

Project.R

2023-12-09

```
# Name: KUSH Patel
# Project
# cohort: GDC TCGA Cervical Cancer (CESC) - HTSeq - Counts
# Exploring Grade Gene Expression (G1 vs G3) Analysis in GDC TCGA Cervical Cancer (CESC) using RNA-Seq Data
# https://xenabrowser.net/datapages/?cohort=GDC+TCGA+Cervical+2013Cancer&z=0&arrow=hub+htseq&data=TCGA+Xena.
# treehouse-gi.ucsc.edu/33443
#####
# Libraries
library(dplyr)

##
## Attaching package: 'dplyr'

##
## The following objects are masked from 'package:stats':
##
## filter, lag

##
## The following objects are masked from 'package:base':
##
## intersect, setdiff, setequal, union

library(ggplot2)

##
## Warning: package 'ggplot2' was built under R version 4.3.2

library(UCSCXenaTools) # needed to retrieve data

##
## Warning: package 'UCSCXenaTools' was built under R version 4.3.2

#####
## UCSCXenaTools version 1.4.8
## Project URL: https://github.com/ropensci/UCSCXenaTools
## Usages: https://cran.r-project.org/web/packages/UCSCXenaTools/vignettes/UCSCXenaTools.html
##
## If you use it in published research, please cite:
## Wang et al., (2019). The UCSCXenaTools R package: a toolkit for accessing genomics data
## From UCSC Xena Platform, From cancer multi-omics to single-cell RNA-seq,
## Journal of Open Source Software, 4(48), 1627. https://doi.org/10.2196/joss.01627
#####
## --Enjoy it--

library(edgeR) # needed for processing, such as TMM

## Loading required package: limma

library(limma) # needed to find DE probes

#####
# 2) Retrieve Data
data(XenaData)
```

```

GDC TCGA Cervical Cancer (CESC)

# limit to desired cohort
cesc <- XenaData %>% filter(XenaCohorts == "GDC TCGA Cervical Cancer (CESC)")

# Get the phenotype / clinical data
cli_query = cesc %>%
  filter(Label == "Phenotype") %>% # select clinical dataset
  XenaGenerate() %>% # generate a XenaHub object
  XenaQuery() %>% # generate the query
  XenaDownload() # download the data

## This will check url status, please be patient.

## All downloaded files will under directory C:\Users\spatel\AppData\Local\Temp\Rtmp04FBd7.

## The 'trans_slash' option is FALSE, keep same directory structure as Xena.

## Creating directories for datasets...

## Downloading TCGA-CESC_GDC_phenotype.tsv.gz

# prepare (load) the data into R
cesc_pheno <- XenaPrepare(cli_query)

# Get the RNA-seq data, including the "probe map"
cli_query <- cesc %>% filter(Label == "HTSeq - Counts") %>%
  XenaGenerate() %>% # generate a XenaHub object
  XenaQuery() %>%
  XenaDownload(download_probeMap = TRUE)

## This will check url status, please be patient.

```

```
## Check ProbeMap urls of datasets.

## All downloaded files will under directory C:\Users\patel\AppData\Local\Temp\Wtsp04FBd7.

## The 'trans_slash' option is FALSE, keep same directory structure as Xena.

## Creating directories for datasets...

## Downloading YCGA-CESC.htseq_counts.tsv.gz

## Downloading gencode.v22.annotation.gene.probeMap
```

```
# prepare (load) the data into R
cesc_counts <- XenaPrepare(c1i_query)

#####
# (2) Data pre-processing: we need to do a fair amount of
# filtering and re-arranging to work with the data, so
# that the expression and phenotype data are aligned
#####
```

```
# First, let's use more manageable names
# - X: expression data, with probes as row names
# - probemap: the probemap
# - Y: the pheno/clinical data

# For X, we need to set the rownames and remove the probe column
# from the data matrix
X <- data.frame(cesc.counts$TCGA.CESC.htseq.counts.tsv.gz)
rownames(X) <- X$ENSEMBL_ID
X <- X[,~1] # remove the probe name column

# probemap = probe names
probemap <- cesc.counts$genome.v22.annotation.gene.probeMap

# Y = pheno data
Y <- cesc_pheno

# The expression and clinical data need to match; currently, e.g.,
# The first column of the expression data does not correspond
# to the first row of the pheno data; the data is also not
# in a consistent format (one has '.', and the other has '-')

# compare sample names between X and Y; they do not match, and are
# not even in the same format!
print(colnames(X)[1])

## [1] "TCGA_DS_A7H1_03A"
```

```
print(Y$submitter_id.samples[1])
```

```
## [1] "TCGA-EA-A97N-01A"
```

```
# 'change '.' to '-' so sample ID format is consistent
colnames(X) <- gsub("\\.", "-", colnames(X))
```

Note that the sample ID is a barcode that has a special meaning:

```
# https://docs.gdc.cancer.gov/Encyclopedia/pages/TCGA_Barcode/
# In particular, the 4th section describes the "Sample" which is
# either tumor (01 - 09) or normal (10-19). For details see:
# https://gdc.cancer.gov/resources-tcga-users/tcga-code-tables/sample-type-codes

# Keep only the '01A' tumor samples
g <- grep("01A$", colnames(X))
X <- X[,g]

# We still need to match the expression data with the clinical data
# Let's do that by first finding the samples that are common
```

```
# Between the expression and clinical data. We can use
# intersect(a,b) to return a vector containing the elements common
# to vectors 'a' and 'b'

common_samples <- intersect(colnames(X), Y$Submitter_Id.samples)

# we then use match(x, y) to get a vector of indices. The value
# x[i] is the index of 't' containing the 1st value of 'x'

nx <- match(common_samples, colnames(X))
my <- match(common_samples, Y$Submitter_Id.samples)
```

```
X <- X[,mx]
Y <- Y[,my,]

# Make sure that the samples match -- if they don't, this will produce an error
stopifnot(all(colnames(X) == Y$submitter_id.samples))

#Total number of samples and probes before processing
n="Total number of samples"
print(ncol(X))
```

```
## [1] 296
```

```
a"Total number of probes"  
print(nrow(X))
```

```
## [1] 60458
```

```
#####
```

```
# Setup step 3: Process the expression data
#####

# convert from log2(count + 1) to count data
X <- round(2**X - 1)

# remove genes with low counts
dge <- DGEList(counts=X)

keep <- filterByExpr(dge,min.prop = .10 )
```

```
## Warning in FilterByExpr.DGEList(dge, min.prop = 0.1): All samples appear to
## belong to the same group.

dge <- dge[keep, keep.lib.sizes=FALSE]

# apply TMM normalization, which computes the normalization
# factors. The actual normalization is done in a later step
dge <- calcNormFactors(dge, method = "TMM")

# Calculate the log CPN values, using the normalization factors:
```

```
# 3 counts are added to each observation to prevent log 0 values
logCPM <- cps(dge, log = TRUE, prior.count = 3)

#Total number of samples and probes after processing
nTotalNumberofSamples"
print(ncol(logCPM))
```

```
#Total number of "probes"
print(nrow(LogCPW))

## [1] 23996

#print(colnames(X)[2])
#print(Y$subscriber_id.samples[1])

#####
```

```
#5. Box plot for the first 10 samples
#####

boxplot(logCPM[, 1:10],
        main = "Boxplot of Normalized Log-CPM Values for First 10 Samples",
        ylab = "Log-CPM Values")
```

Box plot showing CPM Values for 10 genes. The y-axis represents CPM Values from 0 to 20. The x-axis lists genes: GATC, GATC2, GATC3, GATC4, GATC5, GATC6, GATC7, GATC8, GATC9, and GATC10. Each gene has a box plot showing the distribution of CPM values. GATC, GATC2, GATC3, GATC4, GATC5, GATC6, GATC7, GATC8, and GATC9 have medians around 12-14, while GATC10 has a median around 10. All genes show a similar pattern of outliers.

```
#####  
#6 Done  
#7 Extract the column for the grade  
#####  
  
grade <- YSneoplasm_histologic_grade  
grade[is.na(grade)] <- "unknown"  
design <- model.matrix(~1+grade)  
colnames(design) <- c("G1a", "G2", "G3", "G4", "Gx", "Unknown")
```

```
#design
columns_to_keep <- colnames(design)[!(colnames(design) %in% c("G2", "G4", "GK", "Unknown"))]
design <- design[, columns_to_keep]
head(design)

##      G1 G3
## 1 0 0
## 2 1 0
## 3 0 0
## 4 0 1
```

```
## 5 0 0
## 6 0 0

#####
# Use limma to find probes across groups using fdr of 10%
#####

#fit the linear model to each row of the expression matrix
fit <- lmFit(logCPM, design)

#(G3 - G1)
```

```
contrast.matrix <- makeContrasts(G3 - G1, levels = design)
#fit model based on contrasts (e.g., G3 - G1)
fit <- contrasts.fit(fit, contrast.matrix)
#apply the 'eBayes' step to calculate moderated t statistics
fit.de <- eBayes(fit, trend = TRUE)
#FDR < 10%
#FDR < 10% Resulted in 0 Values
tt.1 <- topTable(fit.de, sort.by = "p", p.value = 1, number = 30)
nrow(tt.1)
```

```
# [1] 30
```

```
#Top 30 (All of them)
```

```
tt.1
```

	logFC	AveExpr	t	P.Value	adj.P.Val	
##						
##	ENSG00000225439.2	1.4659931	0.5264893	3.999497	7.928972e-05	0.6467321
##	ENSG00000226810.3	-1.4749898	0.74907813	-3.998740	8.831744e-05	0.6467321
##	ENSG00000154413.12	-3.1950797	0.6866778	-3.995204	8.878497e-05	0.6467321
##	ENSG00000226811.1	-3.0881296	0.6427719	-3.827681	1.484464e-04	0.8567822

#	ENSG00000016589.1	-1	-3.044338	0.688448	0.7	-6.849343	3.950801e-04	0.999822
#	ENSG00000015499.1	1	1.564449	0.5347915	3	6.204608	3.432043e-04	0.999822
#	ENSG0000000985.8	-2	-2.506409	-0.2907249	1	-5.565763	4.341406e-04	0.999822
#	ENSG00000010091.3	3	3.4274138	0.2265728	3	3.537777	6.677144e-04	0.999822
#	ENSG000000066709.2	12	2.743796	1.3771887	18	4.37348	6.771474e-04	0.999822
#	ENSG00000007933.11	-2	-2.446836	-0.5333339	3	-3.732328	8.692256e-04	0.999822
#	ENSG00000010455.12	1	1.521856	0.3857484	3	3.584645	3.432043e-04	0.999822
#	ENSG00000010455.12	1	1.603516	0.8905618	3	3.322655	3.942626e-04	0.999822
#	ENSG0000014057.7	-1	-1.451920	0.1699878	3	-2.87565	2.196262e-03	0.999822
#	ENSG00000124969.3	12	2.266128	1.1795747	3	2.25349	1.392565e-03	0.999822
#	ENSG00000104576.7	12	2.266128	1.1795747	3	2.25349	1.392565e-03	0.999822

#	ENSG0000026361.5	-2.735865	87379967	-1.199171	1.520907e-3	0.999802
#	ENSG00000205890.3	-1.277244	0.4230454	-0.154159	1.783722e-3	0.999802
#	ENSG00000272286.1	-1.216635	0.1427084	-0.145957	1.813994e-3	0.999802
#	ENSG000000777.12	-2.905187	1.1557556	-0.145957	1.813994e-3	0.999802
#	ENSG00000205572.5	0.872601	0.6442598	-0.140011	8.619296e-3	0.999802
#	ENSG00000019620.3	-2.537617	0.1461604	-0.142588	1.848746e-3	0.999802
#	ENSG00000195353.3	-2.910191	-0.1345308	-0.147051	1.880614e-3	0.999802
#	ENSG00000164122.1	-2.872942	0.8097639	-0.130188	1.555555e-3	0.999802
#	ENSG00000240370.5	-0.7043624	0.11535959	-0.135931	8.774326e-3	0.999802
#	ENSG00000247619.6	-0.910944	0.3805924	-0.114517	2.014669e-3	0.999802

[illegible]

```
## ENSG00000044199.10 - 4.433452
## ENSG000000100951.6 - 4.439384
## ENSG000000100973.3 - 4.441886
## ENSG00000060789.12 - 4.450171
## ENSG000000007933.11 - 4.455911
## ENSG000000162571.12 - 4.457637
## ENSG000000175170.13 - 4.459374
## ENSG000000140057.7 - 4.466884
## ENSG000000127249.13 - 4.468565
## ENSG000000176497.9 - 4.470767
## ENSG000000263961.5 - 4.470768
```

```
## ENSG00000205890.3 -4.474742
## ENSG00000271200.1 -4.476197
## ENSG00000090577.13 -4.476292
## ENSG00000272556.1 -4.475440
## ENSG00000106260.3 -4.475478
## ENSG00000197353.3 -4.475631
## ENSG00000164746.12 -4.476093
## ENSG00000240370.5 -4.476024
## ENSG00000241769.6 -4.477783
## ENSG00000172328.12 -4.477885
## ENSG00000116682.13 -4.477815
```

```
## ENSG0000027078.1 -4.478115
## ENSG000000907715.12 -4.470243
## ENSG00000186710.10 -4.470755

#####
#0 Top probe boxplot - LOOK OVER
#####

probe <- rownames(tt.1)[1]
m <- match(probe, rownames(logCPM))
```

```
df <- data.frame(expr = logCPM[m,], grade = grade)

# filter out 'unknown' samples
df <- df[(df$grade %in% c("G2", "G4", "G6", "Unknown")), ]

# convert from logFC to FC
logFC <- tt.$logFC[i]
2**logFC
```

```
## [1] 2.758700

## visualize ##
FC <- paste0("FC = ", round(2**logFC, 2))
main <- paste0("Expression of ", probe, ", ", FC, ", FDR = 100%")

ggplot(dfr, aes(x = grade, y = expr, fill = grade)) + geom_boxplot() +
  ylab(log2(expression)) + ggtitle(main) +
  scale_fill_manual(values = c("pink", "lightblue")) +
  theme_classic() + theme(legend.position = "none")
```

Expression of ENSG00000225439.2, FC = 2.76, FDR = 100%

#10 OF all genes

gene names, probe names, logFC, and adjusted p-values

```

#FDR
m <- match(probe_names, probeMap$id)
gene_names <- probeMap$gene[m]
logFC_Values <- tt$logFC
p_Values <- tt$adj.P.Val

#DF
result_df <- data.frame(
  GeneNames = gene_names,
  ProbeNames = probe_names,
  LogFC = logFC_Values,
  P.Val = p_Values
)

```

```
AdjPValues = p.Values
)

head(result_df, 5)
```

#	GeneNames	ProbeNames	LogFC	AdjPValues
# 1	BOLA3-AS1	ENSG00000225439.2	1.463993	0.6461721
# 2	RNASE2	ENSG00000259818.3	-1.747470	0.6461721
# 3	FUT8	ENSG00000156413.12	-3.195680	0.6461721
# 4	TKO3	ENSG00000125780.11	-3.085128	0.9565702

```
# 5       CUMLE EN50000010059.11 -1.944338 0.9999622

#####
#11 Heat Map - Recheck
#####

m <- match(result_df$ProBNames, rownames(logCPM))
expr <- logCPM[m,]
rownames(expr) <- result_df$GeneNames
# create a color range consisting of 200 values between yellow and blue
col = colorRamp2(1:length(result_df$GeneNames), 0, 1)
```

```
col.heat <- colorRampPalette(c( "yellow", "blue" ))(100)

# set colors for grade
col.grade <- as.integer(as.factor((grade %in% c("Q2", "Q4", "GK", "unknown"))))
col.grade <- c("magenta", "lightgreen")(col.grade)

# Generate the heatmap
heatmap(expr, colSideColors = col.grade, col = col.heat, scale = "none")
```

```
#write.table(genes, row.names = FALSE, quote = FALSE, file = "top_genes.txt")
```