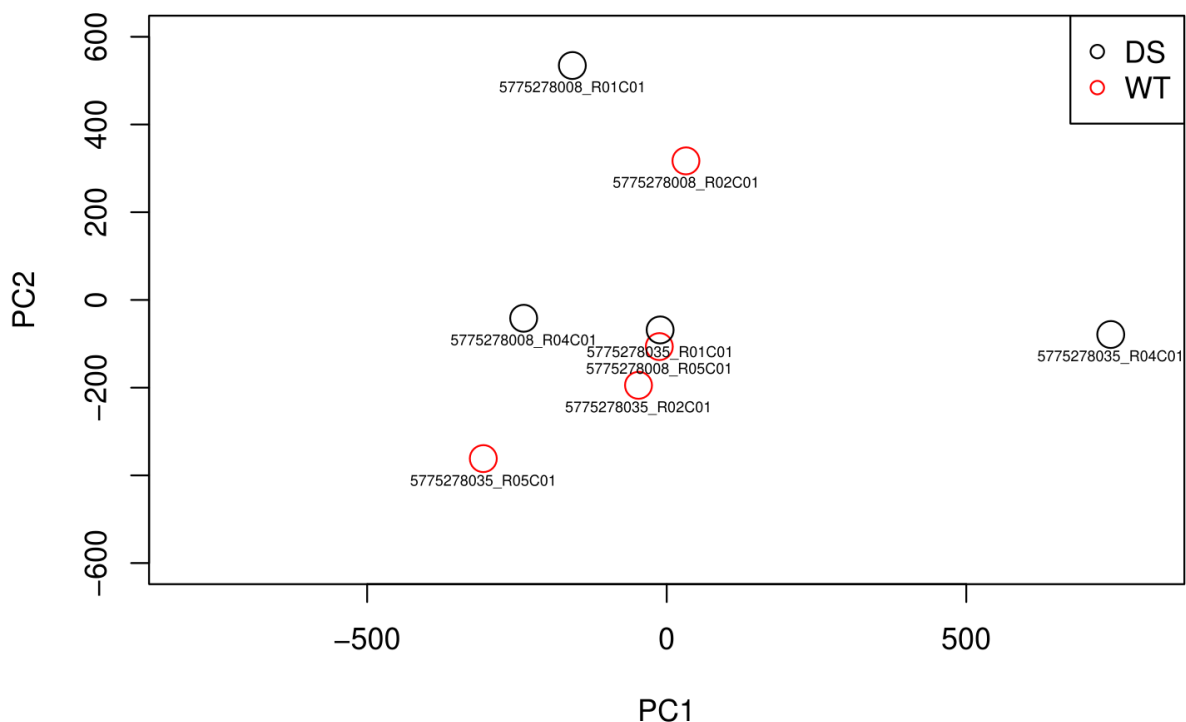
- Plot PC1 and PC2:

```
1  pdf(file="step10.pdf",width=7,height=5)
2  targets$Group <- as.factor(targets$Group)
3  plot(pca$x[,1],pca$x[,2],main="PCA of Norm. Beta
   Values",xlab="PC1",ylab="PC2",xlim=c(-800,800),ylim=c(-600,600),cex=
   2,col=targets$Group)
4  text(pca$x[,1],pca$x[,2],labels=rownames(pca$x),cex=0.5,pos=1)
5  legend("topright",legend=levels(targets$Group),col=c(1:nlevels(targe
   ts$Group)),pch=1)
6  dev.off()
```

## PCA of Norm. Beta Values



**Figure 7:** *5775278008_R02C01* appears to be an outlier in the PCA plot.

# Step 11

## ANOVA

- Perform ANOVA in order to identify differentially methylated probes between groups *WT* and *DS*, using sex as a covariate:

```
1  my_anova_fun <- function(x) {
2    aov_test <- aov(x~ targets$Group+as.factor(targets$Female))
3    return(summary(aov_test)[[1]][[5]][1])
4  }
5
6  p_values_anova <- apply(filt_norm_beta,1,my_anova_fun)
```

# Step 12

## Correction

- Count the number of differentially methylated probes considering a significance threshold of 0.05:

```
1  length(p_values_anova[p_values_anova<=0.05])
2
3  # [1] 30493
```

- Same thing but after Bonferroni correction:

```
1  bonf_p_values_anova <- p.adjust(p_values_anova,"bonferroni")
2  length(bonf_p_values_anova[bonf_p_values_anova<=0.05])
3
4  # [1] 1
```

- After Benjamini-Hochberg:

```
1  benj_p_values_anova <- p.adjust(p_values_anova,"BH")
2  length(benj_p_values_anova[benj_p_values_anova<=0.05])
3
4  # [1] 10
```

- A table showing the probes just before mentioned:

```
1  final_anova <-
   data.frame(filt_norm_beta,p_values_anova,bonf_p_values_anova,benj_p_
   values_anova)
2  final_anova <- final_anova[order(benj_p_values_anova),]
3  final_anova_benj_0.05 <-
   final_anova[final_anova$benj_p_values_anova<=0.05,]
4  View(final_anova_benj_0.05[,9:11])
```

|            | P-values     | Bonferroni  | BH          |
| ---------- | ------------ | ----------- | ----------- |
| cg19246080 | 3.841172e-09 | 0.001856519 | 0.001856519 |
| cg03639185 | 1.606309e-07 | 0.077636264 | 0.031566401 |
| cg23260026 | 1.959344e-07 | 0.094699204 | 0.031566401 |
| cg07867687 | 3.359398e-07 | 0.162366767 | 0.040591692 |
| cg16044109 | 6.457099e-07 | 0.312085155 | 0.044583594 |
| cg25542438 | 5.535240e-07 | 0.267529754 | 0.044583594 |
| cg01086462 | 4.848629e-07 | 0.234344410 | 0.044583594 |
| cg08967584 | 7.640929e-07 | 0.369302139 | 0.046162767 |
| cg06621691 | 9.697929e-07 | 0.468721281 | 0.048031239 |
| cg20348344 | 9.937751e-07 | 0.480312394 | 0.048031239 |

**Table 3:** showing uncorrected and corrected P-values for the ten probes that have a reported FDR adjusted P-value lower than 0.05 with Benjamini-Hochberg.

# Step 13

## Manhattan plot

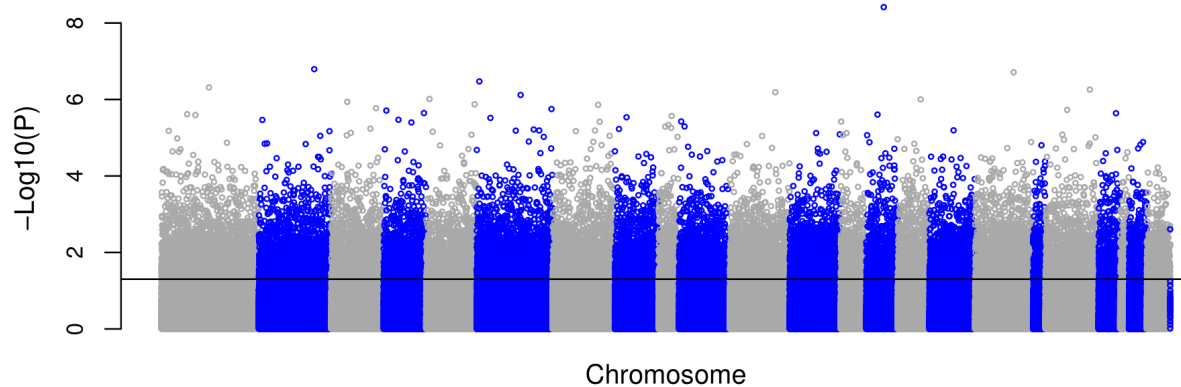- Produce a Manhattan plot of the uncorrected P-values:

```
1  library("gap")
2  final_anova_with_IlmnID <-
   data.frame("IlmnID"=rownames(final_anova),final_anova)
3  final_anova_annotated <-
   merge(final_anova_with_IlmnID,Illumina450Manifest,by="IlmnID")
4
5  chrs <- c(as.character(1:22),"X","Y")
6  db <- data.frame(final_anova_annotated$CHR,
   final_anova_annotated$MAPINFO,
   final_anova_annotated$p_values_anova)
7  db$final_anova_annotated.CHR <-
   factor(db$final_anova_annotated.CHR,levels=chrs)
8
9  pdf(file="step13.1.pdf",width=8,height=4)
10 mhtplot(db,xlab="",ylab="",control=mht.control(colors=rep(c("darkgr
   ey","blue"),12)))
```

```
11   axis(2,cex.axis=0.8)
12   abline(a=-log10(0.05),b=0)
13   title(xlab="Chromosome",line=0.5)
14   title(ylab="-Log10(P)",line=2.5)
15   dev.off()
```



**Figure 8:** the distribution of P-values is mostly uniform across all chromosomes.
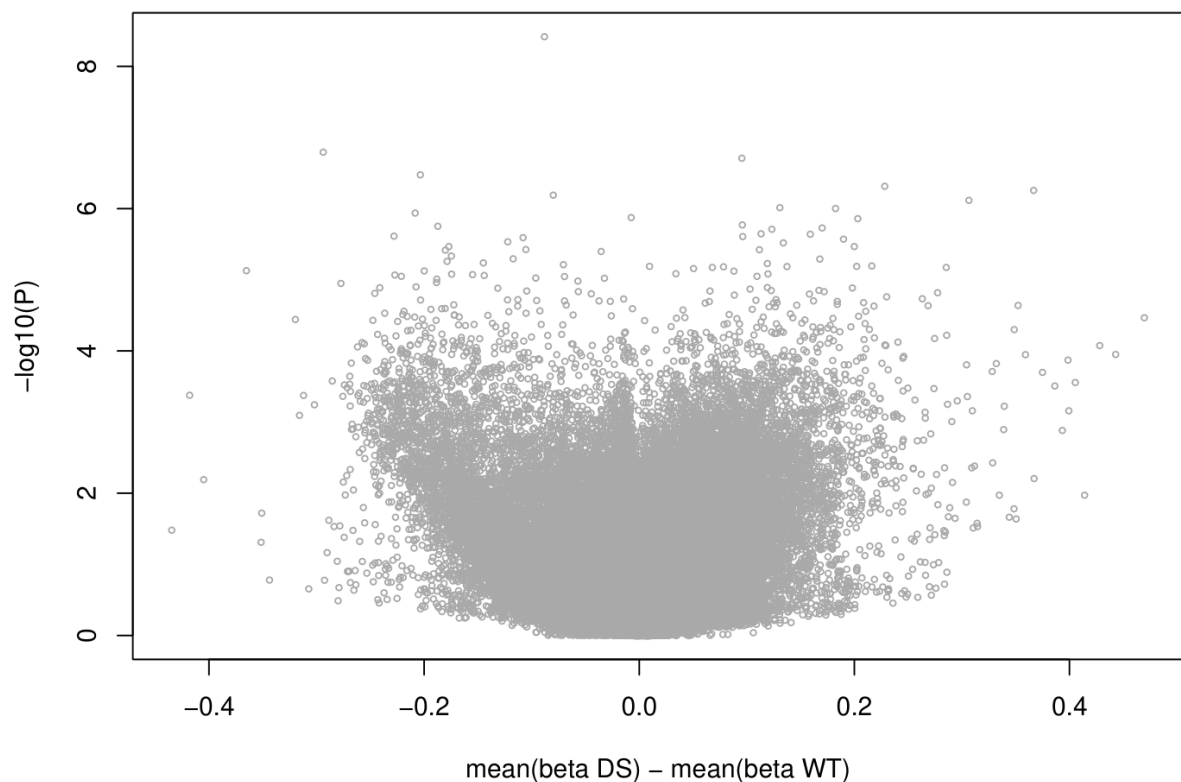
## Volcano plot

- Produce a Volcano plot of the uncorrected P-values:

```
1    to_volc <- data.frame(mean_ds_beta - mean_wt_beta)
2    sites <- rownames(to_volc)
3    to_volc <- data.frame(sites,to_volc)
4
5    colog_p <- data.frame(-log10(p_values_anova))
6    sites <- rownames(colog_p)
7    colog_p <- data.frame(sites,colog_p)
8
9    to_volc <- merge(to_volc,colog_p,by="sites")
10   names(to_volc) <- c("sites","delta","colog_p")
11
12   pdf(file="step13.2.pdf",width=8,height=6)
13   plot(to_volc$delta, to_volc$colog_p, main="", xlab="mean(beta DS) -
     mean(beta WT)", ylab="-log10(P)",cex=0.5,col="darkgrey")
14   dev.off()
```

**Figure 9:** the difference in methylation level doesn't seem to be varying significantly at different reported P-values.
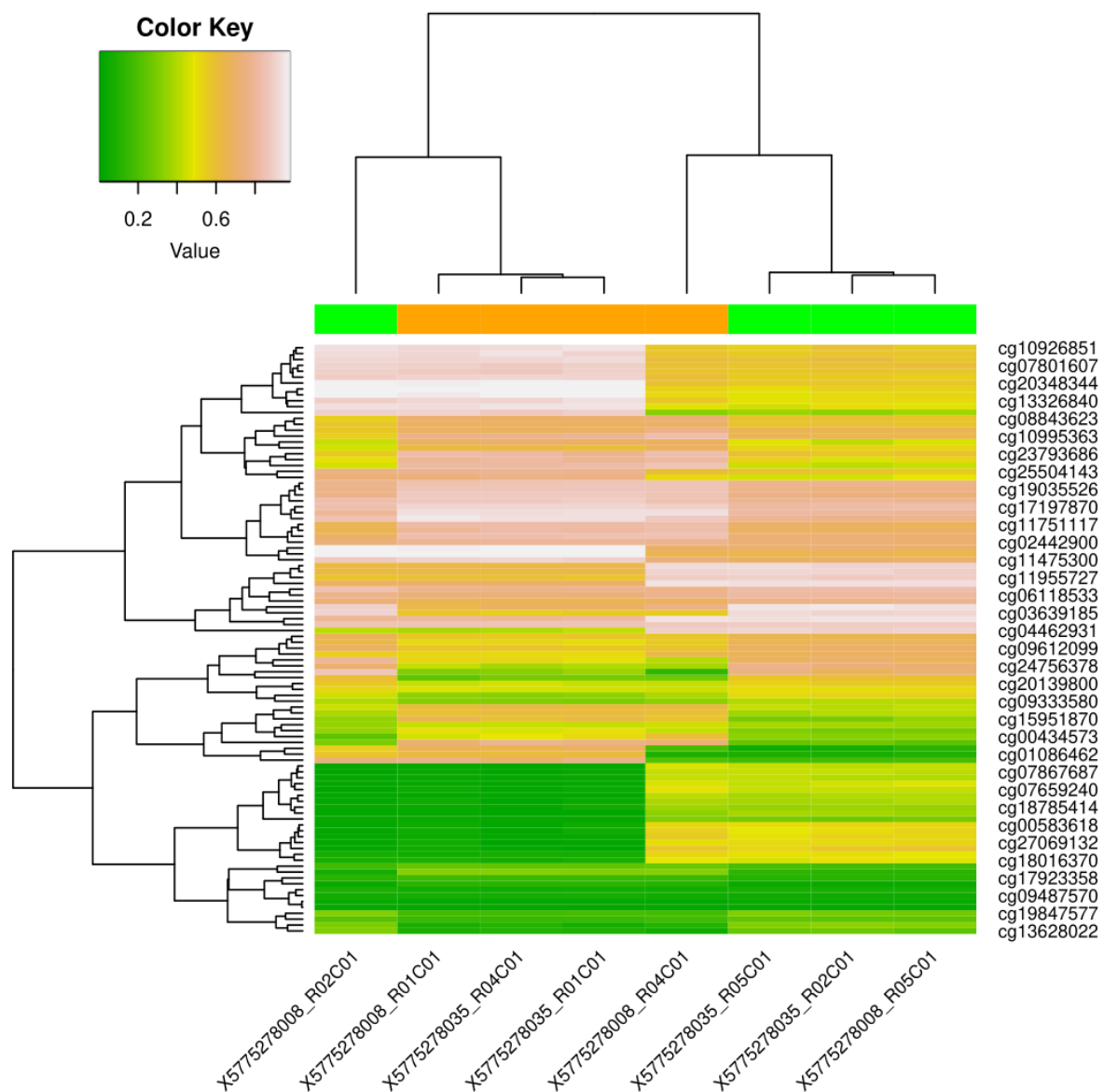
## Step 14

### Heatmap

- Produce a heatmap of the top one hundred methylated probes:

```
1   library(gplots)
2
3   final_anova_annotated <-
    final_anova_annotated[order(final_anova_annotated$p_values_anova),]
4
5   matrix=as.matrix(final_anova_annotated[1:100,2:9])
6   targets$Group
7   # [1] DS WT DS WT DS WT DS WT
8   # Levels: DS WT
9   colorbar <- rep(c("orange","green"),4)
10
11  pdf("step14.complete.pdf")
12  par(mar=c(2,2,2,2))
13  heatmap.2(matrix,col=terrain.colors(100),Rowv=T,Colv=T,
14            dendrogram="both",key=T,ColSideColors=colorbar,
15            density.info="none",trace="none",scale="none",symm=F,
16            cexRow=1,cexCol=1,margins=c(10,8),srtCol=45)
17  dev.off()
```
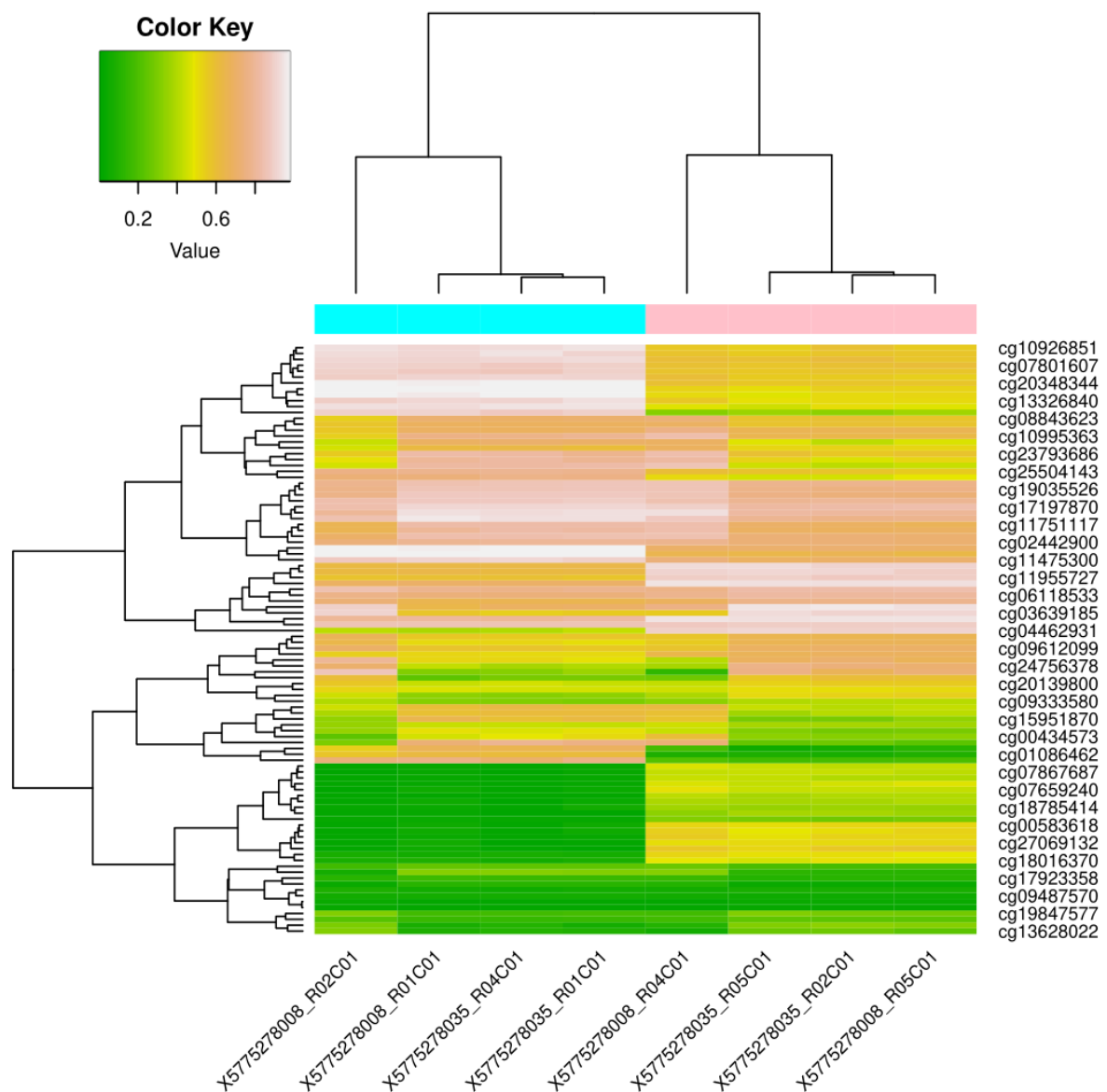
**Figure 10:** the clustering here doesn't seem to be discriminating very well between *WT* and *DS* individuals. Different agglomeration methods were tried other than complete linkage, namely single linkage and unweighted average linkage clustering, but there was no noticeable difference.
It seems that, for this small sample size, correcting for sex has affected more the output than *DS* and *WT* groups, as in fact the algorithm has correctly clusterized the targets on the basis of sex.

```
1  targets$Female
2  # [1] 0 0 1 1 0 1 0 1
3  colorbar <-
   c("cyan","cyan","pink","pink","cyan","pink","cyan","pink")
4
5  pdf("step14.complete.sex.pdf")
6  par(mar=c(15,2,2,2))
7  heatmap.2(matrix,col=terrain.colors(100),Rowv=T,Colv=T,
8            dendrogram="both",key=T,ColSideColors=colorbar,
9            density.info="none",trace="none",scale="none",symm=F,
10           cexRow=1,cexCol=1,margins=c(10,8),srtCol=45)
11 dev.off()
```

**Figure 11:** same heatmap but with sex labelled targets.

## Step 15

### Chromosome 21

- Plot the mean beta density only of those probes that map to chromosome 21 (of which there are 4243):

```
 1  data <- data.frame(rownames(beta),beta)
 2  colnames(data)[1] <- "IlmnID"
 3  data <- merge(data,Illumina450Manifest[,c(1,12)],by="IlmnID")
 4  rownames(data) <- data[,1]
 5  data <- data[,-1]
 6  data <- subset(data,data$CHR=="21")
 7  beta_c21 <- data[,1:8]
 8  rm(data)
 9
10  dim(beta_c21)
11  # [1] 4243      8
12
13  wt_names <- targets[targets$Group=="WT","Basename"]
```
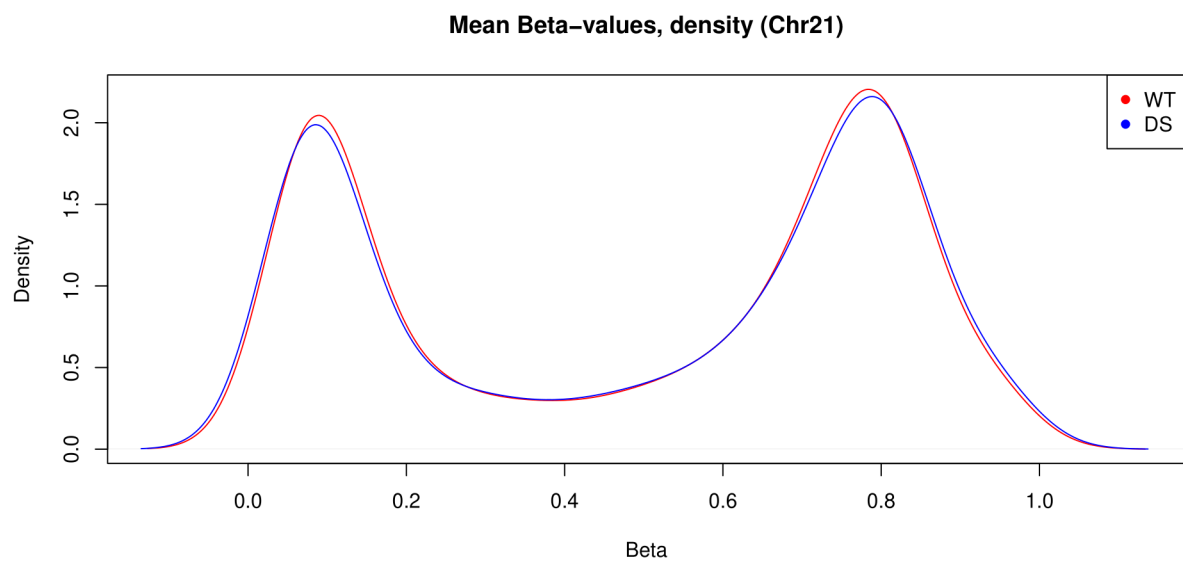
```
14  wt_names <- sub(".+/","X",wt_names)
15  wt_beta_c21 <- beta_c21[,wt_names]
16
17  ds_names <- targets[targets$Group=="DS","Basename"]
18  ds_names <- sub(".+/","X",ds_names)
19  ds_beta_c21 <- beta_c21[,ds_names]
20
21  mean_wt_beta_c21 <- apply(wt_beta_c21,1,mean)
22  mean_ds_beta_c21 <- apply(ds_beta_c21,1,mean)
23
24  pdf(file="step15.pdf",width=10,height=5)
25  plot(density(mean_wt_beta_c21,na.rm=TRUE),col=2,main="Mean Beta-
    values, density (Chr21)",xlab="Beta")
26  lines(density(mean_ds_beta_c21,na.rm=TRUE),col=4)
27  legend('topright',pch=c(16,16),col=c(2,4),legend=c("WT","DS"))
28  dev.off()
```



**Figure 12:** the two lines are so similar here that it doesn't seem there is a significant difference in methylation level between *WT* and *DS* on chromosome 21.