```
Github repositorio: https://github.com/aescos/Escos-Lopez-Alejandra-PEC1.git
====== >>>>> 8d36fb9e244eaba755823a44ce7d45f617bbe1e0
To erase data and start clean:
 rm(list = ls())
To push changes in Git after working on the PEC1.
Los datos que voy a emplear: Descripcion:
"The acompanying dataset has been obtained from a phosphoproteomics experiment that was performed to analyze (3 + 3) PDX models of two
different subtypes using Phosphopeptide enriched samples. LC-MS analysis of 2 technical duplicates has been performed on each sample. The
results set consisted of Normalized abundances of MS signals for ca. 1400 phosphopeptides Goal of the analysis: *search phosphopeptides that
allow differentiation of the two tumor groups This should be made with both Statistical Analysis and visualization. Data have been provided as an
excel file: TIO2+PTYR-human-MSS+MSIvsPD.XLSX
Groups are defined as:
MSS group: Samples M1, M5 and T49, PD group: Samples M42, M43 and M64 with two technical replicates for each sample The first column,
SequenceModification contains abundance values for the distinct phosphopetides. Other columns can be omitted."
 library(readxl)
 data_phospho <- read_excel("TIO2+PTYR-human-MSS+MSIvsPD.XLSX")</pre>
 head(data_phospho)
 ## # A tibble: 6 × 18
    SequenceModifications Accession Description Score M1_1_MSS M1_2_MSS M5_1_MSS
                                          <chr> <dbl> <dbl> <dbl>
                               <chr>
                                                                                  <dbl>
 ## 1 LYPELSQYMGLSLNEEEIR[2]... 000560 Syntenin-1... 48.1 24.3 44476.
                                                                                     0
 ## 2 VDKVIQAQTAFSANPANPAILS... 000560 Syntenin-1... 67.0 0 43139. 2102.
 ## 3 VIQAQTAFSANPANPAILSEAS... 000560 Syntenin-1... 77.7 3413. 172143. 77323.
 ## 4 HADAEMTGYVVTR[6] Oxida... 015264 Mitogen-ac... 44.9 220431. 145657. 104288.
 ## 5 HADAEMTGYVVTR[9] Phosp... 015264 Mitogen-ac... 67.4 18255.
                                                                         8530. 35956.
                                        Claudin-3 ... 63.7 644513. 261938. 187023.
 ## 6 STGPGASLGTGYDR[12] Pho... 015551
 ## # i 11 more variables: M5_2_MSS <dbl>, T49_1_MSS <dbl>, T49_2_MSS <dbl>,
 ## # M42 1 PD <dbl>, M42 2 PD <dbl>, M43 1 PD <dbl>, M43 2 PD <dbl>,
 ## # M64 1 PD <dbl>, M64 2 PD <dbl>, CLASS <chr>, PHOSPHO <chr>
Las anotaciones sobre los datos son las siguientes:
 anotacion <- read excel("TIO2+PTYR-human-MSS+MSIvsPD.XLSX", sheet = "targets")</pre>
 ## New names:
 ## • `Sample` -> `Sample...1`
 ## • `Sample` -> `Sample...2`
Voy a testar a ver si hay duplicados primero para poder ver que utilizo como row names.
 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
  # Duplicados basados en una columna específica
 duplicados_peptide_dplyr <- data_phospho %>%
   filter(duplicated(SequenceModifications) | duplicated(SequenceModifications, fromLast = TRUE))
 print(duplicados_peptide_dplyr)
 ## # A tibble: 2 × 18
      SequenceModifications Accession Description Score M1_1_MSS M1_2_MSS M5_1_MSS
                               <chr> <chr>
                                                      <dbl> <dbl> <dbl> <dbl> <dbl>
 ## 1 GEPNVSYICSR[7] Phospho... P49840 Glycogen s... 54.3 1.18e7 8689448. 5130833.
 ## 2 GEPNVSYICSR[7] Phospho... P49840 Glycogen s... 46.9 1.79e4 17796.
 ## # i 11 more variables: M5_2_MSS <dbl>, T49_1_MSS <dbl>, T49_2_MSS <dbl>,
 ## # M42_1_PD <dbl>, M42_2_PD <dbl>, M43_1_PD <dbl>, M43_2_PD <dbl>,
 ## # M64 1 PD <dbl>, M64 2 PD <dbl>, CLASS <chr>, PHOSPHO <chr>
Aquí podemos observar que hay un peptido duplicado pero cuando vemos el lugar de fosforilacion en la utlima columna es distinto. El primer
lugar de fosforilacion es Y (tirosina) y el segundo es S/T (Serina/Treoina) por lo tanto no esta duplicado el peptido. Por lo tanto voy a unir las dos
columnas por un "-" y asi seria una lista unica de peptidos.
 library(stringr)
 data_phospho.unica <- data_phospho %>%
   mutate(unica = str_c(SequenceModifications, PHOSPHO, sep = "-"))
Ahora querria tener los nombres de los genes, gene symbol, puesto que es mas sencillo de comprender a la hora de interpretar los resultados.
Para ello vamos a utilizar Biomart. Corremos la instalacion del paquete tan solo una vez.
 library(biomaRt)
 # Conectar al servidor de Ensembl
 ensembl <- useMart("ensembl", dataset = "hsapiens_gene_ensembl")</pre>
 # Lista de códigos de Accession (uniprotswissprot)
 accession codes <- data phospho[,2]</pre>
 # Realizar la consulta
 gene_symbols <- getBM(</pre>
   attributes = c("ensembl_gene_id", "uniprotswissprot", "description", "hgnc_symbol", "gene_biotype"), # Obtenemos d
  istintos codigos para el mismo gen.
   filters = "uniprotswissprot",
   values = accession_codes,
   mart = ensembl
 # Mostrar resultados
 head(gene_symbols)
      ensembl_gene_id uniprotswissprot
 ## 1 ENSG00000156603
                                 A0JLT2
 ## 2 ENSG00000188522
                                 A6ND36
 ## 3 ENSG00000196531
                                 E9PAV3
 ## 4 ENSG0000110696
                                 000193
 ## 5 ENSG0000101856
                                 000264
 ## 6 ENSG0000103319
                                 000418
 ##
                                                                                      description
 ## 1
                               mediator complex subunit 19 [Source: HGNC Symbol; Acc: HGNC: 29600]
 ## 2
               family with sequence similarity 83 member G [Source: HGNC Symbol; Acc: HGNC: 32554]
 ## 3 nascent polypeptide associated complex subunit alpha [Source: HGNC Symbol; Acc: HGNC: 7629]
 ## 4
                       chromosome 11 open reading frame 58 [Source: HGNC Symbol; Acc: HGNC: 16990]
 ## 5
                progesterone receptor membrane component 1 [Source:HGNC Symbol;Acc:HGNC:16090]
 ## 6
                     eukaryotic elongation factor 2 kinase [Source: HGNC Symbol; Acc: HGNC: 24615]
      hgnc_symbol gene_biotype
 ## 1
            MED19 protein_coding
 ## 2
            FAM83G protein_coding
 ## 3
              NACA protein_coding
 ## 4
         Cllorf58 protein_coding
 ## 5
            PGRMC1 protein_coding
 ## 6
            EEF2K protein_coding
Unimos ambas listas para que contenga toda la información para cada peptido.
 # Cambiar el nombre de la columna con dplyr
 gene_symbols <- dplyr::rename(gene_symbols, Accession = uniprotswissprot)</pre>
 gene symbols <- gene symbols %>% distinct(Accession, .keep all = TRUE)
 # Union izquierda, en este caso por data phospho.unica
 data <- inner_join(data_phospho.unica, gene_symbols, by = "Accession")</pre>
 # Eliminar duplicados basados en la columna "unica"
 data <- data %>% distinct(unica, .keep_all = TRUE)
 print(data)
 ## # A tibble: 1,437 × 23
       SequenceModifications Accession Description Score M1_1_MSS M1_2_MSS M5_1_MSS
       <chr>
                               <chr>
                                          <chr>
                                                      <dbl> <dbl>
                                                                         <dbl>
                                                                                  <dbl>
 ## 1 LYPELSQYMGLSLNEEEIR[2... 000560
                                         Syntenin-1... 48.1 24.3 44476.
                                         Syntenin-1... 67.0 0
                                                                        43139. 2102.
 ## 2 VDKVIQAQTAFSANPANPAIL... 000560
 ## 3 VIQAQTAFSANPANPAILSEA... 000560
                                          Syntenin-1... 77.7 3413. 172143. 77323.
 ## 4 HADAEMTGYVVTR[6] Oxid... O15264
                                         Mitogen-ac... 44.9 220431. 145657. 104288.
 ## 5 HADAEMTGYVVTR[9] Phos... 015264
                                         Mitogen-ac... 67.4 18255.
                                                                         8530. 35956.
                                         Claudin-3 ... 63.7 644513. 261938. 187023.
 ## 6 STGPGASLGTGYDR[12] Ph... O15551
                                         Prominin-1... 40.7 686820. 331984. 252694.
 ## 7 DHVYGIHNPVMTSPSQH[4] ... O43490
                                         Prominin-1... 58.3 815186. 728701. 267179.
 ## 8 DHVYGIHNPVMTSPSQH[4] ... O43490
                                         Prominin-1... 60.9 1578.
                                                                         9836.
                                                                                  2033.
 ## 9 RMDSEDVYDDVETIPMK[2] ... O43490
 ## 10 RMDSEDVYDDVETIPMK[2] ... O43490
                                         Prominin-1... 47.4 2815.
                                                                       13247.
                                                                                  2121.
 ## # i 1,427 more rows
 ## # i 16 more variables: M5_2_MSS <dbl>, T49_1_MSS <dbl>, T49_2_MSS <dbl>,
 ## # M42_1_PD <dbl>, M42_2_PD <dbl>, M43_1_PD <dbl>, M43_2_PD <dbl>,
 ## # M64_1_PD <dbl>, M64_2_PD <dbl>, CLASS <chr>, PHOSPHO <chr>, unica <chr>,
 ## # ensembl_gene_id <chr>, description <chr>, hgnc_symbol <chr>,
 ## # gene biotype <chr>
Poner los datos en formato SummarizedExperiment:
Instalamos e inicializamos la libreria.
 # Seleccionar columnas de la matriz
 data <-data.frame(data, row.names = "unica")</pre>
 matrix <- dplyr::select(data, M1_1_MSS, M1_2_MSS, M5_1_MSS, M5_2_MSS, T49_1_MSS, T49_2_MSS, M42_1_PD, M42_2_PD, M
 43_1_PD, M43_2_PD, M64_1_PD, M64_2_PD)
 matrix <- matrix %>% as.matrix()
 # Metadatos de las filas (genes)
 data <-data.frame(data)</pre>
 row data <- data.frame(</pre>
   ensembl_gene_id = data[,19],
   description = data[,20],
   row.names = row.names(data),
   Symbol = data[,21]
 # Metadatos de las columnas (muestras)
 col data <- anotacion
  # Crear el objeto SummarizedExperiment
 se <- SummarizedExperiment(</pre>
   assays = list(counts = matrix), # Asignar los datos de expresión
                          # Asignar los metadatos de los genes
   rowData = row_data,
                                # Asignar los metadatos de las muestras
   colData = col_data
 # Mostrar el objeto
 print(se)
 ## class: SummarizedExperiment
 ## dim: 1437 12
 ## metadata(0):
 ## assays(1): counts
 ## rownames(1437): LYPELSQYMGLSLNEEEIR[2] Phospho [9] Oxidation-Y
    VDKVIQAQTAFSANPANPAILSEASAPIPHDGNLYPR[35] Phospho-Y ...
    YQDEVFGGFVTEPQEESEEEVEEPEER[17] Phospho-S/T YSPSQNSPIHHIPSRR[1]
 ## Phospho [7] Phospho-S/T
 ## rowData names(3): ensembl_gene_id description Symbol
 ## colnames(12): M1_1_MSS M1_2_MSS ... M64_1_PD M64_2_PD
 ## colData names(4): Sample...2 Individual Phenotype
 save(se, file = "Phosphoproteomics.Rda")
Estadisticos descriptivos
 library(knitr)
 library(kableExtra)
 ## Attaching package: 'kableExtra'
 ## The following object is masked from 'package:dplyr':
        group_rows
 summary(matrix)
                            M1_2_MSS
                                                M5_1_MSS
                                                                    M5_2_MSS
 ## Min. :
                     0 Min. :
                                         0 Min. :
                                                             0 Min. :
                         1st Qu.:
                                      5488
                                                         2567
                                                                 1st Qu.:
                                                                             3261
     1st Qu.:
                  5651
                                             1st Qu.:
                 30871
                                    27001
                                                        20749
                                                                 Median:
                                                                            26067
     Median:
                         Median :
                                             Median :
                229986
                               : 253312
                                                   : 233072
                                                                      : 261212
                         Mean
                                             Mean
                                                                 Mean
                                             3rd Qu.: 114138
     3rd Qu.: 117475
                         3rd Qu.: 113194
                                                                 3rd Qu.: 130208
            :16719906
                               :43928481
                                                   :15135169
                                                                 Max.
                                                                      :19631820
                         Max.
                                             Max.
        T49_1_MSS
                           T49_2_MSS
                                                M42_1_PD
                                                                    M42_2PD
                                             Min.
                                                                 Min.
                         Min.
                                                             0
     1st Qu.:
                  9293
                         1st Qu.:
                                      8607
                                             1st Qu.:
                                                         5393
                                                                 1st Qu.:
                                                                             4214
     Median :
                55654
                         Median :
                                    46397
                                             Median :
                                                        36887
                                                                 Median :
                                                                            30597
           : 542800
                         Mean : 462909
                                                  : 388693
                                                                 Mean : 333813
                                             Mean
      3rd Qu.: 223267
                         3rd Qu.: 189197
                                             3rd Qu.: 180508
                                                                 3rd Qu.: 152696
            :49218872
                               :29240206
                                                  :48177680
                                                                 Max. :42558111
                         Max.
        M43 1 PD
                            M43_2_PD
                                                M64_1_PD
                                                                    M64_2_PD
                         Min.
                                                                 Min.
     1st Qu.:
                19633
                         1st Qu.:
                                    17226
                                             1st Qu.: 11037
                                                                 1st Qu.:
                                                                             8655
     Median :
                 67737
                                    59598
                                                        52310
                                                                 Median :
                                                                            47454
                         Median :
                                             Median :
                                                   : 470967
           : 349175
                               : 358976
                                                                 Mean
                                                                      : 485038
                         Mean
                                             Mean
      3rd Qu.: 205615
                         3rd Qu.: 201931
                                             3rd Qu.: 210268
                                                                 3rd Qu.: 206426
            :35049402
                         Max.
                                :63082982
                                             Max.
                                                    :71750330
                                                                 Max.
                                                                        :88912734
 groupColors <- c(rep("red", 6), rep("blue", 6)) # Coloreamos por metodos de enriquecimiento de phosphopeptidos.
 boxplot(matrix, col=groupColors, main="Expression values of each sample",
      xlab="Samples",
     ylab="Expression", las=2, cex.axis=0.7, cex.main=0.7)
                                   Expression values of each sample
                                                                              0
    8e+07
                                                                        0
                                                                   0
    6e+07
                                                                              0
 Expression
                                                                              0
                     0
                                                       0
    4e+07
                                                                              0
    2e+07
   0e+00
Aqui es dificil ver la media de valores asi que vamos a transformar los valores a log2.
 logM <-log2(matrix + 1) # sumamos 1 porque los valores son muy bajos y asi evitamos valores proximos a 0.
 groupColors <- c(rep("red", 6), rep("blue", 6)) # Coloreamos por metodos de enriquecimiento de phosphopeptidos.
 boxplot(logM, col=groupColors, main="Expression values of each sample",
     xlab="Samples",
     ylab="Expression", las=2, cex.axis=0.7, cex.main=0.7)
                                   Expression values of each sample
      25
      20
 Expression
      15
      10
                                           4
                                                 4
       5
                           \circ
                                      0
                     0
                                                       0
                                                                        0
                                                                              0
                                M5_2_MSS
                                           Samples
Aqui podemos ver que los duplicados se parecen entre ellos lo cual es lo ideal.
 pcX<-prcomp(t(logM), scale=FALSE) # Ya se han escalado los datos</pre>
 loads<- round(pcX$sdev^2/sum(pcX$sdev^2)*100,1)</pre>
 # Then plot the first two components.
 xlab<-c(paste("PC1",loads[1],"%"))</pre>
 ylab<-c(paste("PC2",loads[2],"%"))</pre>
 plot(pcX$x[,1:2],xlab=xlab,ylab=ylab, col=groupColors,
      main ="Principal components (PCA)")
 #names2plot<-paste0(substr(names(matrix),1,3), 1:6)</pre>
 names2plot <- colnames(matrix)</pre>
 text(pcX$x[,1],pcX$x[,2],names2plot, pos=2, cex=.6)
                              Principal components (PCA)
                  M1_2_MSS O
M1_1_MSS O
      50
          ss o<sup>M5_2_MSS</sup> O
                       T49_2_MSS O
                             T49_1_MSS O
PC2 17.5 %
      0
                                     M64_2_PD O
                                          M64_1_PD O
      -50
                      M42_2_PD 🔾
                          M42_1_PD 🔾
                                                       50
                                                                       100
                       -50
                                          PC1 32 %
Aqui podemos observar que las muestras duplicadas no se replican muy bien y que que como mucho se separan las muestras por PD y MSS
que son dos metodos distintos de enriquecer los fosfopeptidos.
 library(pheatmap)
 logM <-log2(matrix + 1)</pre>
 # Crear un heatmap con etiquetas, colores personalizados y valores mostrados
 heatmap_result <- pheatmap(</pre>
   logM,
   color = colorRampPalette(c("blue", "white", "red"))(50), # Gradiente de colores
   cluster rows = TRUE,
                              # Agrupar genes (filas)
   cluster_cols = TRUE,
                              # Agrupar muestras/condiciones (columnas)
   show_rownames = FALSE,
                             # Mostrar nombres de genes
                               # Mostrar nombres de muestras
   show_colnames = TRUE,
 print(heatmap_result)
                                                                                      25
                                                                                       20
                                                                                       15
                                                                                       10
                                                       M5_2_MSS
Podemos observar que se agrupan el metodo de enriquecimiento de peptidos y por duplicados de muestra. Los heatmaps son muy utiles para
analizar patrones mas profundamente.
```

aescos_omics_PEC1

2024-11-04

<<<<< HEAD