

20.11.13/23.03.14, Pauli Tikka, University of Turku, thesis work

The complete code for handling the gene data

Contents:

mRNA processing

Step 1: Preliminary set-ups,

Step 2-5: filtering,

Step 6: Clustering,

Step 7: Heat Maps,

Step 8: Wilcoxon's tests,

Step 9: Fisher's exact tests,

Step 10: Piano package, or alternative pathway processing

miRNA processing

Step 11: Preliminary set-ups

Step 12: Preprocessing of miRNA data

Step 13: T-tests for miRNAs (defining de-miRNAs)

Step 14: Performing Gene enrichment with Fisher test in a more elaborated fashion

Step 15: Linear modelling of miRNAs

Step 15a: Enhancing the data before the linear modelling scheme

Step 15b: Defining the lists and mirs for linear modelling

Step 15c: Clustering and Enrichment of one of the lists

Step 15d: Linear modelling

Step 15e: Cross-validation

Step 15f: Obtaining results after linear modelling analysis

STEP 1

Setting working directory

setwd("C:/Users/Pauli/Documents/thesis/labour/analysis/TCGA/results")

Reading the expression text file as a table (matrix), loading takes ~3 minutes:

real <- read.table(file="tcga.breastcancer.rnaseq2.gem.txt", header = TRUE, row.names=1, sep="\t")

Transforming the genes to rows:

ret <- t(real)

Reading Vijay's annotation file as a matrix (anno):

anno <- read.table(file="tcga.breastcancer.rnaseq2.gem.sampleannotation2.txt", header = TRUE, row.names=1, sep="\t")

Reading a file for original gene symbols so that the TCGA's expression matrix symbols can be interpreted:

gsymb <- read.table(file="genesymbols.txt", header = TRUE, row.names=NULL, sep="\t")

Notice: row.names=NULL (it caused problems at later stages)

Constructing an expression set (with annotation matrix (anno), and original transposed expression matrix (ret):

```
metadata = data.frame(
  labelDescription = c("Exp", "Sampletype", "TripleNegative", "Gender", "Age.at.Initial.Pathologic.Diagnosis",
    "ER.Status", "PR.Status", "HER2.Final.Status", "Tumor", "Tumor..T1.Coded", "Node", "Node.Coded", "Metastasis",
    "Metastasis.Coded", "AJCC.Stage",
    "Converted.Stage", "Vital.Status", "OS.event", "OS.Time", "PAM50.mRNA", "RPPA.Clusters"),
  row.names=c("No/Exp", "normal/tumour", "ntn/tn", "Male/Female", "age", "er", "pr", "her2", "tumour", "tumourT1code",
    "Node", "NodeCoded", "Metastasis",
    "MetastasisCoded", "AJCCStage", "ConvertedStage", "VitalStatus", "OSEvent", "OSTime", "PAM50mRNA", "RPPA.Clusters"))
```

```
row.names=c("No/Exp", "normal/tumour", "ntn/tn", "Male/Female", "age", "er", "pr", "her2", "tumour", "tumourT1code",
  "Node", "NodeCoded", "Metastasis",
  "MetastasisCoded", "AJCCStage", "ConvertedStage", "VitalStatus", "OSEvent", "OSTime", "PAM50mRNA", "RPPA.Clusters"))
```

```
pheDa=new("AnnotatedDataFrame",data=anno, varMetadata = metadata)
eset <- new("ExpressionSet", exprs = ret, phenoData = pheDa)
```

The log transform is needed for the heat maps!!

ret_log=log2(ret+1)

exprs(eset)=log2(exprs(eset)+1)

Discarding outlier patient samples that are more than 3.5 stds away from normal:

library(WGCNA)

allowWGCNAThreads(60)

par(mfrow=c(3,1))

```

IAC=cor(exprs(eset),use="p")
hist(IAC,sub=paste("Mean=",format(mean(IAC[upper.tri(IAC)]),digits=3)), xlab= ' IAC values')
cluster1= hclust(as.dist(1-IAC), method="average")
n=eset$Sampletype
n=as.matrix(n)
lgc1=as.matrix(n=='tumour')
n[lgc1[,1],1]=1
meanIAC=apply(IAC,2,mean)
sdCorr=sd(meanIAC)
numbersd=(meanIAC-mean(meanIAC))/sdCorr #Where did this formula came from
table(abs(numbersd)>3.6) #FALSE TRUE 915 6
plot(numbersd,ylab='Number of standard deviation')
abline(h=-3.5,col='red')
#Removing 5 outliers:
lgc2=abs(numbersd)>3.6
anno=anno[lgc2==FALSE,]
M_0=M_0[lgc2==FALSE]
e_0=e_0[lgc2==FALSE]

#Removing suspicious annotation:
lgc3=anno[,2]=='normal'
lgc4=anno[,3]==TRUE #TN status
lgc5=lgc3 & lgc4
M_filt_ok=M_filt_ok[lgc5==FALSE]
e_filt_ok=e_filt_ok[,lgc5==FALSE]

# Clustering and right groups defined by heatmap analysis, here is the function for heatmap (hm should be
expression matrix):
plot(cluster1,cex=0.7,labels=n,plotMyHeatmap <- function(filename, hm){
  library(pheatmap)
  e_filt_ok$TripleNegative[hm$TripleNegative=='TRUE']<-'TripleNegative'
  hm$TripleNegative[hm$TripleNegative=='FALSE']<-'Control'
  hm$Vital.Status=as.vector(hm$Vital.Status)
  hm$Vital.Status[is.na(hm$Vital.Status)]<-'NoInfo'
  hm$ER.Status=as.vector(hm$ER.Status)
  hm$ER.Status[is.na(hm$ER.Status)]<-'NoInfo'
  hm$ER.Status[hm$ER.Status=='Indeterminate']<-'Other'
  hm$ER.Status[hm$ER.Status=='Performed but Not Available']<-'Other'
  hm$ER.Status[hm$ER.Status=='Not Performed']<-'Other'
  hm$PR.Status=as.vector(hm$PR.Status)
  hm$PR.Status[is.na(hm$PR.Status)]<-'NoInfo'
  hm$PR.Status[hm$PR.Status=='Indeterminate']<-'Other'
  hm$PR.Status[hm$PR.Status=='Performed but Not Available']<-'Other'
  hm$PR.Status[hm$PR.Status=='Not Performed']<-'Other'
  hm$HER2.Final.Status=as.vector(hm$HER2.Final.Status)
  hm$HER2.Final.Status[is.na(hm$HER2.Final.Status)]<-'NoInfo'
  hm$HER2.Final.Status[hm$HER2.Final.Status=='Not Available']<-'NoInfo'
  hm$AJCC.Stage=as.vector(hm$AJCC.Stage)
  hm$AJCC.Stage[is.na(hm$AJCC.Stage)]<-'NoInfo'
  hm$AJCC.Stage[hm$AJCC.Stage=='[Not Available]']<-'NoInfo'
  hm$PAM50.mRNA=as.vector(hm$PAM50.mRNA)
  hm$PAM50.mRNA[is.na(hm$PAM50.mRNA)]<-'NoInfo'
  hm$Metastasis=as.vector(hm$Metastasis)
  hm$Metastasis[is.na(hm$Metastasis)]<-'NoInfo'
  hm$Tumor=as.vector(hm$Tumor)
  hm$Tumor[is.na(hm$Tumor)]<-'NoInfo'
  hm$Age.at.Initial.Pathologic.Diagnosis=as.vector(hm$Age.at.Initial.Pathologic.Diagnosis)

```

```

hm$Age.at.Initial.Pathologic.Diagnosis[is.na(hm$Age.at.Initial.Pathologic.Diagnosis)]<- 0
hm$OS.Time=as.vector(hm$OS.Time)
hm$OS.Time[is.na(hm$OS.Time)]<- 0

annotation = data.frame(Var1 = hm$Sampletype, Var2 = hm$TripleNegative, Var3 =
    hm$Vital.Status,
    Var4 = hm$ER.Status, Var5 = hm$PR.Status, Var6 = hm$HER2.Final.Status,
    Var7 = hm$AJCC.Stage, Var8 = hm$PAM50.mRNA, Var9 = hm$Metastasis,
    Var10=hm$Tumor, Var11 = hm$Age.at.Initial.Pathologic.Diagnosis, Var12 =
    hm$OS.Time)
rownames(annotation) <- colnames(hm)
#head(annotation)
annotation$Var1 = factor(annotation$Var1, levels = c("normal", "tumour"))
annotation$Var2 = factor(annotation$Var2, levels = c("Control", "TripleNegative"))
annotation$Var3 = factor(annotation$Var3, levels = c("LIVING", "DECEASED", "NoInfo"))
annotation$Var4 = factor(annotation$Var4, levels = c("Positive", "Negative", 'Other',
    "NoInfo"))
annotation$Var5 = factor(annotation$Var5, levels = c("Positive", "Negative", 'Other',
    "NoInfo"))
annotation$Var6 = factor(annotation$Var6, levels = c("Positive", "Negative", 'Equivocal',
    "NoInfo"))
annotation$Var7 = factor(annotation$Var7, levels = c("Stage I", "Stage IA", "Stage IB", "Stage
    II", "Stage IIA", "Stage IIB", "Stage III",
    "Stage IIIA", "Stage IIIB", "Stage IIIC", "Stage IV", "Stage X",
    "NoInfo"))
annotation$Var8 = factor(annotation$Var8, levels = c('Normal-like', 'Basal-like', 'HER2-
    enriched', 'Luminal A', 'Luminal B', 'NoInfo'))
annotation$Var9 = factor(annotation$Var9, levels = c("M0", "M1", "NoInfo"))
annotation$Var10 = factor(annotation$Var10, levels = c("T1", "T2", 'T3', "T4", 'TX', "NoInfo"))
# Specify colors
Var1 = c("lightgrey", "black")
names(Var1) = c("normal", "tumour")
Var2 = c("steelblue1", "red")
names(Var2) = c("Control", "TripleNegative")
Var3 = c("lightgreen", "darkgreen", "grey")
names(Var3) = c("LIVING", "DECEASED", "NoInfo")
Var4 = c("red", "green", "black", "grey")
names(Var4) = c("Positive", "Negative", 'Other', "NoInfo")
Var5 = c("red", "green", "black", "grey")
names(Var5) = c("Positive", "Negative", 'Other', "NoInfo")
Var6 = c("red", "green", "black", "grey")
names(Var6) = c("Positive", "Negative", 'Equivocal', "NoInfo")
Var7 = c("aquamarine", "aquamarine3", "aquamarine4", "deepskyblue", "blue", "navy",
    "pink", "hotpink", "orangered", "red", "magenta", "black", "grey")
names(Var7) = c("Stage I", "Stage IA", "Stage IB", "Stage II", "Stage IIA", "Stage IIB", "Stage
    III",
    "Stage IIIA", "Stage IIIB", "Stage IIIC", "Stage IV", "Stage X", "NoInfo")
Var8 = c("oldlace", "pink4", "seagreen", "royalblue4", "purple3", "grey")

names(Var8) = c('Normal-like', 'Basal-like', 'HER2-enriched', 'Luminal A', 'Luminal B',
    'NoInfo')
Var9 = c("white", "black", "grey")
names(Var9) = c("M0", "M1", "NoInfo")
Var10 = c("yellow", "wheat4", "springgreen", "thistle1", "olivedrab", "grey")
names(Var10) = c("T1", "T2", 'T3', "T4", 'TX', "NoInfo")
Var11 = c("red")
Var12 = c("red")

```

```

ann_colors = list(Var1 = Var1, Var2 = Var2, Var3 = Var3, Var4 = Var4,
                  Var5 = Var5, Var6 = Var6, Var7 = Var7, Var8 = Var8,
                  Var9 = Var9, Var10 = Var10, Var11 = Var11, Var12 = Var12)
pdf(file = filename,width=11, height=11, pointsize=12)
pheatmap(exprs(hm),fontsize=8, scale='column', show_rownames = F, show_colnames =
          F,annotation=annotation, annotation_colors = ann_colors)
dev.off()
}
plotMyHeatmap(filename='legendarium', e_filt_ok) #time to plot (more than 8 minutes): 13:22-30='Clustering of
samples')
#The results of heatmap evaluation maybe checked with kmeans:
km_val_mfl=kmeans(M_filt_ok, 5)
# get cluster means
aggregate(M_filt_ok,by=list(km_val_mfl$cluster),FUN=mean) #14:56-57
# append cluster assignment
mydata <- data.frame(M_filt_ok, km_val_mfl$cluster)
# Principal component analysis can also be used for checking the correct groups:
prcomp(M_filt_ok[1:100,1:100])

# Subset expression set for modified TN and normal groups:
load('breastCancerTCGAtriplenegativeVSnormal.Rdata')

#### STEP 2 ####
#Discarding the zero values:
par(mfrow=c(1,2))
hist(unlist(ret))
hist(unlist(ret_log))
## Two approaches follows:
pdf("mean_exp.pdf")
mean_exp=apply(ret_log,1,mean)      # check this by eye -> take the tale away using proper cut-off such as 1.5, #
see sd and var also
hist(mean_exp)                    # mem=mean(mean_exp) ~ 6.3 cannot be used here, because then too
much values are discarded
dev.off()
# within one gene (i.e. ABC in row 123, i.e. "1" in apply), how many patient samples, have expressed that gene above
threshold value of 1.5:
sum=apply(ret_log > 1.5, 1, sum)
# There can be e.g. 175 genes expressed above 1.5 across the whole 1100 patient samples, so the percentage is 175
/ 1100 * 100, and in general:
a = sum/dim(ret_log)[2] * 100
# Finally select only those GENES (in rows (r) in e.g. matrix A(r,c)) that have expression across the all patient samples
above 50%:
M_0 = ret_log[ a > 10,]
dim(M_0)
#[1] 15337 1100
# Apply the same functions also to your gene symbol list (so that the row numbers keep the same) and expression
set matrix
G_0 = gsymb[ a > 10,]
e_0 = eset[ a > 10,]
hist(unlist(M_0))

#### STEP 4 ####
# IQR filtration
GeneIQR=apply(M_0,1,IQR)
GeneIQRg=apply(G_0,1,IQR)
GeneIQRm=mean(GeneIQR)

```

```

GeneIQRm
#[1] 1.313247
## or check the plot
pdf("GeneIQR.pdf")
hist(GeneIQR)
dev.off()
# All off the values which gives 'TRUE' (according to the statement) in
#'which' command (that produces logical values (true/false)) are discarded (-)
M_IQR <- M_0 [ - which (GeneIQR <= 2) , ]
dim(M_IQR)
#2597 1100
G_IQR = G_0[ - which (GeneIQR <= 2) ,]
e_IQR = e_0[ - which (GeneIQR <= 2) ,]
#this is kind of much so I do also variance filtering

#### STEP 5 ####
# Variance filtering
GeneVar=apply(M_IQR,1,var)
GVM=mean(GeneVar)
GVM
#2.888428
## or check the plot
pdf("GeneVar.pdf")
hist(GeneVar)
choose 4
dev.off()
M_filt_ok <- M_IQR [ - which (GeneVar <= 4) , ]
dim(M_filt_ok)
#[1] 1157 1100
G_filt_ok = G_IQR[ - which (GeneVar <= 4) ,]
e_filt_ok = e_IQR[ - which (GeneVar <= 4) ,]

## STEP 6 ###
# CLUSTERING for filtered genes:
method="complete"
distance="euclidean"
pdf('clust_all.pdf')
clust=hclust(dist(t(exprs(e_filt_ok))),method=distance),method=method) #ok
plot(clust, labels=F)
dev.off()

## STEP 7 ###
# Heat map generated with a function, if you do not want to see the names of
# the patient samples and filtered genes, insert F for rownames and columnnames in the pheatmap:
plotMyHeatmap <- function(filename, heat){
  library(pheatmap)
  heat$TripleNegative[heat$TripleNegative=='TRUE']<-'TripleNegative'
  heat$TripleNegative[heat$TripleNegative=='FALSE']<-'Control'
  annotation = data.frame(Var1 = heat$TripleNegative, Var2 = heat$Sampletype)
  rownames(annotation) <- colnames(heat)
  #head(annotation)
  annotation$Var1 = factor(annotation$Var1, levels = c("Control", "TripleNegative", "NA"))
  annotation$Var2 = factor(annotation$Var2, levels = c("normal", "tumour"))
  # Specify colours:
  Var1 = c("steelblue1", "red", "yellow")
  names(Var1) = c("Control", "TripleNegative")
  Var2 = c("lightgrey", "black")

```

```

names(Var2) = c("normal", "tumour")
ann_colors = list(Var1 = Var1, Var2=Var2)
pdf(file = filename,width=11, height=11, pointsize=12)
pheatmap(exprs(heat),fontsize=8, scale='column', show_rownames = T, show_colnames = T,annotation=annotation,
annotation_colors = ann_colors, main="21.11.2013, Pauli Tikka, Heat map")
dev.off()
}

```

Insert a name of the file, and the filtered expression set matrix (e_filt_ok) for the arguments of this heat map function:

```

plotMyHeatmap('Tikka_heat_map_2.pdf',e_filt_ok)          # with this function one cannot add "x, and y labels",
they need to be added with picture editor.

```

STEP 8

Wilcox's tests, (t-tests similarly, with loop)

Gene symbol list with only true values:

```
g_e=gsymb[(eset$Exp == T),]
```

Normal samples:

```
esn=exprs(esetSub[,esetSub$Sampletype=='normal'])
```

#Triple negative samples:

```
est=exprs(esetSub[,esetSub$TripleNegative=='TRUE'])
```

with real (r) groups (est, esn (normi)):

```
wr=1:20531 * 0
```

```
for(i in 1:20531){hr=wilcox.test(est[i,],esn[i,]);wr[i]=hr[[3]];
```

```
hist(log10(wr),100)
```

Bonferroni method: every p-value is multiplied by the number of performed tests (here: 20,000 (genes))

Do this analysis again!:

```
swr=sum(wr[is.nan(wr)==F]<1e-25)
```

```
s_r=wr[is.nan(wr)==F]<1e-25
```

For gene list:

```
gr=g_e[s_r,]
```

```
grf=gr[,1] # Only Entrez number is required at next tests
```

STEP 9

Fisher's exact tests

This function works but the matrix is not good and it needs the real pathway names, and not the KEGG pathway names

```
pwy_func2=function(pathway.list2,grfx,s_r) {
```

```
  for (i in 1:length(pathway.list2)) {
```

```
    de_tot=length(grfx)
```

```
    g_tot=length(s_r)
```

```
    p.val=1:length(pathway.list2) * 0;
```

```
    pwy_oxp_l = length(pathway.list2[[i]]);
```

```
    pwy_oxp_de=sum(grfx %in% pathway.list2[[i]]);
```

```
    pwy_de = de_tot - pwy_oxp_de;
```

```
    pwy_oxp_n = pwy_oxp_l - pwy_oxp_de;
```

```
    pwy_nn = g_tot - de_tot - pwy_oxp_n;
```

```
    fpgs = matrix(c(pwy_oxp_de, pwy_de, pwy_oxp_n, pwy_nn), 2, 2);
```

```
    p.val = fisher.test(fpgs);
```

```
    at=print(cbind(names(pathway.list2[i]),p.val[1]))
```

```
    #if (p.val <= 0.05) {cbind(names((pathway.list2[i])),p.val[1])}
```

```
    #if (as.matrix(p.val) <= 0.05) return(p.val)
```

```
    #print(as.vector(names(pathway.list2[i])))
```

```

##### HERE SHOULD BE RETURN #####
}
}

#### STEP 10 ####
# Piano package, or alternative pathway processing
# One can apply gsea, and page with runGSA function.
# Gene sets are pathways. Instead of g2pwyid, there could be some other list of genes.
myGsc <- loadGSC(g2pwyid)
GSA = runGSA(...)
# qusage package could be also applied but the enrich internet page gives you all this more easily, and faster
(http://amp.pharm.mssm.edu/Enrichr/index.html)

#### miRNA processing ####
#### STEP 11 ####
# Preliminary steps for microRNA data:
# This is just the directory address where you have saved your TCGA data:
path_to_files = "C:/cygwin64/home/Pauli/miRNA/Level_3"
# Reading all the file names from this directory
files = list.files(path_to_files, pattern="mirna.quantification",all=F)
# Just reading the first file to know how many rows (i.e how many miRNA) are there in the file
tmp <- read.table(paste(path_to_files, files[ 1 ] ,sep = "" ), header=T,as.is=T,sep="\t")
# Creating an empty matrix where to combine all the files
data <- matrix(nrow = nrow (tmp))
# Naming the rows of the data matrix
rownames(data) <- tmp$miRNA_ID
# Removing the tmp variable that we wont use anymore
rm(tmp)
# Loop over all the files to read
for(i in 1:length(files)) {
  tab <- read.table(paste(path_to_files, files[ i ] ,sep = "" ), header=T,as.is=T,sep="\t") ## This reads the first file and
so on over the loop of all files
  if(nrow(tab) == nrow(data)) {
    tab <- tab[order(match(rownames(data),rownames(tab))),]
    data = cbind(data, tab$reads_per_million_miRNA_mapped) ## This will read only the column that you would like to
have i.e read_per_million
  } }
data <- data[, -1]
# Naming the columns with the patient names (IDs)
colnames(data) <- substr(files,1,15)
## OR ##
# The way to import all the (miRNA) data (including miRNA isoform expressions) from
# a directory's files to an R list is to use foreign package # http://www.ats.ucla.edu/stat/r/pages/read\_multiple.htm
library(foreign)
# This is the directory:
path_to_files = "C:/cygwin64/home/Pauli/miRNA/Level_3/"
# Reading all the file names from this directory
files = list.files(path_to_files, pattern="13.isoform.quantification",all=F)
# Construct a variable that has the path of directory and the file names to read as a table (or list in this case) in R
fa <- file.path(path_to_files, files)
# Lapply command has several functions such as read.dta or read.csv, check the options from above internet site
da <- lapply(fa, read.table)
# Selecting certain values of this list: da_test=c()
for (i in 1:length(dat))( da_test=qpcR::cbind.na(da_testi, dat[[i]][6]), dat[[i]][4]) #so far so good...
# Then renaming the matrix
sum_mir_dvl=da_test[,-1]
# Preparing the sum matrix from sum_mir_dvl

```

```

# Remove star/mature beginnings from names
# First making the empty matrix: su_mmat=matrix(nrow=7702, ncol=1540) su_mmat = "[<-"(su_mmat,value=0)
# This loop is done three times (so i times and then repeated 3 times), where the pattern is first mature, and then
star
for (i in 1:770){#su_mmat[, (i*2)]=gsub(pattern="star," , replacement="", su_mmat[, (i*2)], ignore.case = FALSE, perl =
FALSE, fixed = FALSE, useBytes = FALSE)
su_mmat[, ((i*2)-1)] = sum_mir_dvl[, ((i*2)-1)]
#Changing the colnames of the new matrix:
colnames(su_mmat)=colnames(sum_mir_dvl)
# This is the long wanted sum matrix for mirnas with corresponding mature names and their sum expression values:
# (note: agstest=matrix(nrow=1200, ncol=1540) agstest = "[<-"(agstest,value=0) =
for (i in 1:770) (agstest=qpCR::cbind.na(agstest,aggregate(su_mmat[, (i*2-1)], list(su_mmat[, (i*2)]), sum)) )
# Defining the number of rows and their names in agstest matrix according to unique MIMAT-names in agstest in
overall:
# http://stackoverflow.com/questions/12154814/r-get-a-single-column-from-many-columns
result <- data.frame(Wave = unlist(atf,use.names=FALSE))
agstest_rn=unique(result) ## jee
# rownames should not have na-values:
rownames(agstest)=agstest_rn[is.na(agstest_rn)==F,1]
# Then one needs to merge the correct rownames (of agstest matrix) and agstest-mimat names (at every second
column) to obtain the overall matrix dtt:
#defining the new variable before the looping solution:
dtt=matrix(nrow=1048,ncol=1540)
dtt = "[<-"(dtt,value=0) #from art of r programming
dtt=data.frame(dtt)
# The use of merging function was found at (see below), and then applied for every column:
# http://r.789695.n4.nabble.com/How-to-compare-match-two-columns-from-diferent-dataframe-and-assign-values-
from-one-datafram-to-the-r-td2544573.html
for (i in 1:770) (
dtt[, ((i*2-1):(i*2))]=merge(data.frame(x=rownames(agstest)),
data.frame(x=agstest[, (i*2-1)], y=agstest[, (i*2)]), by='x', all.x=T) )
# Because dtt was an empty matrix the rownames, and column names must be inserted again, and also removing not
needed extra-mimat-names, that where
# needed for the merging
dtt=dtt[, (rep(c(FALSE,TRUE),770))]
rownames(dtt)=rownames(agstest)
# Some repetition from previous stages is good:
path_to_files = "C:/cygwin64/home/Pauli/miRNA/Level_3/"
files = list.files(path_to_files, pattern="13.isoform.quantification", all=F)
colnames(dtt) <- substr(files,1,15)

#Tuschen, new:
tush_n <- read.table("TS-PT-PI_21TRUE_complete.txt", fill=TRUE)
#Selecting just the miR-infos:
m <- grep("hsa-*", tush_n[,1], value=FALSE)
tush_n=tush_n[m,]
# Preparing the rownames, and deleting an extra column, and NA rows:
x <- tush_n[,1]
y <- gsub(":", "", x, fixed=T)
rownames(tush_n)=y
tush_n=tush_n[, -1]
#Replacing empty spots as nas (this is important)
tush_n[tush_n==""] <- 'na'
# make a list of tush_n
tush_nl=lapply(seq_len(nrow(tush_n)), function(i) tush_n[i,])
for (i in 1:length(tush_nl)) (
tush_nl[[i]]=as.vector(unlist(tush_nl[[i]])) ) # Ok, now it looks the same

```



```
#Na-values (just values, not rows) away:
tush_n1 <- lapply(tush_n1, function(x) x[!is.na(x)]) # wei hou
# Naming the list values
names(tush_n1)=rownames(tush_n)
```

STEP 12

```
# Preprocessing of miRNA data #
# Selecting nas away from data:
data_in = data
data_in[is.na(data_in)] <- 0
# Selecting zeros away from data
v=rep(1,dim(data_in)[1])
for (i in 1:dim(data_in)[1]){
  sum=sum(data_in[i,]> 10)/dim(data_in)[2] # 10 is ~1% of mean(data_in)
  v[i]<-sum}
datv=data_in[v>0.1,] # sample values should be above 10 for more than 10% of a particular miRNA
```

Annotation and expression redefined

```
x <- rownames(annoi)
y <- colnames(datv)
matches <- which (x %in% y)
a1 <- annoi[matches,]
e1 <- expv[,matches]
# Ordering the datas:
e1=e1[order(colnames(e1))]
a1=a1[order(rownames(a1)),]
datv=datv[,order(colnames(datv))]
# Matching the datav to the annotation (a1)
y <- rownames(a1)
x <- colnames(datv)
matches <- which (x %in% y)
datv2 <- datv[,matches]
```

#Discovering the non-equivalent samples and discarding them

```
colnames(datv2)==rownames(a1)
colnames(datv2)[117:127] # checking the FALSE sites
colnames(e1)[117:127] # comparing FALSE sites to the ok sites
datv2=datv2[,-119] # deleting the FALSE (duplicate) samples
```

STEP 13

Wilcox tests are needed for skewed data

```
# annotation, expression, and mirna data (after preprocessing of 0s and NAs and matching):
```

```
dim(a1)
dim(e1)
dim(datv2)
# Expression set for the ok anno, exp, and data groups:
metadata2 =
  data.frame(labelDescription =
c("Exp","Sampletype","TripleNegative","Gender","Age.at.Initial.Pathologic.Diagnosis","ER.Status","PR.Status","HER2.
Final.Status","Tumor","Tumor..T1.Coded","Node","Node.Coded","Metastasis","Metastasis.Coded","AJCC.Stage","Co
nverted.Stage","Vital.Status","OS.event","OS.Time","PAM50.mRNA","RPPA.Clusters"),
```

```
row.names=c("No/Exp","normal/tumour","ntn/tn","Male/Female","age","er","pr","her2","tumour","tumourT1code
d","Node","NodeCoded","Metastasis","MetastasisCoded","AJCCStage","ConvertedStage","VitalStatus","OSevent","O
STime","PAM50mRNA","RPPAClusters"))
pheDa2=new("AnnotatedDataFrame",data=a1, varMetadata = metadata2)
eset_i <- new("ExpressionSet", exprs = datv2, phenoData = pheDa2)
```

```

# own mirna data experssion
e_d=as.matrix(colnames(exprs(eset_i[1,]))) #567, maybe not needed: e_d=e_d[order(e_d)]
# vijay normal column names
e_n=(as.matrix(colnames(exprs(esetSub[,esetSub$Sampletype=='normal'])))) #87
e_n=e_n[order(e_n)]
# Vijay tn
e_t=(as.matrix(colnames(exprs(esetSub[,esetSub$TripleNegative==TRUE]))))
e_t=e_t[order(e_t)]
# for normals and tns (FINALLY ok)
matches = (e_d %in% e_n)
normals=exprs(eset_i[,matches])
tneg=exprs(eset_i[,matches])
# The de-miRNA:)
normals_lg=log10(normals+1)
tneg_lg=log10(tneg+1)

## Eli tää toimii (Wilcox skewille datalle login jälkeen duplikatet poistettuna ja dimmin sijaan
# loopissa on vain dim-arvot sellaisenaan)
# check also that tneg_lg and normals_lg are as numeric (which maybe applied with as.matrix...) (and not data.frame)
# P-values for all of the genes are also needed, i.e. Wilcox for normals (nrm_i) and tns (t_n3) at e11_lg
matches = (e_d %in% e_n)
nrm_i=e11_lg[,matches]
t_n3=e11_lg[,matches]

# changing data.frames as numeric with as.matrix
nrm_i=as.matrix(nrm_i)
t_n3=as.matrix(t_n3)
# Differentially expressed genes with Wilcox test :
wto=as.vector(1:16458)
wto="[<-"(wto,1:16458,value=0)
for(i in 1:length(wto))
{hr=(wilcox.test(t_n3[i,],nrm_i[i,])); wto[i]=hr[[3]];}
# Bonferroni correction: Take the p-value of each gene and multiply it by the number of genes in the gene list:
wto2=wto*16458
hist(log10(wto2), breaks=25)
# Checking the de_genes:
rownames(t_n3)
de_gen=cbind(rownames(t_n3),as.numeric(wto2))
colnames(de_gen)<-c('gen','pval')
# Selecting the best differentially expressed genes:
de_gen_lgi = de_gen[,2] < 0.05
# dif.exp. mir:s that are ok, according to the pval (in the second column):
de_geno=de_gen[de_gen_lgi,] # dim(de_geno)[1] is needed for Fisher Test

### STEP 14 ###
# Performing Gene enrichment with Fisher test in a more elaborated fashion
# Checking if the dtt or dttv2 data was skewed:
dtt2=as.matrix(dtt) # Somehow my data was not numeric, so I changed so whit as.matrix command
dttv3=as.matrix(dttv2)
hist(log10(dttv3+1), breaks=50) # ok, the data was skewed so using Wilcox was ok. Log. trans needed:
dttv3_lg=log10(dttv3+1)

# Starting points for miRNA enrichment:
de_miro3 # differentially expressed mature miRNAs
targets_ok # the miRNA targets (change the rownames, and na. info, and delete the first column)
gsymb2 # gsymbols for entrez to genes
a11 # Annotation, See 6114_pt, and 141213_pt files (not sure if needed, if I have tn/nrm grps)

```

```

e11          # Expression matrix, See 6114_pt, and 141213_pt files (no log trans done yet)
e11_lg=log10(e11+1) # log -transformation for expression matrix
normals2_lg   # log transformed normal group
normals3      # normals group with mimat names, and vijay 70 patient nrms
tneg2_lg      # log transformed tn group

# Processing the matrices
targets_ok=read.table("targets_ok.txt", header=T, sep="\t")
# rowname checks for targets_ok and e11_lg
rownames(targets_ok)= targets_ok[,1]
rownames(targets_ok)= sub("(r)", "\\U\\1", rownames(targets_ok), perl=TRUE)
targets_ok=targets_ok[,-1]
rownames(e11_lg) = sub("X", "", rownames(e11_lg), perl=TRUE)

# Up and downregulation for de-Genes, Pre-values
matches = (e_d %in% e_t)
nrm_e=e11[,matches]
tneg_e=e11[,matches]
nrm_e=as.matrix(nrm_e)
tneg_e=as.matrix(tneg_e)
# The regulation of genes:
h=as.vector(1:dim(nrm_e)[1])
h = "[<-"(h,1:dim(nrm_e)[1],value=0)
for (i in 1:dim(nrm_e)[1]) (
  h[i]= (median(tneg_e[i,]) - median(nrm_e[i,]))
gen_rgl=cbind(gsyimb2[,1],h)

# de_miRNA list values (demot2) must be compared to gene names of the good p-values (de_genos)
# Function for this
f_match_down=function(de_genos, demot2) {
  big=as.character(de_genos[,1])
  tmd=NULL
  #smalls
  for (i in 1:length(demot2)) (
    tmd[i]=list(big[(big %in% as.character(unlist(demot2[i])))]))
  return(tmd) } #jee

f_match_down2= function(hero, gen_rgl) {
  h=NULL
  tv=gen_rgl[,2]<0
  big=as.character(gen_rgl[tv,1])

  #smalls
  for (i in 1:length(hero)) (
    h[i]=list(big[(big %in% as.character(unlist(hero[i])))]))
  return(sapply(h, length))
}

f_match_down3=function(de_genos, demot2) {
  big=as.character(gen_rgl[,2])
  #smalls
  for (i in 1:length(demot2)) (
    tmd[i]=list(big[(big %in% as.character(unlist(demot2[i])))]))
  return(tmd)}

# genes up (when mirs down)
f_match_down4 =function(hero4, gen_rgl) {

```

```

h=NULL
tv=gen_rgl[,2]>0
big=as.character(gen_rgl[tv,1])
#smalls
for (i in 1:length(hero4)) {
  h[i]=list(big[(big %in% as.character(unlist(hero4[i])))]))
}
return(sapply(h, length))

# This function might be useful:
f_miR_enrich_fish = function (tush_nl,de_miro3,de_geno,gen_rgl) {
#You need a list of miRNA targets (tush_nl), differentially expressed miRNAs (de_miro3),
# differentially expressed mRNAs (de_geno), and finally log fold change values for genes (only I would say); gen_rgl

# Finding targets (trg) of de_mat_mirnas (dmm):
trg_dmm2=tush_nl[de_miro3[,1]]
f_not_z2=function(trg_dmm2) {
  h=NULL
  for(i in 1:length(trg_dmm2))
    (if (length(trg_dmm2[[i]])!=0) h=c(h,i))
  return(trg_dmm2[h])}
trg_dmm2=f_not_z2(trg_dmm2)
# Changing the names of targets:
names(trg_dmm2)=f_mir_to_mature_vec(names(trg_dmm2))

# matching of low p-value matrix, and gen_rgl<0 matrix
# below zero value genes
ge_rgl_z=gen_rgl[,2]<0
ge_rgl_z=gen_rgl[ge_rgl_z,]
big=as.character(ge_rgl_z[,1])
# de_genes
small=as.character(de_geno[,1])
#matching
matches=big %in% small
gen_ok_f=big[matches]

# matching of low p-value matrix, and gen_rgl>0 matrix
# above zero value genes
ge_rgl_z=gen_rgl[,2]>0
ge_rgl_z=gen_rgl[ge_rgl_z,]
big=as.character(ge_rgl_z[,1])
# de_genes
small=as.character(de_geno[,1])
#matching
matches=big %in% small
gen_ok_fu=big[matches] # this is upregulated,

# Selecting upregulated/downregulated miRs (in general):
u=as.matrix(gen_rgl_mir[,2])>0
u_mir=(gen_rgl_mir[u,1])
d=as.matrix(gen_rgl_mir[,2])<0
d_mir=(gen_rgl_mir[d,1])
# Selecting upregulated mimats from my list
um=intersect(names(trg_dmm2),u_mir)
# Selecting downregulated mimats from my list
dm=intersect(names(trg_dmm2),d_mir)
## de_miRNA list values (trg_dmm2) must to be compared to gene names of the good p-values (de_geno)
hero2=f_match_down(de_geno, trg_dmm2) #using ready made function

```

```

names(hero2)=names(trg_dmm2)
# mirs up, genes down
hero3=hero2[um] #ok
t_mru_d=f_match_down2(hero3, gen_rgl)
# mirs down
hero4=hero2[dm] #ok
t_mrd_u=f_match_down4(hero4, gen_rgl)

# Mirs up, genes down Fish:
fi=as.vector(matrix(0,nrow=length(t_mru_d)))
for(i in 1:length(t_mru_d))
{hr=fisher.test(matrix(c(
  t_mru_d[i],
  as.vector(sapply(hero3[i],length))-t_mru_d[i],
  length(gen_ok_f)-t_mru_d[i],
  dim(gen_rgl)[1]-(length(gen_ok_f)-t_mru_d[i])), nrow=2)); fi[i]=hr[[1]]);}

# MIRNA UP (=mu), GENE DOWN (=gd) => mugd
mugd=cbind(names(hero3),fi)
colnames(mugd)<-c('Mature miRNA','Fisher test p-value')
mugd=mugd[order(as.numeric(mugd[,2])),]
mugd[,2]=as.numeric(mugd[,2])
mugd[,1]=f_mature_to_mir_vec(mugd[,1])

# mirs down, genes up
fit=as.vector(matrix(0,nrow=length(t_mrd_u)))
for(i in 1:length(t_mrd_u))
{hr=fisher.test(matrix(c(
  t_mrd_u[i],
  as.vector(sapply(hero4[i],length))-t_mrd_u[i],
  length(gen_ok_fu)-t_mrd_u[i],
  dim(gen_rgl)[1]-(length(gen_ok_fu)-t_mrd_u[i])), nrow=2)); fit[i]=hr[[1]]);}

# MIRNA DOWN (=md), GENE UP (=gd) => mdgu
mdgu=cbind(names(hero4),fit)
colnames(mdgu)<-c('Mature miRNA','Fisher test p-value')
mdgu=mdgu[order(as.numeric(mdgu[,2])),]
mdgu[,1]=f_mature_to_mir_vec(mdgu[,1])

# p-values below 0.05, and their combined matrix (str(mdgu2): chr):
#mir d gen up
a=as.numeric(mdgu[,2])
b=a<0.05
c=mdgu[b,]

#mir u gen d
a=as.numeric(mugd[,2])
b=a<0.05
c1=mugd[b,]

c1c=NULL
c1c=rbind(c1,c)
c1c=cbind(c(rep('d',dim(c)[1]), rep('u',dim(c1)[1])),c1c) #THIS IS not ok
colnames(c1c)<-c('regul_miR','miRNA name','Fisher test p-value')
write.table(c1c, "all_good_ver_alp.txt", sep="\t", row.names=F)
return(c1c)
}

```

```

mir_reg=f_miR_enrich_fish(tush_nl,de_miro3,de_geno,gen_rgl) #ok
# Selecting the below 0.05 p-value miRs
trg_dmm4=trg_dmm3[mir_reg[,2]]

f_m_enr_fish2 = function (tush_nl,de_miro3,de_geno,gen_rgl,gen_rgl_mir) {
  #You need a list of miRNA targets (tush_nl), differentially expressed miRNAs (de_miro3),
  # differentially expressed mRNAs (de_geno), and finally log fold change values for genes (only I would say); gen_rgl

  # Finding targets (trg) of de_mat_mirnas (dmm):
  trg_dmm2=tush_nl[de_miro3[,1]]
  f_not_z2=function(trg_dmm2) {
    h=NULL
    for(i in 1:length(trg_dmm2))
      (if (length(trg_dmm2[[i]])!=0) h=c(h,i))
    return(trg_dmm2[h])}
  trg_dmm2=f_not_z2(trg_dmm2)
  # Changing the names of targets:
  names(trg_dmm2)=f_mir_to_mature_vec(names(trg_dmm2))

  # matching of low p-value matrix, and gen_rgl<0 matrix
  # below zero value genes
  ge_rgl_z=gen_rgl[,2]<0
  ge_rgl_z=gen_rgl[ge_rgl_z,]
  big=as.character(ge_rgl_z[,1])
  # de_genes
  small=as.character(de_geno[,1])
  #matching
  matches=big %in% small
  gen_ok_f=big[matches]

  # matching of low p-value matrix, and gen_rgl>0 matrix
  # above zero value genes
  ge_rgl_z=gen_rgl[,2]>0
  ge_rgl_z=gen_rgl[ge_rgl_z,]
  big=as.character(ge_rgl_z[,1])
  # de_genes
  small=as.character(de_geno[,1])
  #matching
  matches=big %in% small
  gen_ok_fu=big[matches] # this is upregulated,

  # Selecting upregulated/downregulated miRs (in general):
  u=as.matrix(gen_rgl_mir[,2])>0
  u_mir=(gen_rgl_mir[u,1])
  d=as.matrix(gen_rgl_mir[,2])<0
  d_mir=(gen_rgl_mir[d,1])
  # Selecting upregulated mimats from my list
  um=intersect(names(trg_dmm2),u_mir)
  # Selecting downregulated mimats from my list
  dm=intersect(names(trg_dmm2),d_mir)
  ## de_miRNA list values (trg_dmm2) must to be compared to gene names of the good p-values (de_geno)
  hero2=f_match_down(de_geno, trg_dmm2) #using ready made function
  names(hero2)=names(trg_dmm2)

  # mirs up, genes down
  hero3=hero2[um] #ok
  t_mru_d=f_match_down2(hero3, gen_rgl)

```

```

# mirs down
hero4=hero2[dm] #ok
t_mrd_u=f_match_down4(hero4, gen_rgl)

# Mirs up, genes down Fish:
fi=as.vector(matrix(0,nrow=length(t_mru_d)))
for(i in 1:length(t_mru_d))
{hr=fisher.test(matrix(c(
  t_mru_d[i],
  as.vector(sapply(hero3[i],length))-t_mru_d[i],
  length(gen_ok_f)-t_mru_d[i],
  dim(gen_rgl)[1]-(length(gen_ok_f)-t_mru_d[i])), nrow=2)); fi[i]=hr[[1]]);}

# MIRNA UP (=mu), GENE DOWN (=gd) => mugd
mugd=cbind(names(hero3),fi)
colnames(mugd)<-c('Mature miRNA','Fisher test p-value')
mugd=mugd[order(as.numeric(mugd[,2])),]
mugd[,2]=as.numeric(mugd[,2])
mugd[,1]=f_mature_to_mir_vec(mugd[,1])

# mirs down, genes up
fit=as.vector(matrix(0,nrow=length(t_mrd_u)))
for(i in 1:length(t_mrd_u))
{hr=fisher.test(matrix(c(
  t_mrd_u[i],
  as.vector(sapply(hero4[i],length))-t_mrd_u[i],
  length(gen_ok_fu)-t_mrd_u[i],
  dim(gen_rgl)[1]-(length(gen_ok_fu)-t_mrd_u[i])), nrow=2)); fit[i]=hr[[1]];}

# MIRNA DOWN (=md), GENE UP (=gd) => mdgu
mdgu=cbind(names(hero4),fit)
colnames(mdgu)<-c('Mature miRNA','Fisher test p-value')
mdgu=mdgu[order(as.numeric(mdgu[,2])),]
mdgu[,1]=f_mature_to_mir_vec(mdgu[,1])

# p-values below 0.05, and their combined matrix (str(mdgu2): chr):
# and BH correction
ao=as.numeric(mdgu[,2])
ao=cbind(rep('d',length(ao)), mdgu)

ap=as.numeric(mugd[,2])
ap=cbind(rep('u',length(ap)), mugd)

at=rbind(ao,ap)
colnames(at)<-c('regul_miR','miRNA name','Fisher test p-value')

#BH correction:
c_enmir=p.adjust(at[,3], method="BH")
at[,3]=c_enmir
se=as.numeric(at[,3])<0.05
at=at[se,]

write.table(at, "enr_mir_BH_cor.txt", sep="\t", row.names=F)
return(at)
}

```

#Values of confusion matrix extracted

```
f_enr_mir_cmat= function (t_mru_d,hero3,gen_ok_f,gen_rgl,c1) {  
  t=vector("list",length(t_mru_d))  
  names(t)=f_mature_to_mir_vec(names(hero3))  
  a=NULL  
  o=NULL  
  for(i in 1:length(t_mru_d)) (  
    t[[i]] =list(t_mru_d[i], as.vector(sapply(hero3[i],length))-t_mru_d[i],  
      length(gen_ok_f)-t_mru_d[i],dim(gen_rgl)[1]-(length(gen_ok_f)-t_mru_d[i])) )  
  a=do.call(rbind,t)  
  o=match(as.vector(as.vector(c1[,1])),rownames(a))  
  a=a[o,]  
  colnames(a)=c("A","B","C","D")  
  write.table(a,"enr_mir_rgl-uod.txt", row.names=F, sep="\t")  
  return(a)}  
tush_new_miru_cm=f_enr_mir_cmat(t_mru_d,hero3,gen_ok_f,gen_rgl,c1) #ok
```

#Enhanced version from the previous to get all calculated information at once:

```
f_enr_mir_cmate= function (t_mru_d,t_mrd_u,hero3,hero4, gen_ok_f,gen_ok_fu,gen_rgl,c1c) {  
  t=vector("list",length(t_mru_d))  
  names(t)=f_mature_to_mir_vec(names(hero3))  
  ti=vector("list",length(t_mrd_u))  
  names(ti)=f_mature_to_mir_vec(names(hero4))  
  a=NULL  
  b=NULL  
  c=NULL  
  o=NULL  
  oi=NULL  
  for(i in 1:length(t_mru_d)) (  
    t[[i]] =list(t_mru_d[i], as.vector(sapply(hero3[i],length))-t_mru_d[i],  
      length(gen_ok_f)-t_mru_d[i],dim(gen_rgl)[1]-(length(gen_ok_f)-t_mru_d[i])) )  
    for(i in 1:length(t_mrd_u)) (  
      ti[[i]] =list(t_mrd_u[i], as.vector(sapply(hero4[i],length))-t_mrd_u[i],  
        length(gen_ok_fu)-t_mrd_u[i],dim(gen_rgl)[1]-(length(gen_ok_fu)-t_mrd_u[i]))  
    )  
  a=do.call(rbind,t)  
  b=do.call(rbind,ti)  
  o=match(as.vector(as.vector(c1c[,2])),rownames(a))  
  oi=match(as.vector(as.vector(c1c[,2])),rownames(b))  
  a=a[o,]  
  a=a[is.na(rownames(a))!=T,]  
  b=b[oi,]  
  b=b[is.na(rownames(b))!=T,]  
  c=rbind(a,b)  
  c=as.matrix(cbind(as.data.frame(c1c[,3]),c))  
  colnames(c)=c("P-value","A","B","C","D")  
  
  write.table(c,"enr_mir_rgl-uode.txt", row.names=F, sep="\t")  
  return(c)}
```

Clustering for enriched miRs

```
d <- dist(tsem, method = "euclidean") # distance matrix  
fot <- hclust(d, method="average")  
par(ps=15)  
plot(fot, hang = -1) # display dendrogram  
pv2 = pvclust(t(tsem),method.hclust=method,method.dist=distance,nboot = 500) #notice: nboot=500
```



```
pvrect(pv2, alpha = 0.95 ,pv="au", type="geq", max.only=FALSE, border=4) #See the borders and edit in an image program
```

```
##Important function for getting the targets of miR-list, used below;
```

```
f_miR_targs= function (list_mir) {  
  aus=vector("list", length=length(tush_nl))  
  names(aus)=names(tush_nl)  
  cond=c(1:length(tush_nl))[is.na(match(names(tush_nl),list_mir)==T)==F]  
  aus=tush_nl[cond]  
  return(aus)  
}
```

```
# A function to get the names of up/down regulated genes (in a group list) for enriched miRNA (down/up=d/u)
```

```
f_trg_ud_em_dg = function (trg_dmm4,group2) {  
  gr1=f_gsymb_ent(names(group2)) #length(group3)  
  ug1=gen_u[gen_u[,1] %in% gr1,] #dim(ug1)  
  # corresponding down enr_mirs:  
  ug1m=f_mult_g_mir(ug1[,1]) #str(ug1m)  
  ug1m_enr=f_enr_mir_in_lst(trg_dmm4,ug1m)  
  reg_g1_mir=unique(unlist(ug1m_enr))  
  mR_enr_d1=mR_enr_d[(mR_enr_d[,1] %in% reg_g1_mir),]  
  #The function is used to find the targets of these enr_mirs  
  mR_enr_d1_targ=f_miR_targs(mR_enr_d1[,1])  
  d1_trg_un=unique(unlist(mR_enr_d1_targ))  
  #what targets are the same as groups dysregulated genes:  
  g1_ug_denrm=d1_trg_un[d1_trg_un %in% ug1[,1]]  
  #g1_ug_denrm_ks=f_ent_gsymb(g1_ug_denrm) #Gene names, length(g1_ug_denrm_ks)  
  write.table(g1_ug_denrm, "sel_ug_denrm_ks.txt",sep="\t",row.names=F, col.names=F)
```

```
# downregulated names(group1),  
dg1=gen_d[unique(gen_d[,1]) %in% gr1,] #dim(dg1)  
# corresponding up enr_mirs:  
ug1md=f_mult_g_mir(dg1[,1]) #str(ug1md)  
ug1md_enr=f_enr_mir_in_lst(trg_dmm4,ug1md)  
reg_g1_mird=unique(unlist(ug1md_enr))  
mR_enr_u1=mR_enr_u[(mR_enr_u[,1] %in% reg_g1_mird),]  
#The function is used to find the targets of these enr_mirs  
mR_enr_u1_targ=f_miR_targs(mR_enr_u1[,1])  
u1_trg_un=unique(unlist(mR_enr_u1_targ))  
#what targets are the same as groups dysregulated genes:  
g1_dg_uenrm=u1_trg_un[u1_trg_un %in% dg1[,1]]  
#g1_dg_uenrm_ks=f_ent_gsymb(g1_dg_uenrm) #Gene names, length(g1_dg_uenrm_ks)  
write.table(g1_dg_uenrm, "sel_dg_uenrm_ks.txt",sep="\t",row.names=F, col.names=F)
```

```
return(list(g1_ug_denrm_ks,g1_dg_uenrm_ks)) }
```

```
tg1=f_trg_ud_em_dg(trg_dmm4,group1)
```

```
# for normal up/down (when normal mirs up/down)
```

```
f_Trg_ud_M_dg = function (trg_dmm4,de_geno2) {  
  gr1=de_geno2[,1] #length(group3)  
  ug1=gen_u[gen_u[,1] %in% gr1,] #dim(ug1)  
  # corresponding down mirs:  
  ug1m=f_mult_g_mir(as.character(ug1[,1])) #str(ug1m)  
  f_not_z2(ug1m)  
  #ug1m_enr=f_enr_mir_in_lst(trg_dmm4,ug1m)  
  reg_g1_mir=unique(unlist(ug1m))
```

```

mR_d1=miR_d[(miR_d[,1] %in% reg_g1_mir),]
#The function is used to find the targets of these mirs
mR_d1_targ=f_miR_targs(mR_d1[,1])
d1_trg_un=unique(unlist(mR_d1_targ))
#what targets are the same as groups dysregulated genes:
g1_ug_denrm=d1_trg_un[d1_trg_un %in% ug1[,1]]
g1_ug_denrm_ks=f_ent_ksymb(g1_ug_denrm) #Gene names, length(g1_ug_denrm_ks)
write.table(g1_ug_denrm_ks, "sel_ug_dm_ks.txt",sep="\t",row.names=F, col.names=F)

# downregulated names(group1),
dg1=gen_d[unique(gen_d[,1]) %in% gr1,] #dim(dg1)
# corresponding up mirs:
ug1md=f_mult_g_mir(dg1[,1]) #str(ug1m)
#ug1md_enr=f_enr_mir_in_lst(trg_dmm4,ug1md)
reg_g1_mird=unique(unlist(ug1md))
mR_u1=miR_u[(miR_u[,1] %in% reg_g1_mird),]
#The function is used to find the targets of these enr_mirs
mR_enr_u1_targ=f_miR_targs(mR_u1[,1])
u1_trg_un=unique(unlist(mR_enr_u1_targ))
#what targets are the same as groups dysregulated genes:
g1_dg_uenrm=u1_trg_un[u1_trg_un %in% dg1[,1]]
g1_dg_uenrm_ks=f_ent_ksymb(g1_dg_uenrm) #Gene names, length(g1_dg_uenrm_ks)
write.table(g1_dg_uenrm_ks, "sel_dg_um_ks.txt",sep="\t",row.names=F, col.names=F)

return(list(g1_ug_denrm_ks,g1_dg_uenrm_ks)) }
de_gt_udnm=f_Trg_ud_M_dg(trg_dmm4,de_genos) #This takes time if de_genos is long (de genes)

#Obtain enrichment information from internet with previously defined lists of genes
path_to_files =
"C:/Users/Pauli/Documents/thesis/labour/analysis/TCGA/results/COLLECTED_OKS/enrichment/247_go_mf_bp_mirs
_genecodis/"
# For e.g. MFs
files = list.files(path_to_files, pattern="mf_genecodis",all=F)
enr_gomf_cc <- file.path(path_to_files, files)
enr_gomf_cc_ext <- lapply(enr_gomf_cc, sep="\t", quote = "", row.names = NULL,stringsAsFactors = FALSE,read.csv)
#japadapa duuu!!
write.table(enr_gomf_cc_ext, "247_gomf_cc.txt", sep="\t",row.names=F)

##Kegg pwys for up_genes (d_mirs, normal)
path_to_files = "C:/Users/Pauli/Documents/thesis/labour/analysis/TCGA/results/COLLECTED_OKS/enrichment &
clusterings/"
# For e.g. MFs
files = list.files(path_to_files, pattern="_ks",all=F)
sel_gm <- file.path(path_to_files, files)
sel_gmj <- lapply(sel_gm, sep="\t", quote = "", row.names = NULL,stringsAsFactors = FALSE,read.csv)
sel_ug_dm_ks=do.call(rbind,sel_gmj[2])
sud=sel_ug_dm_ks[,1]
write.table(sud, "sud.txt",sep="\t",row.names=F)
sud=read.table(file="sud.txt")
sel_ug_dm_ks=as.matrix(f_ksymb_ent(sud[,1]))
my_pwys_ud=pwy_func5(pathway.list2, sel_ug_dm_ks, e11_l2) #no result is ok result

#Get certain information from a list quickly to same source (path.files):
# 247_U_gen_D_enrmirs BP
path_to_files = "C:/Users/Pauli/Documents/thesis/labour/analysis/TCGA/results/COLLECTED_OKS/enrichment &
clusterings/Genecodis/List_247_min_2_enrmir/Enr mirs/247_U_gen_D_enrmirs/"
# For e.g. MFs

```

```

files = list.files(path_to_files, pattern="bp",all=F)

#function to get info from genecodis files: bp,mf,kegg
f_prs_1_fil= function(path_to_files,files) {

udem_bp <- file.path(path_to_files, files)
udem_bp_l <- lapply(udem_bp, sep="\t", quote = "", row.names = NULL,stringsAsFactors = FALSE,read.csv)
#dim(udem_bp_l[[1]])
cond2=udem_bp_l[[1]][,9]<0.01 & !is.na(udem_bp_l[[1]][,9]<0.05)
udem_bp_l2=udem_bp_l[[1]][cond2,]
udem_bp_l3=udem_bp_l2[order(udem_bp_l2[,9]),]
#dim(udem_bp_l3)
t_fap=function(a,b,c,d) ((a/b)/(c/d))
v1=NULL
for (i in 1:dim(udem_bp_l3)[1]) {
v1[[i]]=t_fap(udem_bp_l2[i,c(4)],udem_bp_l2[i,c(5)],udem_bp_l2[i,c(6)],udem_bp_l2[i,c(7)]) )
udem_bp_l4=cbind(udem_bp_l3,v1)
udem_bp_l4=udem_bp_l4[,c(2,3,9,4,11)]
colnames(udem_bp_l4)=c("GO:ID","GO:Term","Q-value","Sig.genes","Odds ratio")
x <- udem_bp_l4[,2]
y <- gsub(" (BP)","",x,fixed=T)
udem_bp_l4[,2] = y
udem_bp_l4[,3]=round(udem_bp_l4[,3],digits=4)
udem_bp_l4[,5]=round(udem_bp_l4[,5],digits=1)
setwd(path_to_files)
write.table(udem_bp_l4, "gc_CHECK.txt", sep="\t",row.names=F)
setwd("C:/Users/Pauli/Documents/thesis/labour/analysis/TCGA/results")
return(udem_bp_l4)
}

```

```

#function to get info from genecodis files: bp,mf,kegg
t_prs_1_fil= f_prs_1_fil(path_to_files,files) #ok

```

Step 15

Linear modelling of miRNAs

Step 15a

Enhancing the data before the linear modelling scheme:

miR normalisation

```
testia=justvsna(as.matrix(dtt2c)) #dtt2c=dtt2b
```

miR z-scoring

```
data_zs=zNorm(testia)
```

Check the distribution with the box_plot

```
boxplot(log2(data_zs[,100:125])) # about ok
```

exp log-2 (Vijay also did log-2)

```
e11_l2=log2(e11+1)
```

e11_l2 z-scoring

```
e11_zs=zNorm(e11_l2)
```

Check the distribution with the box_plot

```
boxplot(log2(e11_zs[,1:35])) #ok
```

A function for finding the corresponding gene symbol to Entrez name

```
f_ent_gsymb = function (list) {
```

```
mm=match( as.matrix(names(list)),as.character(gsymb2[,1]))
```

```
return (gsymb2[mm,2])}
```

```
#### Step 15b ####
```

```
# Defining the lists and mirs for linear modelling
```

```
# list 1:
```

```
list11a=unlist(trg_dmm4) #all genes: length(list11a): [1] 8868
```

```
list11=unique(unlist(trg_dmm4)) #unique genes: length(list11): [1] ~4896
```

```
# To obtain all the miRNAs that the target gene (x) of interest (from tush_nl list) has:
```

```
f_mir_limoo = function(tush_nl, x) {  
  x=as.character(x)  
  tmd=NULL #note that tush_nl is a list  
  #Matching with %in% command  
  for (i in 1:length(tush_nl)) {  
    if (sum ( (as.character(unlist(tush_nl[i])) %in% x)==T ) > 0 )  
      tmd[i]=names(tush_nl)[i] }  
  return(tmd[!is.na(tmd)]) }
```

```
# Finding miRNAs for these genes
```

```
f_mult_g_mir= function(list) {  
  #list is a gene list  
  gene_nam=vector(mode="list", length=length(list))  
  mirs=vector(mode="list", length=length(list))  
  for (i in 1:length(list)) {  
    gene_nam[[i]] = as.vector((list)[i])  
    mirs[[i]] = f_mir_limoo(tush_nl, gene_nam[[i]]) }  
  return(mirs) } #or, using matrix
```

```
f_mult_g_mir2= function(matrix) {  
  #list is a gene list  
  gene_nam=vector(mode="list", dim(matrix)[1])  
  mirs=vector(mode="list", dim(matrix)[1])  
  for (i in 1:dim(matrix)[1]) {  
    gene_nam[[i]] = matrix[i,1]  
    mirs[[i]] = f_mir_limoo(tush_nl, gene_nam[[i]]) }  
  return(mirs) }
```

```
mirs_list1=f_mult_g_mir(list11)
```

```
#find the number of enriched targets in miRNA list
```

```
f_enr_mir_in_lst=function(trg_dmm4,mirs_list3r) {  
  x=NULL  
  for (i in 1:length(mirs_list3r)) {  
    x[[i]]=mirs_list3r[[i]][unlist(mirs_list3r[[i]]) %in% unlist(names(trg_dmm4))]  
  }  
  return(x)  
}
```

```
en_mir_deva2t=f_enr_mir_in_lst(trg_dmm4,mirs_list3r)
```

```
# list2:
```

```
int_gen23=read.table("interesting_genes_zdk_pt_4214.txt",sep="\t")
```

```
int_gen23=unique(int_gen23[,1])
```

```
# Then you should check:
```

```
# All the (miRNA target) gene names that you have in your list of interest /two, "stp3"/:
```

```
f_mir_limo2 = function(targets_ok3, stp3) {  
  tmd=NULL # stp must be in vector have same names as the genes (so double)  
  tmd=stp3[stp3 %in% as.character(unlist(targets_ok3))]  
  return(tmd) }
```

#matching gene symbols to my Entrez list (the list have slightly less values than in id_converter database) (I should use some other method then)

```
f_g symb_ent= function (list) {  
  mm=match( as.matrix((list)),as.character(g symb2[,2]))  
  return (g symb2[mm,1])} #notice that some gene symbols do not have EntrezNames, check also:  
http://idconverter.bioinfo.cnio.es/
```

int_gen23=f_g symb_ent(int_gen23) # some ent symbols where not available, and they where collected from internet inividually:

```
ig23l=c("7490","383","5009","4842","2271","189","5563","2645","5163","5313","5106")
```

Replace NA values with other values:

```
y <- which(is.na(int_gen23)==TRUE) # This is how you get the indexes
```

```
int_gen23[y]=ig23l
```

```
mirs_list3=f_mult_g_mir(int_gen23)
```

#It seems that I need to check from data_zs_t that I got all the miRs also

```
f_c_ok_mirs=function (data_zs_t,lists_s_mirs) {  
  mirs_ig=NULL  
  for (i in 1:length(lists_s_mirs)) (  
    mirs_ig[[i]]=lists_s_mirs[[i]][which (lists_s_mirs[[i]] %in% rownames(data_zs_t))] )  
  return(mirs_ig)  
}
```

```
mirs_list3r=f_c_ok_mirs(data_zs_t,lists_s_mirs)
```

Check how many miRs you have in your miRNA list:

```
f_enr_mir_nro=function(mirs_list3r) {  
  x=NULL  
  for (i in 1:length(mirs_list3r)) (  
    x[[i]]=length(mirs_list3r[[i]]) )  
  return(x)  
}
```

```
mirs_list3r_nro=f_enr_mir_nro(mirs_list3r)
```

#Removing nas

```
f_not_z2=function(list_mirs) {  
  h=NULL  
  for (i in 1:length(list_mirs)) (  
    if ( length(list_mirs[[i]]) != 0) h=c(h,i) )  
  return(list_mirs[h])  
}
```

```
mirs_list3r=f_not_z2(mirs_list3r_nro)
```

#Adding and parsing some elements to a list... (see tikka_complete_extras.r):

```
oks4=rownames(e11_zs)[rownames(e11_zs) %in% as.vector(int_gen23y)]
```

Then the previous ordering might be done as:

```
mir_oks4=vector("list", length(oks4))
```

```
mir_oks4=f_mult_g_mir(oks4)
```

```
mir_oks4=f_c_ok_mirs(data_zs_t,mir_oks4)
```

```
names(mir_oks4)=names(oks4)
```

```
mir_oks4=f_not_z2(mir_oks4)
```

#ordering of the list

```
pito=f_enr_mir_nro(mir_oks4)
```

```
mir_oks4_ord = (mir_oks4[rev(order(pito))])
```

```
list_2m=cbind(names(mir_oks4_ord), pito)
```

```
#write.table(list_2m,"list2m.txt",row.names=F,col.names=F,sep="\t")
```

```

#list 3:
#Finding the gene list from overlapping miRNA's information:
f_overlap2 = function (x) {
  ai=unique(unlist(x)) #4321
  overlap_mat2 = matrix(nrow=length(x),ncol=length(ai))
  overlap_mat2 = "[<-"(overlap_mat2,value=0) # all values are zero
  rownames(overlap_mat2)=names(x)
  colnames(overlap_mat2)=ai
  for (j in 1:length(x)) {
    for (k in 1:length(x[[j]])) {
      for (i in 1:dim(overlap_mat2)[2]) {
        if (colnames(overlap_mat2)[i]==as.character(x[[j]][k]))
          (overlap_mat2[j,i]=1) )))
  return(overlap_mat2)}
overlap_mat3=f_overlap2(trg_dmm3)
# Or if you have done this matrix in another computer then you may download it:
overlap_mat3=read.table("overlap_mat.txt", sep="\t", col.names=T)

# One does need to limit the amount of genes so for this a list of most overlapping genes is done:
over_sum=apply(overlap_mat2,2,sum)
over_sum_ok=over_sum > 1
over_lap=over_sum[over_sum_ok]
over_ok=as.matrix(over_lap[rev(order(over_lap))])
## are all the names in the over_ok at e11_zs?
dim(e11_zs[(rownames(e11_zs) %in% rownames(over_ok)),]) # vs. dim(over_ok)
## aha! they do not match! So
over_ok_mod = over_ok[(rownames(over_ok) %in% rownames(e11_zs)),1]
list=over_ok_mod[1:150]

# This matrix only for differentially expressed genes:
f_overlap3 = function (de_genos, trg_dmm4) {
  ai=unique(unlist(de_genos[,1])) #1963
  overlap_mati2 = matrix(nrow=length(trg_dmm4),ncol=length(ai))
  overlap_mati2 = "[<-"(overlap_mati2,value=0) # all values are zero dim(overlap_mati2)
  rownames(overlap_mati2)=names(trg_dmm4)
  colnames(overlap_mati2)=ai
  for (j in 1:length(trg_dmm4)) {
    for (k in 1:length(trg_dmm4[[j]])) {
      for (i in 1:dim(overlap_mati2)[2]) {
        if (colnames(overlap_mati2)[i]==as.character(trg_dmm4[[j]][k]))
          (overlap_mati2[j,i]=1) )))
  return(overlap_mati2)}
overlap_mati3=f_overlap3(de_genos, trg_dmm4)

#Renaming the list (notice that list name should be "3" and not list22...
list22=list
names(list22)=names(ibtaai)
list22=as.matrix(list22)
write.table(list22,"list22.txt",sep="\t", col.names=F) #ok

Combining the lists:
#1: list_11ii
#2: list2ok_nam
#3: list

#For list one (list_11ii) and three, Selecting only the ones that are not in list three
m_l31=match(names(list),rownames(list_11ii))

```

```

l1_red=list_11ii[-m_l31,]
lir1=rownames(l1_red[1:50,])
lir11=rownames(l1_red[1:73,])

#For list 2 (list2ok_nam) and three
m_l21=match(names(list2ok_nam),names(list)) #voitaneen lisätä tällaisenaan
# checking: table(names(list) %in% as.character(list2ok_nam)) #, ok

#All 3 lists (except list nr1 as such in order: 1,2,3)
list_s_tot=c(lir1,as.character(list2ok_nam),names(list))
lists_s_tot=c(lir11,as.character(list2ok_nam),names(list)) #Better to have slightly bigger list and round number in
total 250
#you need names for your list for the function
names(lists_s_tot)=lists_s_tot
#Check that you have every gene:
table(rownames(e11_zs_t) %in% lists_s_tot)
match(lists_s_tot,rownames(e11_zs_t))

#Checking again
lists_s_mirs=f_mult_g_mir2(lists_s_tot)
names(lists_s_mirs)=lists_s_tot

#This is what you need:
lmt_m_ok=f_c_ok_mirs(data_zs_t,lists_s_mirs)
lmt_m_nro=f_enr_mir_nro(lmt_m_ok)
#There must be at least one miRNA in the model (I corrected the f_limo_restr2 accordingly)
rm_lt=c(1:250)[is.na(match(lmt_m_nro,0))==T]==F]
lists_s_tot=lists_s_tot[-rm_lt] #Removing zero miRNA value genes from list
lir11_add=rownames(l1_red[c(75:76),]) # adding nonzero miRNA genes
lists_s_tot=c(lists_s_tot,lir11_add) #combining the reduced list and the newly added genes
names(lists_s_tot)[249:250]=lir11_add #naming all the entries

#So finally you may use your functions (and there should not be zeros any more):
lmt_m_ok=f_c_ok_mirs(data_zs_t,lists_s_mirs)
lmt_m_nro=f_enr_mir_nro(lmt_m_ok) #Alternative way: lists_mir_lengths=sapply(lmt_m_ok,length), and
median(lists_mir_lengths)

# The big driving in HKI computers, first combining list 1-3 (as such):
list_big=c(rownames(l1_red),as.character(list2ok_nam),names(list))

### Step 15c ###
# Clustering and Enrichment of one of the lists (list3, previously: (just top 150) list):
# constructing the function for hammington distance matrix
f_ham_dis= function (ibt2) {
  #ibt is the gene (and miR) information in list of lists
  abcd <- vector("list", length(ibt2))
  abcd=NULL
  for (i in 1:length(ibt2)) {
    abcd[[i]]=c(rownames(ibt2[[i]]))
  }
  abcd=unlist(abcd)
  abcd=as.matrix(unique(abcd))
  abcd=abcd[(abcd[,1]!="beta_0"),1] # jee

  #the ham matrix constructed
  gen_mir_ham=matrix(0,nrow=length(ibt2),ncol=length(abcd))
  rownames(gen_mir_ham)=names(ibt2)
  colnames(gen_mir_ham)=abcd

```

```

for (i in 1:length(ibt2)) (
for (j in 1:(length(ibt2[[i]][[2]]))) (
if (ibt2[[i]][[2]][j,1]>0) (gen_mir_ham[i,(abcd %in%
as.matrix(rownames(ibt2[[i]][[2]]))[-match("beta_0",rownames(ibt2[[i]][[2]])),1][j])=1) ) )
return(gen_mir_ham) }

```

```

# The distanced maybe thus calculated:
f_calc_dist = function (gen_mir_ham2) {
n <- nrow(gen_mir_ham2)
m <- matrix(nrow=n, ncol=n)
rownames(m)=rownames(gen_mir_ham2)
colnames(m)=rownames(gen_mir_ham2)
for(i in seq_len(n - 1))
  for(j in seq(i, n))
    m[j, i] <- m[i, j] <- sum(gen_mir_ham2[i,] != gen_mir_ham2[j,])
return(m)
}

```

```

# Clusterings of Hammington distances of top 150 miRNA overlapping genes correlations
pdf('clust_all.pdf')
method="complete"
distance="euclidean"
test_clust=hclust(dist(mo, method=distance),method=method)
par(ps=3)
plot(test_clust,par) # you may also save as such in the Rstudio as pdf
dev.off()

```

```

# Information for cluster groupings:
# Ward Hierarchical Clustering was found to be good one for clustering:
d <- dist(moi2, method = "euclidean") # distance matrix
rownames(moi2)
fit <- hclust(d, method="ward")
par(ps=3)
plot(fit) # display dendrogram
groups <- cutree(fit, k=9) # cut tree into 9 clusters
# draw dendrogram with red borders around the 9 clusters
rect.hclust(fit, k=9, border="red")
#Group selection from the image is a delicate process:
group1=groups[(groups==1)] # ok
(groups)["BACH2"] ##in group 4
group2=groups[(groups==4)]
(groups)["KLF12"] #in group 2
(groups)["CELF2"] # in group 6
gr3= (groups==2) | (groups==6) # combining two groups (was needed according to image)
group3=groups[gr3]
(groups)["TRPS1"] #in group 5
group4=groups[(groups==5)]
(groups)["IKZF4"] #in group 3
group5=groups[(groups==3)]
(groups)["SLC1A2"] #in group 8
group6=groups[(groups==8)]
(groups)["PHACTR2"] #in group 9
group7=groups[(groups==9)]
(groups)["RC3H1"] #in group 9
group8=groups[(groups==7)]

```



```

#Collecting information of total list (here list 247):
path_to_files =
"C:/Users/Pauli/Documents/thesis/labour/analysis/TCGA/results/COLLECTED_OKS/enrichment/247_kegg_pwys_gen
ecodis/extra/"
# Fore every miR
files = list.files(path_to_files, pattern="_miRs",all=F)
enr_k_cc <- file.path(path_to_files, files)
enr_k_cc_ext <- lapply(enr_k_cc, sep="\t", quote = "", row.names = NULL,stringsAsFactors = FALSE,read.csv)
#japadapa duuu!!
write.table(enr_k_cc_ext, "247_kegg_cc.txt", sep="\t",row.names=F)

mir_info2[[1]][1:4, 1:7]
a=sub(".txt", "", files, perl=TRUE) #jee
names(mir_inf)=a
names(mir_inf2)=a

#Better to have 8 tables where extract information:
grp_tot_nam2=list(names(group1),
names(group2),names(group3),names(group4),names(group5),names(group6),names(group7),names(group8))
names(grp_tot_nam2)=c("group1", "group2","group3", "group4","group5","group6","group7","group8")
lapply(names(grp_tot_nam2),function(x, grp_tot_nam2)
write.table(grp_tot_nam2[[x]], paste(x, ".txt", sep = ""),col.names=FALSE, row.names=FALSE, sep="\t",
quote=FALSE),grp_tot_nam2)

#It is maybe better to do this 8 times for every group the same info:
path_to_files = "C:/Users/Pauli/Documents/thesis/labour/analysis/TCGA/results/cluster_150_enr/"
g_inf=vector("list",8)
files=vector("list",8)
gi=vector("list",8)
pattern=vector("list",8)

# Somehow one loop solution did not work:
for (i in 1:8) (
  pattern[[i]]=paste("_g",i, sep="") )
for (i in 1:8) (
  files[[i]] = list.files(path_to_files, pattern[[i]],all=F) )
for (i in 1:8) (
  gi[[i]] <- file.path(path_to_files, files[[i]]) )
for (i in 1:8) (
  g_inf[[i]] <- lapply(gi[[i]], sep="\t", quote = "", row.names = NULL,stringsAsFactors = FALSE,read.csv) )
names(g_inf)=c("group1", "group2","group3", "group4","group5","group6","group7","group8")

#Enrichment information for these groups from: http://amp.pharm.mssm.edu/Enrichr/index.html
#Selecting certain things for observations:
# Find out which list factor is which:
g_inf[[1]][[1]][1:4, 1:7] ## cancer cell en
g_inf[[1]][[2]][1:4, 1:7] ## go_bp
g_inf[[1]][[3]][1:4, 1:7] ## go_mol_func
g_inf[[1]][[4]][1:4, 1:7] ## kegg
g_inf[[1]][[5]][1:4, 1:7] ## mir

#Now changing the last module you may see all of the values that you want:
f_wanted_gr= function (g_inf) {
  abc=NULL
  for (i in 1:8)
    (abc[[i]]=g_inf[[i]][[4]][1:5, c(1,2,4)])
  return(abc) }

```

```

#Defining wanted:
sel_grp=f_wanted_gr(g_inf)
names(sel_grp)=names(g_inf)

#Clustering of samples according to the miRNAs:
#First some filtering: (#testia2 is justvsrn normalized data, after filtering of zero expressed miRNAs: dim(testia2) 332
567)
# Variance:
mean_mir_expr=apply(testia2,1,mean)
quantile(mean_mir_expr)
m_val=mean_mir_expr[8.270042<=mean_mir_expr] #length(m_val)
m_val_expr=testia2[names(m_val),] # dim(m_val_expr): 83 567

#Next hierachichal clustering of miRNAs/mRNAs
# miRs SAMPLES!!:
install.packages('sparcl')
library(sparcl)
dim(m_val_expr) #83miRNAs
d <- dist(t(m_val_expr), method = "euclidean") # distance matrix
fit <- hclust(d, method="ward")
method = "ward"
distance = "euclidean"
pv2 = pvclust((m_val_expr),method.hclust=method,method.dist=distance,nboot = 10)

# e_n
yen=c(1:dim(e11_zs)[2])[is.na(match(colnames(e11_zs),e_n))==T]==F]
yen1=c(1:dim(e11_zs)[2])[yen]=1
# e_t
load("e_t.rdata")
yet=c(1:dim(e11_zs)[2])[is.na(match(colnames(e11_zs),e_t))==T]==F]

# all -e_n - e_t
yetn=c(yen,yet)
ym=c(1:dim(e11_zs)[2])[is.na(match(c(1:dim(e11_zs)[2]),yetn))==T)==F]
yea=c(1:dim(e11_zs)[2])[-ym]
y=c(yea,yet,yen) #length(y)
# Perform hierarchical clustering
hc <- hclust(dist(x),method="complete")

# Plot
method = "average" #?
distance = "euclidean"
pv2 = pvclust((m_val_expr),method.hclust=method,method.dist=distance,nboot = 10)
par(ps=22)
y2 = matrix(ncol=1,nrow=567)
y2 = "[<-"(y2,value=0) #
y2[yen,]=1
y2[yet,]=2
y2[yea,]=3
ColorDendrogram(fit,y=y2,main="",branchlength=150)#,label=F)
pvrect(pv2, alpha = 0.99,pv="au", type="geq", max.only=FALSE, border=8)#, col="WHITE")

#now for genes:
mean_gen_expr=apply(e11_l2,1,mean)
quantile(mean_gen_expr)
ge_val=mean_gen_expr[9.9610451<=mean_gen_expr] #length(ge_val)
ge_val_expr=e11_l2[names(ge_val),] # dim(ge_val_expr), too high: [1] 4115 567

```

```

#Also variance:
GeneVar=apply(ge_val_expr,1,var)
GVM=mean(GeneVar)
ge_val_expr2 <- ge_val_expr [ - which (GeneVar <= 0.8) , ] #dim(ge_val_expr2), 710->ok,

#the rest similarly

#### Step 15d ####
# Linear modelling

# Formulating the problem for one (or more genes)
# Samples (i.e. error values): i, # miRNAs: m
# RHS (right-hand-side, ie. equations, or rows) : n (= i x 2)
# LHS (left-hand-side, i.e. columns): i + m + 1

#First, finding restriction value for miRs
quantile(lmt_m_nro,probs = seq(0, 1, 0.05)) #Ok, let the miRNA in this case be 20% quantile
plot(quantile(lmt_m_nro,probs = seq(0, 1, 0.05)))
# According to mir_amount decide the TimeLimit:
mir_amount=22 # 20% quantile
#Estimate the mean time needed to calculate values between 20-90% quantile, by running the basic function once
with miRNA value above mir_amount:
TL=333 #TL= TimeLimit
TL=333 #TL= TimeLimit

# This is the enhanced function for making linear modelling using also a restriction analysis (and training groups if
needed):
f_limo_restr2 = function (e11_zs_t, e11_zs_t2, data_zs_t, data_zs_t2, gene_nam, mirs, mir_amount, TL) {
  mae=max(e11_zs_t)
  mid=min(abs(data_zs_t))
  mult=ifelse(mae/mid>300,300,mae/mid)
  e_lim=ifelse(round(max(e11_zs_t)-mult*(min(data_zs_t)))>25,25,round(max(e11_zs_t)-mult*(min(data_zs_t))))

  mir_b_max = mult
  #mir_amount = 153 #ok, with correction of data_zs finding right miRs cor~-0.1124, and list1,1/1000/227: -0.047,
now without bounds: 0.2058, but it uses only one miR?
  #ok, now the correlation is 0.2 (after just bounding miRs with equations, but so many zeros...), with 17 mirs: no
correlation. The earlier best celf2: cor:na

  #miRNA list (checking what I can find)
  ee=length(rownames(data_zs_t)[rownames(data_zs_t) %in% as.vector(unlist(mirs))])
  mir_val = if (ee>1) (t(data_zs_t[(rownames(data_zs_t) %in% as.vector(unlist(mirs))),])) else
(matrix(t(data_zs_t[(rownames(data_zs_t) %in% as.vector(unlist(mirs))),])))

  colnames(mir_val) = if (dim(mir_val)[2]<=1)
rownames(data_zs)[c(1:length(rownames(data_zs)))[is.na(match(rownames(data_zs),unlist(mirs))==T)==F]] else
colnames(mir_val)

#dim(mir_val)
#colnames(mir_val)
#mir_val[1:5,]

#Maximum (or minimum) miRNA expression:
m=mir_b_max
#The amount of miRNA you want to use:
w=mir_amount

```

```
w=ifelse(w>dim(mir_val)[2], dim(mir_val)[2],w)
```

```
nro1 = dim(mir_val)[1]
nro2 = dim(mir_val)[2]
nro3 = (nro2+1)
```

```
#rows:
```

```
#rhs = (nro1*2+nro2*4+1)
```

```
rhs = (nro1*2+nro2*2+1)
```

```
#columns:
```

```
lhs = (1+nro2+nro1+nro2)
```

```
#The (row size ) limit of the "normal matrix" (i.e. after that comes beta restrion constrain equations area):
```

```
#oe=rhs-nro2*4-1
```

```
oe=rhs-nro2*2-1
```

```
# The (column size ) limit of the "normal matrix"
```

```
# (after that the binary variables):
```

```
ae=lhs-nro2
```

```
# The constraint matrix should then look like:
```

```
A_mat_f = matrix(0,nrow=rhs, ncol=lhs, byrow=T)
```

```
#dim(A_mat_f)
```

```
#Constraints for betas (part 2a) matrix. Building column by column
```

```
#zero (column) for constants
```

```
A_mat_cont1 = matrix(0,nrow=2*nro2, ncol=1, byrow=T)
```

```
# values for miRNAs
```

```
# Constraints for betas (part 2b)
```

```
be=rbind(-m,m)
```

```
Lt=c()
```

```
for(i in 1:(nro2))(Lt[i]=list(be)) # The "be" needs to be in a list
```

```
b2=data.matrix(bdiag(Lt))
```

```
#dim(b2)
```

```
#b2[1:,1:4]
```

```
A_mat_cont2 = matrix(0,nrow=2*nro2, ncol=nro2, byrow=T)
```

```
# I need looping:
```

```
for (j in 1:dim(A_mat_cont2)[2]) (
```

```
  A_mat_cont2[(b2[,j]!=F),j]=rbind(1,1) )
```

```
#dim(A_mat_cont2)
```

```
#A_mat_cont2[1:4,1:4]
```

```
#zero (column)s for errors
```

```
A_mat_cont3 = matrix(0,nrow=2*nro2, ncol=nro1, byrow=T)
```

```
#dim(A_mat_cont2)
```

```
#putting part 2a and 2b together:
```

```
A_mat_cont=cbind(A_mat_cont1,A_mat_cont2,A_mat_cont3,b2)
```

```
#dim(A_mat_cont)
```

```
#A_mat_cont[1:6,570:574]
```

```
#bounds (part3) #not needed => needed (17.2.2014), maybe not needed (20.2.2014)
```

```
A_mat_contb = matrix(0,nrow=2*nro2, ncol=lhs, byrow=T)
```

```
A_mat_cont2a = matrix(0,nrow=2*nro2, ncol=nro2, byrow=T)
```

```
# I need looping:
```

```
for (j in 1:dim(A_mat_cont2a)[2]) (
```

```
  A_mat_cont2a[(b2[,j]!=F),j]=rbind(1,1) )
```

```

A_mat_contb[,2:(nro3)]=A_mat_cont2a
#dim(A_mat_contb)
#A_mat_contb[1:6,570:574] #ok

#restriction of x_s (part4)
A_mat_res = matrix(0,nrow=1, ncol=dim(A_mat_cont)[2], byrow=T)
A_mat_res[, (ae+1):lhs]=1

#the additional
#old:
#A_mat_more=rbind(A_mat_cont,A_mat_contb,A_mat_res)
#dim(A_mat_more)
#"new old:
A_mat_more=rbind(A_mat_cont,A_mat_res)

# making of vector 1)
v_e=rep(c(1,-1), nro2)
#v_i=matrix(1, nrow=nro1)
#dim(v_i)
#v_e = as.vector(rbind(v_i,-v_i))
#length(v_e)
#miRNA list (check: sortmiro)
#mir_val=t(data_zs_t[(rownames(data_zs_t) %in% mirs),])
#dim(mir_val)
# Make a matrix (mir_neg_pos) that has every other row
# as a negative value of the previous matrix (mir_val):
#mir_neg_pos=matrix(nrow=(oe), ncol=nro2)
#for (i in 1:(oe)) (
# if (i %% 2 !=0) (mir_neg_pos[i,]=mir_val[(1+(i-1)/2),])
# else (mir_neg_pos[i,]=-mir_val[(i/2),]) ) # ok
#if (sum(is.na(mir_neg_pos)!=F)>0) (mir_neg_pos[is.na(mir_neg_pos)] <- 0)

mir_neg_pos=matrix(nrow=oe, ncol=nro2)
#for (i in 1:(oe)) (
# mir_neg_pos[i,] = ifelse (i %% 2 !=0, mir_val[(1+(i-1)/2)], -mir_val[(i/2)]) )

mirval2=mir_val[rep(1:nrow(mir_val),each=2),]
mirval2=rep(c(1,-1),dim(mir_val)[1])*mirval2

mir_neg_pos=mirval2

#mir_neg_pos[1:6,1:6]
#dim(mir_neg_pos)
## The zero list (b):
a=-rbind(1,1) # this is the double. And before the loop an empty list is created:
Lst=c()
for(i in 1:nro1)(Lst[i]=list(a)) # The "a" needs to be in a list
b=data.matrix(bdiag(Lst)) # this is how you use your list to construct this matrix.

#column 1
A_mat_f[1:length(v_e),1]=v_e
# columns 2:(nro2+1)
A_mat_f[1:(oe),2:(nro3)]=mir_neg_pos
# lhs = (1+nro2+nro1+nro2)
A_mat_f[1:(oe),(nro3+1):(ae)]=b #see comments about the type of matrix above
#dim(b)
#dim(A_mat_f[1:(oe),(nro3+1):(ae)])

```

```

# Final addition row-wise:
A_mat_f[(oe+1):rhs,]=A_mat_more #finally
#dim(A_mat_f[(oe+1):rhs,])

# Making the model list values:
model_f <- list()
# Making of the objective function
obj_f <- vector(mode = "numeric", length = lhs)
obj_f = "[<-"(obj_f,1:lhs,value=0)
obj_f = "[<-"(obj_f,((nro3+1):ae),value=1)

model_f$obj = obj_f
model_f$model_sense <- "min"
#model_f$sense <- c(rep('<=', oe), rep(rbind('<=', '>='),2*nro2), '=')
model_f$sense <- c(rep('<=', oe), rep(rbind('<=', '>='),nro2), '=')
#length(model_f$sense)
model_f$lb <- c(-m,rep(-m,nro2),rep(-e_lim,nro1),rep(0,nro2))
model_f$sub <- c(m,rep(m,nro2),rep(e_lim,nro1),rep(1,nro2))
#length(model_f$sub)
model_f$vtype <- c(rep('C',ae), rep('B',nro2))
#length(model_f$vtype)
#lhs
# Making the Right side of equations:
#model_f$rhs <- as.vector(c(rbind((e11_zs_t[gene_nam,]),(-e11_zs_t[gene_nam,])),(rep(0,2*nro2)),(rep(c(m,-
m),nro2)),w))
model_f$rhs <- as.vector(c(rbind((e11_zs_t[gene_nam,]),(-e11_zs_t[gene_nam,])),(rep(0,2*nro2)),w))
#length(model_f$rhs)

# When you got A_mat_f ok, then this works (if row and column numbers are ok)
model_f$A <- A_mat_f
#mir_val[1:4,1:4]
#e11_zs_t[gene_nam,]

#write.table(A_mat_f, "amf.txt", sep="\t")
#params_f <- list(MIPGap=0.01,TIME_LIMIT=12000, ResultFile='model_f.lp') #If doing restriction analysis
#params_f <- list(MIPGap=0.05,TIME_LIMIT=300, ResultFile='model_f.lp') #If doing restriction analysis (small)
params_f <- list(TIME_LIMIT=TL, ResultFile='model_f.lp')
#params_f <- list(Method=-1,ResultFile='model_f.lp') #normally
result_f <- gurobi(model_f,params_f)
#result_f <- gurobi(params_f)

#more info
limo_mir=cbind(result_f$x[1:(nro2+1)]) #ok, I think this is working now :)
#max(result_f$x[nro3:ae])
colnames(limo_mir)=c('lin.mod.coefficient')
rownames(limo_mir)=c("beta_0",colnames(mir_val))

# Ordering the miRs according to coefficient
limo_mir2=limo_mir
limo_mir2=as.matrix(limo_mir[rev(order(limo_mir)),])
colnames(limo_mir2)=c('lin.mod.coefficient')
# Finding overlapping ones to the all miR list
en_mir=as.matrix(limo_mir2[rownames(limo_mir2) %in% names(trg_dmm4),])
colnames(en_mir)=c('lin.mod.coefficient')

#Multiplication:

```

```

b_m_i_f=data_zs_t2[(rownames(data_zs_t2) %in% rownames(limo_mir)),]
#table(rownames(data_zs_t2) %in% c("hsa-miR-30d-5p"))
limo_miri=limo_mir[rownames(limo_mir) %in% rownames(data_zs_t2)]
#dim(b_m_i_f)
#length(limo_miri)
ert=dim(data_zs_t2)[2]
# matrix (or vector) multiplication is needed :)
g_pred_f_i=(b_m_i_f*as.numeric(limo_miri))
auo=if (length(g_pred_f_i) > ert) dim(g_pred_f_i)[2] else length(g_pred_f_i)
#dim(g_pred_f_i)

g_pred_fo=c(1:auo)
#length(g_pred_fo)
#dim(g_pred_f_i)[2]

g_pred_fo = "[<-"(g_pred_fo,1:auo,value=result_f$x[1])

g_pred_f=rbind(g_pred_fo,g_pred_f_i)
if (sum(is.na(g_pred_f)!=F)>0) (g_pred_f[is.na(g_pred_f)] <- 0)
#g_pred_f[1:3,1:3]

# but the equations continue...:
g_pred_val_f=NULL
for (j in 1:dim(g_pred_f)[2])
  (g_pred_val_f[j]=sum(g_pred_f[,j]))
#g_pred_val_f[1:3]
#length(g_pred_val_f)

# And the correlation
y=as.numeric(e11_zs_t2[gene_nam,])
x=as.numeric(g_pred_val_f)
cor=cor(y, x, method = "pearson") #different values if I

return(list(result_f, limo_mir2, en_mir, cor))
}

#Ok, lets make a function for all:
fi_loop2=function ( e11_zs_t, e11_zs_t2, data_zs_t, data_zs_t2,mir_amount, list) {
  library("gurobi")
  library("Matrix")

  # Making the initial vector:
  mg_list=vector(mode="list", length=length(list))
  names(mg_list) <- f_ent_g symb(list)

  gene_nam=vector(mode="list", length=length(list))
  mirs=vector(mode="list", length=length(list))

  mi=mir_amount

  for (i in 1:length(list)) {
    gene_nam[[i]] = as.vector((list)[i])
    mirs[[i]] = f_mir_limoo(tush_nl, as.vector(unlist(gene_nam[[i]])))
    mg_list[[i]] = f_limo_restr2(e11_zs_t, e11_zs_t2,data_zs_t, data_zs_t2,as.vector(unlist(gene_nam[[i]])),
as.vector(unlist(mirs[[i]])),mi) }
  #here are the values:
  return(mg_list)
}

```

```

}
#Test:
l_tot_s_test1=fi_loop2( e11_zs_t, e11_zs_t2, data_zs_t, data_zs_t2,mir_amount, lists_s_tot)
ltot1_req=f_requested_limooo(l_tot_s_test1,lists_s_tot) #See description later

#Get the values for your validation:
f_train_validate2= function (lim_anno, e11_zs, data_zs, train1, train2, validate, list) {
  # The most important:
  # Selecting appropriate values (e.g. e_n, e_t) for the function to train:
  e_nrtnc=c(train1,train2)
  data_zs_t=data_zs[,colnames(data_zs) %in% e_nrtnc]
  e11_zs_t=e11_zs[,colnames(e11_zs) %in% e_nrtnc]

  #lim_anno is the limited annotation: dim: 567 21
  # group is the validation group name in annotation (column 20)
  # Selecting Luminal A group for the validation:
  a11=lim_anno
  log_a11_na=is.na(a11[,20]) #so nas are evaluated
  log_a11_la=(a11[,20]==as.character(validate)) #la comes from luminal A but it could be something else as well,
group is evaluated
  lA=a11[(log_a11_na!=TRUE & log_a11_la),20] #no nas but all group members
  a11_rn=rownames(a11)
  la_rn=a11_rn[log_a11_na!=TRUE & log_a11_la] #for names in the column

  #values for the validation:
  e11_zs_t2=e11_zs[,la_rn]
  data_zs_t2=data_zs[,la_rn]

  return(list((e11_zs_t), (data_zs_t), (e11_zs_t2), (data_zs_t2)))
}

ftv2o=f_train_validate2(lim_anno=a11, e11_zs, data_zs, train1=e_n, train2=e_t, validate = 'Luminal A', list)
e11_zs_to=as.matrix(ftv2o[[1]]) #e11_zs_to[1:5,1:5] #ok e11_zs_t[1:5,1:5]
data_zs_to=as.matrix(ftv2o[[2]]) #data_zs_to[1:5,1:5] #ok data_zs_t[1:5,1:5]
e11_zs_t2o=as.matrix(ftv2o[[3]]) #e11_zs_t2o[1:5,1:5] #ok e11_zs_t2[1:5,1:5]
data_zs_t2o=as.matrix(ftv2o[[4]]) #data_zs_t2o[1:5,1:5] #ok data_zs_t2[1:5,1:5] #oks

#Test:
l_tot_s_test2=fi_loop2( e11_zs_t, e11_zs_t2, data_zs_t, data_zs_t2,mir_amount, lists_s_tot)
ltot1_req2=f_requested_limooo(l_tot_s_test2,lists_s_tot) #See description later

# Shuffling the patient names to check if my results are arising from random or not
f_shuffle_rename = function (e11_zs, n_times) {
  #sample vector
  c_nam_e11_zs=vector("list", n_times)

  for (i in 1:n_times) (
    c_nam_e11_zs[[i]]=colnames(e11_zs)[sample(1:dim(e11_zs)[2], dim(e11_zs)[2], replace=F)]
  )
  return(c_nam_e11_zs)
}

# Testing the model with random shuffling of the patient names:
suffles=f_shuffle_rename(e11_zs, 20)

#This is for the random all model:
f_limos_rnd2 = function (e11_zs_t, e11_zs_t2, data_zs_t, data_zs_t2, mir_amount, list, suffles) {

```



```

tot=vector("list", length(suffles))
for (i in 1:length(suffles)) (
  tot[[i]]=fi_loop2(e11_zs_t[suffles[[i]]], e11_zs_t2[suffles[[i]]], data_zs_t, data_zs_t2, mir_amount, list)
)
return(tot)
}
}
shuffling_comp_new=f_limos_rnd2(e11_zs_t, e11_zs_t2, data_zs_t, data_zs_t2, mir_amount, list, suffles)

```

15e: cross-validation:

#Validation groups e11_zst etc:

```

f_cval_grp= function (sample_nro,div_nro) {
  ab=sample_nro # sample nro vector in 1:i, e.g. 1 2 3 ... i
  max =div_nro #amount of samples in the sample pool
  x <- seq_along(ab)
  d1 <- split(ab, ceiling(x/max))
#d1 #http://stackoverflow.com/questions/3318333/split-a-vector-into-chunks-in-r
#str(d1)

```

```

a=NULL
b=NULL
for (i in 1:length(d1)) (
  a[[i]]= d1[[i]])
for (i in 1:length(d1)) (
  b[[i]]=c(1:567)[(c(1:567) %in% d1[[i]])!=T] )
e_val1=NULL
e_val2=NULL
d_val1=NULL
d_val2=NULL

```

```

for (i in 1:length(d1)) (
  e_val1[[i]] = e11_zs[,b[[i]]] )
for (i in 1:length(d1)) (
  e_val2[[i]] = e11_zs[,a[[i]]] )
for (i in 1:length(d1)) (
  d_val1[[i]] = data_zs[,b[[i]]] )
for (i in 1:length(d1)) (
  d_val2[[i]] = data_zs[,a[[i]]] )

```

```

return(list(e_val1,e_val2,d_val1,d_val2))}

```

#cv codes

```

f_cm_walk_opwj = function (a, mir_amount, list,TL,MG) {
  library("gurobi")
  library("Matrix")

```

#Cross-validation values:

```

e11_zs_t=NULL
e11_zs_t2=NULL
data_zs_t=NULL
data_zs_t2=NULL

```

Making the initial vector:

```

mg_list=vector(mode="list", length=length(a[[1]]))
#names(mg_list) <- f_ent_g symb(list)
names(mg_list) <- c("G1","G2","G3","G4","G5")

```

```

# for the restriction linear model:
gene_nam=vector(mode="list", length=length(list))
mirs=vector(mode="list", length=length(list))

mi=mir_amount
#str(a)
for (j in 1:length(a[[1]])) {
  e11_zs_t[[j]]=a[[1]][[j]]
  e11_zs_t2[[j]]=a[[2]][[j]]
  data_zs_t[[j]]=a[[3]][[j]]
  data_zs_t2[[j]]=a[[4]][[j]]

for (i in 1:length(list)) {
  gene_nam[[i]] = as.vector((list)[i])
  mirs[[i]] = f_mir_limoo(tush_n1, as.vector(unlist(gene_nam[[i]])))
  mg_list[[j]][[i]] = f_limo_restr2(e11_zs_t[[j]], e11_zs_t2[[j]],data_zs_t[[j]], data_zs_t2[[j]],
    as.vector(unlist(gene_nam[[i]])), as.vector(unlist(mirs[[i]])),mi,TL,MG)
  names(mg_list[[j]][[i]]) <- f_ent_gsymb(list[[i]])}
}
#here are the values:
return(mg_list)
}

```

```

f_cv_cor1 = function (a,tsc5,cv_0.5) {
#data, list and cv result
t_i=vector("list",length=length(tsc5))
for (j in 1:length(t_i)) (

for (i in 1:length(a[[1]])) (
  t_i[[j]][[i]]=cv_0.5[[i]][[j]][[4]] ) )

at=as.matrix(mapply(function(y) mean(y[1:5]), t_i))
rownames(at)=f_ent_gsymb((tsc5))
colnames(at)=c('cor_coefficient') #,'tot_miRs_used','miRs_below_zero','miRs_enr')
at=as.matrix(at[rev(order(at[,1])),])
atm=apply(at,2,median) #median of mean values of correlations

```

```

#mean(unlist(sapply(t_i[1:5], '[', 'element'))))
return(list(at, atm)) }

```

```

# looping genes in the list, a is the sample groups
f_cmo_m= function (a, mir_amounts, list, TL,MG) {
tot=vector("list", length(mir_amounts))
mir_amount=NULL
cor_cofs=NULL
for (i in 1:length(mir_amounts)) {
  mir_amount[[i]]=mir_amounts[[i]];
  tot[[i]]=f_cm_walk_opwj(a, mir_amount[[i]], list,TL,MG);
  cor_cofs[[i]]=f_cv_cor1(a,list,tot[[i]]) }
return(cor_cofs)
}

```

Step 15f

Obtaining results after linear modelling analysis

1) Requested information 1:

```
# the correlation coefficients,
# how many miRNAs you have used for the predictions,
# how many of them are above zero
# how many of all of the miRNAs are enriched ones
```

```
# A function for this:
```

```
f_requested_limooo = function (sl_valid_test2,list) {
  aeu <- vector("list", length(list))

  for (i in 1:length(aeu)) {
    aeu[[i]][1]=as.matrix(sl_valid_test2[[i]][[4]])
    aeu[[i]][2]=(length(sl_valid_test2[[i]][[2]])-1)
    aeu[[i]][3]=sum(sl_valid_test2[[i]][[2]]>0)
    aeu[[i]][4]=(length(sl_valid_test2[[i]][[3]]))
  }
  aeu3=aeu

  #The additional information
  x=f_ent_g symb((list))
  aeu_et=do.call(rbind, aeu3)
  colnames(aeu_et)=c('cor_coefficient','tot_miRs_used','miRs_mod_nz','miRs_selected')
  rownames(aeu_et)=x
  aeu_et=aeu_et[rev(order(aeu_et[,1])),]
  write.table(aeu_et, "top_150_cor_gene_alpat3_1.txt", sep="\t")
  return(aeu_et)
}
```

```
# 2) All correlations for random validation:
```

```
f_cor_rnds = function (shuff_comp_new,comp_list_length) {
  #In case of top 150 list, comp_list_length was 150...
  cor_suf2a=vector("list",comp_list_length)
  for (i in 1:length(shuff_comp_new)) (
    for (j in 1:length(cor_suf2a)) (
      cor_suf2a[[j]][[i]]=shuff_comp_new[[i]][[j]][[4]] )
  names(cor_suf2a)=names(shuff_comp_new[[1]])
  hn1=(unlist(cor_suf2a))
  return(hn1)
}
```

```
hn1=f_cor_rnds(shuff_comp_new,150)
```

```
# 3) Calculating means, and Doing Histograms of linear model correlations, and t-testing these correlations (plot all in same picture with powerpoint):
```

```
#rnd
mean(hn1)
#valid
hn2=(l3_r_oklma[,1])
mean(hn2) #[1] 0.2983685
#nrm
hn3=(l3_r_ok[,1])
mean(hn3) #[1] 0.5517583
```

```
# Histogram Grey Color
```

```
hn1a=hist(hn1, col=rgb(0.1,0.1,0.1,0.5),xlim=c(-0.2,0.8), ylim=c(0,250),breaks=25, main="Overlapping Histogram")
hn1b=hist(hn2, col=rgb(0.8,0.8,0.8,0.5), add=T,breaks=25,xlim=c(-0.2,0.8),ylim=c(0,25))
hn1c=hist(hn3, col=rgb(0.4,0.4,0.4,0.4), add=T,breaks=25,xlim=c(-0.2,0.8),ylim=c(0,25))
```

```

#Scaling (to the smallest curve i.e. here hn1c, according to its counts, without changing it)
hn1a[[2]]=hn1a[[2]]*max(hn1c$counts)/max(hn1a$counts)
hn1b[[2]]=hn1b[[2]]*max(hn1c$counts)/max(hn1b$counts)
#Plotting
plot(hn1a,col=rgb(0.4,0.1,0.1,0.1),xlim=c(-0.2,0.9),ylim=c(0,15),xlab='Correlation coef.',ylab='Genes with this
correlation (scaled)',main="")
plot(hn1b,col=rgb(0.1,0.1,0.1,0.5), add=T,xlim=c(-0.2,0.8),ylim=c(0,15))
plot(hn1c,col=rgb(0.8,0.8,0.8,0.8),add=TRUE)
lines(density(hn1))
lines(density(hn2))
lines(density(hn3))

# Just densities:
plot(density(rndl_cor_f2),xlim=c(-0.3,0.9),ylim=c(0,8),xlab='Correlation coefficient',ylab='Genes with this correlation
(scaled)',main="")
lines(density(hn2))
lines(density(hn3))

#Testing the correlations with each other using t-test:
#rnd against norm:
cor_t2_rnd_nrm=t.test(hn1,hn3)
cor_t2_rnd_nrm[3] #[1] 9.346306e-74

#rnd against validation:
cor_t2_rnd_val=t.test(hn1,hn2)
cor_t2_rnd_val[3] #[1] 4.492909e-10

#validation against norm:
cor_t2_val_nrm=t.test(hn2,hn3)
cor_t2_val_nrm[3] #[1] 2.988722e-22

# Boxplotting interesting gene , enriched miRNA miRNA combinations

gene_name=c("LUZP1")
f_bp_eval=function (lst_fnlm,nrm_tst_f,gene_name) {
  #for the enriched miRNAs
  auk=as.character(gene_name)
  auk=as.character(f_gsymb_ent(auk))
  nrou=c(1:length(lst_fnlm))[lst_fnlm %in% auk]

  qkireg2=gen_rgl_mir2[gen_rgl_mir2[,1] %in% names(nrm_tst_f[[nrou]][[3]][nrm_tst_f[[nrou]][[3]][,1]<0,]),] #"hsa-
miR-3609","hsa-miR-29b-3p"
  qkireg2=qkireg2[rev(order(qkireg2[,2])),]
  #qki=gen_rgl[gen_rgl[,1] %in% auk,]

  # for 2 best enr mirs
  umt=data_zs[(rownames(data_zs) %in% qkireg2[1,1]),(colnames(data_zs) %in% e_t)]
  umn=data_zs[(rownames(data_zs) %in% qkireg2[1,1]),(colnames(data_zs) %in% e_n)]
  umt2=data_zs[(rownames(data_zs) %in% qkireg2[2,1]),(colnames(data_zs) %in% e_t)]
  umn2=data_zs[(rownames(data_zs) %in% qkireg2[2,1]),(colnames(data_zs) %in% e_n)]
  # for the gene
  dt=e11_zs[(rownames(e11_zs) %in% auk),(colnames(e11_zs) %in% e_t)]
  dn=e11_zs[(rownames(e11_zs) %in% auk),(colnames(e11_zs) %in% e_n)]

  return(boxplot(umn,umt,umn2,umt2,dn,dt,col = rep(c("blue","red"),3), ylab="z-scored expression level")) }

bp_eval=f_bp_eval(lst_fnlm,nrm_tst_f,gene_name)

```

```

#the most regulated enriched miRNAs (according to the limo):
f_bp_eval2=function(lst_fnlm,nrm_tst_f,gene_name,emrs) {
  auk=as.character(gene_name)
  auk=as.character(f_g symb_ent(auk))
  nrou=c(1:length(lst_fnlm))[lst_fnlm %in% auk]

  qkireg2=gen_rgl_mir2[gen_rgl_mir2[,1] %in% names(nrm_tst_f[[nrou]][[3]][nrm_tst_f[[nrou]][[3]][,1]<0,]),] #"hsa-
miR-3609","hsa-miR-29b-3p"
  qkireg2=qkireg2[rev(order(qkireg2[,2])),]
  qkireg2= if (dim(qkireg2)[2]<emrs) (qkireg2=qkireg2) else (qkireg2=qkireg2[c(1:emrs),])
  #dim(qkireg2) #traditional worked ok...
  #qki=gen_rgl[gen_rgl[,1] %in% auk,]
  umn=NULL
  umt=NULL
  for (i in 1:emrs) {
    umn[[i]]=data_zs[(rownames(data_zs) %in% qkireg2[i,1]),(colnames(data_zs) %in% e_n)]
    umt[[i]]=data_zs[(rownames(data_zs) %in% qkireg2[i,1]),(colnames(data_zs) %in% e_t)]
  }
  #For the genes
  dt=e11_zs[(rownames(e11_zs) %in% auk),(colnames(e11_zs) %in% e_t)]
  dn=e11_zs[(rownames(e11_zs) %in% auk),(colnames(e11_zs) %in% e_n)]
  dt=as.data.frame(dt)
  dt=as.list(dt)
  dn=as.data.frame(dn)
  dn=as.list(dn)

  c_tot=c(umn,umt)
  #Selecting normals and tumours (mirs)
  cu=NULL
  for(i in 1:emrs) {
    cu[[i]]=list(1+(i-1),emrs+(i))
  }

  cu=unlist(cu)
  c_tot=c_tot[cu]
  #combining miRs and genes
  c_toti=c(c_tot,dn,dt)

  #boxplot(x = as.list(umt)) #this works but I need them all
  #http://stackoverflow.com/questions/11346880/r-plot-multiple-box-plots-using-columns-from-data-frame
  ae=boxplot(c_toti,col = rep(c("blue","red"),3), ylab="Z-scored expression level")
  return(list(qkireg2,ae)) }

# bp to luzp1 with best up.reg.mirs to next function f_bp_eval3
c(1:247)[names(nrm_tst_f) %in% c("LUZP1")] #27
rownames(nrm_tst_f[[27]][2][[1]])
cond_m=c(1:332)[gen_rgl_mir2[,1] %in% rownames(nrm_tst_f[[27]][2][[1]])]
luz_mir_up=gen_rgl_mir2[cond_m,]
luz_mir_up=luz_mir_up[rev(order(luz_mir_up[,2])),]
luz_mir_upi=luz_mir_up[1:2,]

#most up-regulated (select by hand in limo) mirs are checked with gene
f_bp_eval3=function(lst_fnlm,nrm_tst_f,gene_name,emrs,luz_mir_upi) {
  qkireg2=luz_mir_upi #if you want to select mirs then use: luz_mir_upi
  umn=NULL

```

```

umt=NULL
for (i in 1:emrs) {
  umn[[i]]=data_zs[(rownames(data_zs) %in% qkireg2[i,1]),(colnames(data_zs) %in% e_n)]
  umt[[i]]=data_zs[(rownames(data_zs) %in% qkireg2[i,1]),(colnames(data_zs) %in% e_t)]
}

#For the genes
dt=e11_zs[(rownames(e11_zs) %in% auk),(colnames(e11_zs) %in% e_t)]
dn=e11_zs[(rownames(e11_zs) %in% auk),(colnames(e11_zs) %in% e_n)]
dt=as.data.frame(dt)
dt=as.list(dt)
dn=as.data.frame(dn)
dn=as.list(dn)

c_tot=c(umn,umt)
#str(c_tot)
cu=NULL
for(i in 1:emrs) {
  cu[[i]]=list(1+(i-1),emrs+(i))
}

cu=unlist(cu)
c_tot=c_tot[cu]
c_toti=c(c_tot,dn,dt)
ae=boxplot(c_toti,col = rep(c("blue","red"),3), ylab="Z-scored expression level")
return(list(qkireg2,ae))
}

```