# RNA-Seq

## Aula 3:
## Normalização e análise diferencial

Edgar Kozlova

Gabriela Luiz

# O Curso

**Pré requisitos obrigatorio:**

**Notebook, WiFi**, Notepad++, R e Rstudio

**Programação das aulas:**

1. Banco de dados:
NCBI/SRA
NCBI/GEO

2. RStudio e Instalação de pacotes
edgeR, limma, pheatmap, gplots, ROTS

3. Normalização e Análise Diferencial
voom, RPKM, FPKM, TPM, CPM, counts

4. Análise Diferencial e Visualização
Script, MAplot, VolcanoPlot, Heatmap, Venn

# Objetivo

Introduzir os principais conceitos de normalizacao de contagems.

Aplicar estes metodos sobre dados reais.

# Carregar bioconductor / dados e definir os grupos da análise

```
Untitled1* ×
                                    Source on Save  Q  🪄 ▾  ▤                    ⇥ Run  ⤴⇥  ⇥ Source  ▾  ≡

 8   ### LOAD BIOCONDUCTOR
 9 ▾ ###############################
10   source("https://bioconductor.org/biocLite.R")
11   biocLite()
12   biocLite("ROTS")
13   biocLite("limma")
14   biocLite("pheatmap")
15   biocLite("gplots")
16   biocLite("edgeR")
17   biocLite("RColorBrewer")
18 ▾ ###############################
19   ### LOAD PACKAGES
20 ▾ ###############################
21   library(pheatmap)
22   library(ROTS)
23   library(limma)
24   library(gplots)
25   library(edgeR)
26   library(RColorBrewer)
27 ▾ ###############################
28   ### LOAD DATA
29 ▾ ###############################
30   dados<-read.table("GSE107218_CBPB-hg19-counts.txt",sep="\t",header=TRUE)
31   ### groups
32   group <- as.factor(c(rep("CB_CD34",3),rep("CB_BFUE",3),rep("CB_CFUE",3),rep("CB_PRO",3),
33                        rep("CB_EBASO",3),rep("CB_LB",3),rep("CB_POLY",3),rep("CB_ORTHO",3),
34                        rep("PB_CD34",3),rep("PB_BFU",3),rep("PB_CFU",3),rep("PB_PRO",3),
35                        rep("PB_EBASO",3),rep("PB_LB",3),rep("PB_POLY",3),rep("PB_ORTHO",3)))
36 ▾ ###############################
```

# Tratamento dos reads/counts prévios a normalização

$$RPKM = \frac{\text{Num. de reads mapeados a um gene x } 10^3 \text{ x } 10^6}{\text{Total de reads mapeados x tamanho do gene}}$$

$$RPM = \frac{\text{Num. de reads mapeados a um gene x } 10^6}{\text{Total de reads mapeados}}$$

$$TPM = \frac{\text{Num. de reads mapeados a um gene x tamanho do read x } 10^6}{\text{Total de reads mapeados x tamanho do gene}}$$

$$CPM = \frac{\text{Num. de reads mapeados a um gene}}{\text{Total de reads mapeados}} \text{ x } 10^6$$

https://doi.org/10.1186/gb-2014-15-2-r29

# Transformar os dados de reads/counts para CPM e filtrar genes com linhas vazias

```
36   ##################################
37   ### CPM https://doi.org/10.1186/gb-2014-15-2-r29
38   ##################################
39   x<-dados[,7:54]
40   rownames(x)<-dados$Geneid
41   cpm <- cpm(x)
42   lcpm <- cpm(x, log=TRUE)
43   ### zeroes removal
44   table(rowSums(x==0)==48)
45   ### filter
46   keep.exprs <- rowSums(cpm>1)>=3
47   x <- x[keep.exprs,]
48   ##################################
```

# Figura para verificar a filtragem dos genes com linhas vazias

```
45    x <- x[keepexprs,]
46 ▾  ##############################
47    ### plot to compare to unfiltered data
48 ▾  ##############################
49    nsamples <- ncol(x)
50    col <- brewer.pal(nsamples, "Paired")
51    par(mfrow=c(1,2))
52    plot(density(lcpm[,1]), col=col[1], lwd=2, ylim=c(0,0.21), las=2,
53        main="", xlab="")
54    title(main="A. Raw data", xlab="Log-cpm")
55    abline(v=0, lty=3)
56 ▾  for (i in 2:nsamples){
57      den <- density(lcpm[,i])
58      lines(den$x, den$y, col=col[i], lwd=2)
59    }
60    legend("topright", legend=group, text.col=col, bty="n")
61    ###
62    lcpm <- cpm(x, log=TRUE)
63    plot(density(lcpm[,1]), col=col[1], lwd=2, ylim=c(0,0.21), las=2,
64        main="", xlab="")
65    title(main="B. Filtered data", xlab="Log-cpm")
66    abline(v=0, lty=3)
67 ▾  for (i in 2:nsamples){
68      den <- density(lcpm[,i])
69      lines(den$x, den$y, col=col[i], lwd=2)
70    }
71    legend("topright", legend=group, text.col=col, bty="n")
```

**A. Raw data**

**B. Filtered data**

CB_CD34
CB_CD34
CB_CD34
CB_BFUE
CB_BFUE
CB_BFUE
CB_CFUE
CB_CFUE
CB_CFUE
CB_PRO
CB_PRO
CB_PRO
CB_EBASO
CB_EBASO
CB_EBASO
CB_LB
CB_LB
CB_LB
CB_POLY
CB_POLY
CB_POLY
CB_ORTHO
CB_ORTHO
CB_ORTHO
PB_CD34
PB_CD34
PB_CD34
PB_BFU
PB_BFU

# Normalização
# TMM – Trimmed Mean of M values

g   = gene

k   = libraria

Lg = tamanho do gene

Nk = total number of reads

Sk = total rna sample

$$E[Y_{gk}] = \frac{\mu_{gk} L_g}{S_k} N_k$$

$$\text{where } S_k = \sum_{g=1}^{G} \mu_{gk} L_g;$$
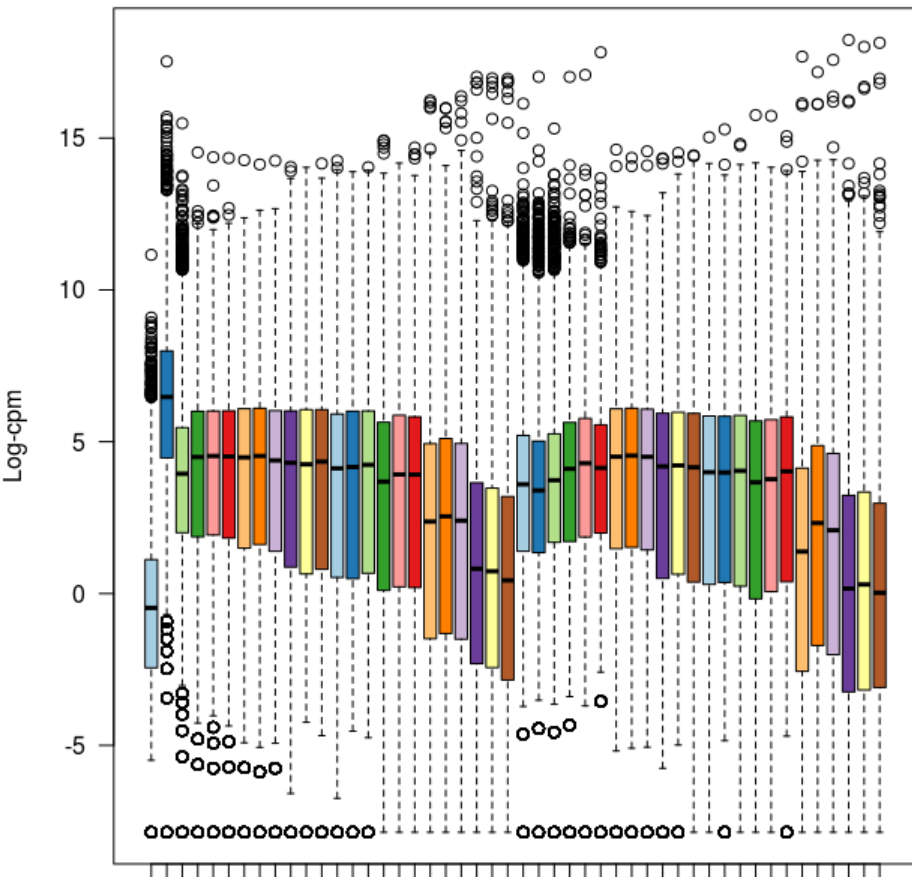
$$M_g = \log_2 \frac{Y_{gk}/N_k}{Y_{gk'}/N_{k'}}$$

$$A_g = \frac{1}{2}\log_2\left(Y_{gk}/N_k \bullet Y_{gk'}/N_{k'}\right) \text{ for } Y_{g\bullet} \neq 0$$
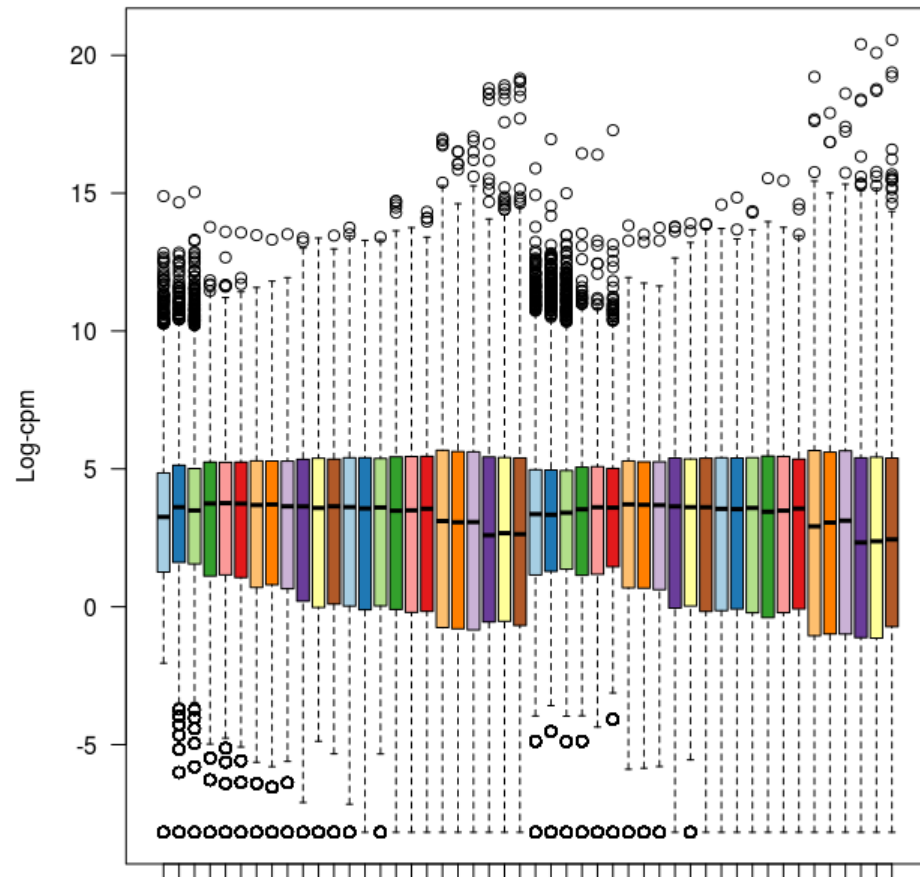
# NORMALIZAÇÃO "TMM" method

```
72 ▾ #############################
73   ### normalize modified data
74 ▾ #############################
75   d.cpm.x <- DGEList(counts=x,group=group) #Create a DGEList object
76   d.cpm.x <- calcNormFactors(d.cpm.x, method = "TMM") #method TMM
77   #d.cpm.x$samples$norm.factors # to see the factors
78 ▾ #############################
79   ### comparison to unormalized data
80 ▾ #############################
81   ### makes a not normalized dataset
82   d.cpm.x2 <- d.cpm.x
83   d.cpm.x2$samples$norm.factors <- 1
84   d.cpm.x2$counts[,1] <- ceiling(d.cpm.x2$counts[,1]*0.05)
85   d.cpm.x2$counts[,2] <- d.cpm.x2$counts[,2]*5
86 ▾ #############################
87   ### plot the unnormalized data and the normalized together
88 ▾ #############################
89   par(mfrow=c(1,2))
90   lcpm <- cpm(d.cpm.x2, log=TRUE)
91   boxplot(lcpm, las=2, col=col, main="")
92   title(main="A. Example: Unnormalised data",ylab="Log-cpm")
93   d.cpm.x2 <- calcNormFactors(d.cpm.x2,method = "TMM")
94   d.cpm.x2$samples$norm.factors
95   lcpm <- cpm(d.cpm.x2, log=TRUE)
96   boxplot(lcpm, las=2, col=col, main="")
97   title(main="B. Example: Normalised data",ylab="Log-cpm")
98 ▾ #############################
```

# Normalização "TMM"



A. Example: Unnormalised data

B. Example: Normalised data

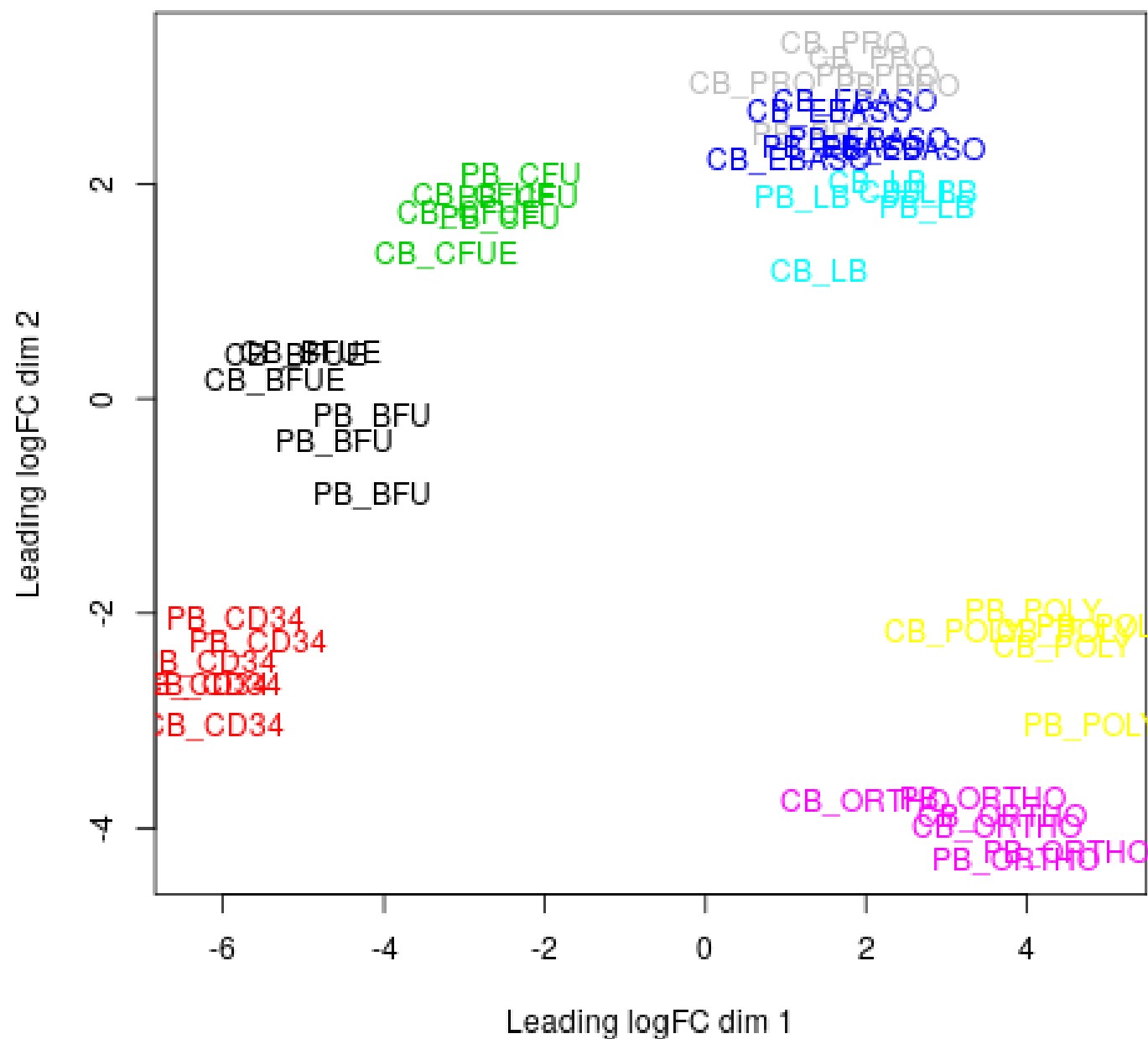# MultiDimensionalScaling

```r
 91  boxplot(lcpm, las=2, col=col, main="")
 92  title(main="A. Example: Unnormalised data",ylab="Log-cpm")
 93  d.cpm.x2 <- calcNormFactors(d.cpm.x2,method = "TMM")
 94  d.cpm.x2$samples$norm.factors
 95  lcpm <- cpm(d.cpm.x2, log=TRUE)
 96  boxplot(lcpm, las=2, col=col, main="")
 97  title(main="B. Example: Normalised data",ylab="Log-cpm")
 98  dev.off()
 99  ###############################
100  ### MDS
101  ###############################
102  lcpm <- cpm(d.cpm.x, log=TRUE)
103  plotMDS(lcpm, labels=group, col=as.numeric(group))
104  title(main="MDS - Sample groups")
105  ###############################
106  ### Voom
107  ###############################
108  design = model.matrix( ~ 0 + group, data=d.cpm.x$samples)
109  colnames(design) <- levels(group)
110  d.cpm.x = estimateCommonDisp(d.cpm.x, verbose=TRUE)
111  d.cpm.x = estimateTagwiseDisp(d.cpm.x)
112  par(mfrow=c(1,2))
113  v <- voom(d.cpm.x, design, plot=TRUE)
114  ###############################
```

MDS - Sample groups

# Limma - Voom !!!



voom: Mean-variance trend

a

b
gene-wise
mean-variance trend
lowess fit

c
gene-wise
mean-variance trend
sqrt standard
deviation for
observation 2, $\sqrt{s^*_{g2}}$
sqrt standard deviation
observation 1, $\sqrt{s^*_{g1}}$
log2 count for
observation 2,
$\hat{\kappa}_{g2}$
log2 count for
observation 1,
$\hat{\kappa}_{g1}$

Sqrt( standard deviation )
Average log2(count size + 0.5)
Fitted log2(count size + 0.5)



(a) Equal library sizes

voom
limma trend
PoissonSeq
edgeR classic
edgeR glm
DESeq
DSS
baySeq
limma notrend
t-test
TSPM

(b) Unequal library sizes

voom
limma trend
edgeR glm
edgeR classic
PoissonSeq
baySeq
DESeq
limma notrend
t-test
DSS
TSPM

False discoveries
Genes chosen

# Limma – Voom !!!

Everything before limma is foreplay
Edgar 2018

```
107    ### Voom
108 ▾  ################################
109    design = model.matrix( ~ 0 + group, data=d.cpm.x$samples)
110    colnames(design) <- levels(group)
111    d.cpm.x = estimateCommonDisp(d.cpm.x, verbose=TRUE)
112    d.cpm.x = estimateTagwiseDisp(d.cpm.x)
113    par(mfrow=c(1,2))
114    v <- voom(d.cpm.x, design, plot=TRUE)
115 ▾  ################################
```

**voom: Mean-variance trend**

Sqrt( standard deviation )

log2( count size + 0.5 )

**Final model: Mean Variance Trend**

sqrt(sigma)

Average log-expression