

PEC 1: Análisis de Datos Ómicos

Daniel Camacho Montaña

2024-11-01

Contents

1. Introducción	1
2. Lectura de los datos	2
2.1 Preprocesado de los datos	4
2.2 Guardado de datos. Creación del repositorio de GitHub.	5
3. Análisis de los datos	6
3.1 Estudio de la distribución	6

1. Introducción

El presente análisis se basa en un conjunto de datos de fosfoproteómica obtenido a partir de un experimento que investigó dos subtipos tumorales diferentes utilizando modelos de xenoinjertos derivados de pacientes (PDX). Las muestras fueron enriquecidas con fosfopéptidos y luego analizadas mediante espectrometría de masas acoplada a cromatografía líquida (LC-MS) en condiciones de duplicados técnicos, lo que permitió medir las abundancias normalizadas de señales MS de aproximadamente 1,400 fosfopéptidos.

El objetivo principal de este análisis es identificar fosfopéptidos que puedan diferenciar entre los dos grupos tumorales mediante el uso de métodos estadísticos y visualizaciones gráficas. Los datos proporcionados están organizados en un archivo de Excel, denominado TIO2+PTYR-human-MSS+MSIvsPD.XLSX, donde la columna SequenceModifications contiene las modificaciones y sus ubicaciones en los péptidos.

Los dos grupos tumorales se definen como:

- Grupo MSS: Incluye las muestras M1, M5 y T49.
- Grupo PD: Incluye las muestras M42, M43 y M64.

Antes de empezar a trabajar, tendremos que crear un directorio git, en el que guardaremos todos los datos que vayamos generando. Los repositorios de gitHub ofrece múltiples beneficios para el desarrollo de software, ya que facilita el control de versiones, permitiendo realizar un seguimiento de cambios y colaborar con otros de manera eficiente. Su interfaz intuitiva y herramientas de gestión de proyectos mejoran la comunicación entre equipos.

En gitHub, podemos crear directamente el nuevo repositorio, al que llamaremos “Camacho-Montana-Daniel-PEC1”, y el cual estará vinculado con nuestra carpeta en el terminal local, de manera que, al actualizar nuestros datos, podremos actualizar el repositorio al momento.

Para crearlo, primero tendremos que inicializar el repositorio (`git init`), y configuraremos nuestros datos del directorio de gitHUB.

```
system("git init") # Inicializa un nuevo repositorio de Git
```

```
## [1] 0
```

```
system("git config --global user.name 'Daniel Camacho'")
```

```
## [1] 0
```

```
system("git config --global user.email 'dcamacmon@uoc.edu'")
```

```
## [1] 0
```

Las respuestas de 0, indican que no ha habido ningún problema al generar dichas configuraciones. A continuación, tendremos que vincular el repositorio local con el repositorio remoto.

```
repo_url <- "https://github.com/dcamacmon/Camacho-Monta-o-Daniel-PEC1"  
system(paste("git remote add origin", repo_url))
```

```
## [1] 3
```

Con la respuesta afirmativa, podemos aceptar que el directorio remoto se ha vinculado correctamente al local, por lo que podemos empezar a trabajar.

2. Lectura de los datos

Para iniciar con el estudio de los datos, primero debemos cargar la información del dataset de metabolómica desde github. El archivo, al tratarse de XLSM, usaremos el paquete `readxl` y cargaremos la información directamente desde el repositorio remoto.

Como importamos los datos desde un directorio web, primero debemos crear un archivo temporal con la extensión `.xlsx` (`temp_file`) y descargar el archivo desde la URL (`download.file()`) y finalmente cargar los datos en nuestro data frame (`read_excel()`). Con

```
library(readxl)  
  
url <- paste0("https://github.com/nutrimetabolomics/metaboData/raw/main/",  
             "Datasets/2018-Phosphoproteomics/TIO2%2BPTYR-human-MSS%2BMSIvsPD.XLSX")  
  
temp_file <- tempfile(fileext = ".xlsm") #Creación del archivo temporal  
download.file(url, temp_file, mode = "wb") #Descarga del archivo desde el repositorio  
fosfo_data <- read_excel(temp_file) #Cargamos el documento
```

Con el archivo ya cargado, podemos ver su estructura con `str()`.

```
summary(fosfo_data) #Estructura de los datos
```

```
## SequenceModifications  Accession      Description      Score
## Length:1438           Length:1438     Length:1438      Min.   : 19.51
## Class :character       Class :character  Class :character  1st Qu.: 38.96
## Mode  :character       Mode  :character  Mode  :character  Median : 47.48
##                               Mean   : 51.30
##                               3rd Qu.: 60.06
##                               Max.   :132.19
##      M1_1_MSS           M1_2_MSS           M5_1_MSS           M5_2_MSS
## Min.   :      0      Min.   :      0      Min.   :      0      Min.   :      0
## 1st Qu.:    5653     1st Qu.:    5497     1st Qu.:    2573     1st Qu.:    3273
## Median :   30682     Median :   26980     Median :   20801     Median :   26241
## Mean   :  229841     Mean   :  253151     Mean   :  232967     Mean   :  261067
## 3rd Qu.:  117373     3rd Qu.:  113004     3rd Qu.:  113958     3rd Qu.:  130132
## Max.   :16719906     Max.   :43928481     Max.   :15135169     Max.   :19631820
##      T49_1_MSS           T49_2_MSS           M42_1_PD           M42_2_PD
## Min.   :      0      Min.   :      0      Min.   :      0      Min.   :      0
## 1st Qu.:    9306     1st Qu.:    8611     1st Qu.:    5341     1st Qu.:    4216
## Median :   55641     Median :   46110     Median :   36854     Median :   30533
## Mean   :   542449     Mean   :  462616     Mean   :  388424     Mean   :  333587
## 3rd Qu.:  223103     3rd Qu.:  189141     3rd Qu.:  180252     3rd Qu.:  152088
## Max.   :49218872     Max.   :29240206     Max.   :48177680     Max.   :42558111
##      M43_1_PD           M43_2_PD           M64_1_PD           M64_2_PD
## Min.   :      0      Min.   :      0      Min.   :      0      Min.   :      0
## 1st Qu.:   19641     1st Qu.:   17299     1st Qu.:   11038     1st Qu.:    8660
## Median :   67945     Median :   59607     Median :   52249     Median :   47330
## Mean   :   349020     Mean   :  358822     Mean   :  470655     Mean   :  484712
## 3rd Qu.:  205471     3rd Qu.:  201924     3rd Qu.:  209896     3rd Qu.:  206036
## Max.   :35049402     Max.   :63082982     Max.   :71750330     Max.   :88912734
##      CLASS              PHOSPHO
## Length:1438           Length:1438
## Class :character       Class :character
## Mode  :character       Mode  :character
##
##
##
```

En un análisis preeliminar, vemos que hay 18 columnas de datos.

- La primera, contiene el tipo de las modificaciones en la secuencia, así como la ubicación de la modificación (la columna de identificación).
- La tercera columna, tenemos la descripción de la secuencia analizada.
- La cuarta columna tenemos el Score, que evalúa la confiabilidad de la identificación de una secuencia peptídica.
- Las siguientes 12 columnas, tratan los datos obtenidos de LC-SM, de los grupos tumorales con las dos réplicas técnicas.
- Las dos últimas columnas contienen dos variables, CLASS y PHOSPHO, de tipo categóricas

Como las dos últimas variables son categóricas, las convertiremos a tipo factor para posibles análisis futuros.

```
fosfo_data$CLASS<-(as.factor(fosfo_data$CLASS))
str(fosfo_data$CLASS)
```

```
## Factor w/ 2 levels "C","H": 2 2 2 2 2 2 2 2 2 ...
```

```
fosfo_data$PHOSPHO<-(as.factor(fosfo_data$PHOSPHO))
str(fosfo_data$PHOSPHO)
```

```
## Factor w/ 2 levels "S/T","Y": 2 2 2 2 2 2 2 2 2 ...
```

2.1 Preprocesado de los datos

Una vez cargados los datos, crearemos un objeto de clase `SumarizedExperiment` (una extensión de `ExpressionSet`), ya que permiten el manejo y almacenamiento de datos de fosfoproteómica, ya que puede organizar datos de abundancia junto con los metadatos relevantes para las filas (entradas). Gracias a este formato, podemos asociar metadatos a filas y columnas de manera más clara, y permite datos con réplicas (necesario para nuestro análisis).

Para crear el `SummarizedExperiment`, tendremos que cargar el paquete `BiocManager`

```
if (!requireNamespace("SummarizedExperiment", quietly = TRUE)) {
  BiocManager::install("SummarizedExperiment")
}
library(SummarizedExperiment)
```

Una vez tenemos el paquete cargado, debemos extraer los datos de “abundancia”, es decir, las lecturas de cada entrada, las cuales se ubican entre la columna 5 y la 16 (incluidas las dos).

Del mismo modo, creamos un dataframe para los metadatos, los cuales ya conocemos por el análisis preliminar (filas y columnas). Recordemos que tenemos 3 grupos por cada clase tumoral y tenemos 2 réplicas por cada grupo.

```
abundances <- fosfo_data[, 5:16]
row_data <- fosfo_data[, c("SequenceModifications", "Accession",
                           "Description", "Score", "CLASS", "PHOSPHO")]
col_metadata <- DataFrame(
  SampleID = colnames(abundances),
  Group = c(rep("MSS", 6), rep("PD", 6))
)
```

Una vez hemos pre-procesado los datos para que tengan el formato adecuado, podemos crear el `SummarizedExperiment`.

```
se <- SummarizedExperiment(
  assays = list(counts = as.matrix(abundances)), # Datos de abundancia como matriz
  rowData = row_data, # Metadatos para las filas
  colData = col_metadata # Metadatos para las columnas
)
```

2.2 Guardado de datos. Creación del repositorio de GitHub.

Ahora que tenemos el SE, lo exportaremos en formato .RDA para poder incluirlo en nuestro repositorio de github. Del mismo modo, extraeremos los datos de expresión (assay data), los metadatos de las muestras (column data) y los metadatos de las variables (row data) en formato de texto.

```
save(se, file = "Fosfodatos_se.Rda") #Guardamos el SE en formato .Rda

counts<-assay(se,"counts") #Obtenemos los datos de assay data
write.table(counts, file = "datos_assay.txt", sep = "\t",
            quote = FALSE, row.names = TRUE, col.names = TRUE) #Obtenemos el .txt

sample_metadata <- colData(se) #Obtenemos los metadatos de las muestras
write.table(sample_metadata, file = "metadatos_muestras.txt", sep = "\t",
            quote = FALSE, row.names = TRUE, col.names = TRUE) #Se guardan en un .txt

gene_metadata <- rowData(se) #Obtenemos los metadatos de las variables
write.table(gene_metadata, file = "metadatos_genes.txt", sep = "\t",
            quote = FALSE, row.names = TRUE, col.names = TRUE) #Se guardan en un .txt
```

En ocasiones, conviene mantener los metadatos por separado, por lo que crearemos un archivo que contenga, en formato .md, los metadatos de nuestro dataset.

```
markdown_content <- "
# Metadatos del Dataset de Fosfoproteómica

## Descripción General
Este documento describe los metadatos asociados al dataset de fosfoproteómica utilizado en
el análisis.El dataset incluye información sobre las modificaciones post-traduccionales de
secuencias peptídicas en diferentes grupos tumorales.

## Estructura del Dataset
El dataset consta de las siguientes columnas:
```

Columna	Descripción
`SequenceModifications`	Tipo de modificaciones en la secuencia y ubicación.
`Accession`	Identificador de la secuencia.
`Description`	Descripción de la secuencia analizada.
`Score`	Score que evalúa la confiabilidad de la identificación.
`CLASS`	Clase de la muestra (H o C).
`PHOSPHO`	Aminoácido donde ocurre la fosforilación (Y, S/T).
`Group`	Grupo tumoral (MSS o PD).

```
## Información sobre los Grupos
- **Grupo MSS**: Incluye las muestras M1, M5 y T49.
- **Grupo PD**: Incluye las muestras M42, M43 y M64.

## Consideraciones
- Las columnas `CLASS` y `PHOSPHO` son categóricas y pueden ser convertidas a factores para
análisis estadísticos.
"
```

```
writeLines(markdown_content, "metadatos_dataset.md")
```

Toda esta información, el objeto SE, así como los metadatos y el resto de datos asociados al SE, se tendrán que almacenar en un repositorio de GitHub. Con el repositorio local previamente creado, podemos cargar los datos previamente guardados. Primero tendremos que agregar los archivos al directorio git. En caso de que los archivos ya están guardados, realizamos un commit, que implica “cargar” las modificaciones que hemos generado en los ficheros, para finalmente “empujar” (push) hacia el directorio remoto.

```
# Agregar todos los archivos  
system("git add .")
```

```
## [1] 0
```

```
# Hacer commit de los cambios  
system("git commit -m 'Añadir_archivos'")
```

```
## [1] 0
```

```
# Hacer push para subir los cambios a GitHub  
system("git push origin main") # Asegúrate de que la rama sea correcta
```

```
## [1] 0
```

El mensaje de retorno [1] 0 indica que los comandos se ejecutaron correctamente, lo que significa que todos los cambios en el directorio de trabajo han sido añadidos al área de preparación (staging area) de Git.

3. Análisis de los datos

3.1 Estudio de la distribución

Con el SE, podremos visualizar los datos para entender su distribución y variabilidad. Tenemos dos grupos principales (las clases tumorales), pero dentro de cada grupo podemos clasificar 3 muestras en cada uno.

Como tenemos dos réplicas de cada muestra, debemos analizar si la varianza es significativa entre ambas lecturas, ya que dicha diferencia podría estar causada por un error de la técnica o la recolección de datos. Usaremos el paquete `cary` analizaremos M1 (M1_1 y M1_2).

```
library(car)  
abundancias_M1 <- abundances[, c("M1_1_MSS", "M1_2_MSS")]
```

Crearemos un nuevo data frame con solo los datos de M1, con los datos de la abundancia (tanto de la réplica 1 y la réplica 2), y con la información de qué replica como factor de dos niveles. A continuación, visualizaremos las 6 primeras observaciones, así como la estructura de los datos.

```
# Crear un dataframe con los nombres de las columnas adecuados  
datos_M1 <- data.frame(  
  Abundancia = c(abundancias_M1[["M1_1_MSS"]], abundancias_M1[["M1_2_MSS"]]),  
  Replica = factor(rep(c("Replica 1", "Replica 2"), each = nrow(abundancias_M1)))  
)  
head(datos_M1)
```

```
##      Abundancia  Replica
## 1      24.29438 Replica 1
## 2       0.00000 Replica 1
## 3     3412.60332 Replica 1
## 4    220431.17880 Replica 1
## 5     18254.77813 Replica 1
## 6    644513.31840 Replica 1
```

```
str(datos_M1)
```

```
## 'data.frame': 2876 obs. of 2 variables:
## $ Abundancia: num 24.3 0 3412.6 220431.2 18254.8 ...
## $ Replica : Factor w/ 2 levels "Replica 1","Replica 2": 1 1 1 1 1 1 1 1 1 1 ...
```

Finalmente, realizaremos el test de varianzas levene entre las dos réplicas.

```
resultados_levene <- leveneTest(Abundancia ~ Replica, data = datos_M1)
print(resultados_levene)
```

```
## Levene's Test for Homogeneity of Variance (center = median)
##      Df F value Pr(>F)
## group 1 0.2676 0.605
##      2874
```

Vemos que el test de homogeneidad de varianzas acepta la hipótesis nula, por lo que podemos asumir que no hay una variabilidad significativa entre ambas lecturas.

Para obtener el resultado de test de Levene para todas las réplicas, primero crearemos un nuevo dataframe vacío que contendrá el resultado del test, tanto el p-valor como el identificador de la muestra.

```
abundance_data<-assay(se, "counts")
muestras <- c("M1", "M5", "T49", "M42", "M43", "M64") # Obtener los nombres de las columnas para las réplicas

resultados_levene <- data.frame(Muestra = character(),
                                p_valor = numeric(),
                                stringsAsFactors = FALSE)
```

Seguidamente, iteraremos las muestras en que el sufijo “_MSS” o “_PD” indicará el grupo, y el sufijo “_1” o “_2” indicará la réplica. Entonces, ejecutaremos un bucle `if else`, en que se verifica si existe una columna MSS (if, col1 y col2) o PD (else, col3 o col4) y evaluará los datos de la varianza de la réplica 1 (_1) respecto a la réplica 2 (_2) con el test de Levene. Finalmente, guardará los datos en el dataframe que hemos creado anteriormente (vacío hasta el momento).

```
# Empezamos por iterar sobre las muestras
for (muestra in muestras) {
  # Definimos los nombres de las columnas de las réplicas de MSS
  col1 <- paste0(muestra, "_1_MSS")
  col2 <- paste0(muestra, "_2_MSS")

  # Definimos los nombres de las columnas de las réplicas de PD
  col3 <- paste0(muestra, "_1_PD")
  col4 <- paste0(muestra, "_2_PD")
```

```

# Verificamos si ambas columnas MSS existen
if (col1 %in% colnames(abundance_data) & col2 %in% colnames(abundance_data)) {

  abundancias <- abundance_data[, c(col1, col2)] # Extraemos las abundancias de las dos réplicas MSS
  # Crear un dataframe para el test de Levene
  datos <- data.frame(
    Abundancia = c(abundancias[, 1], abundancias[, 2]),
    Replica = factor(rep(c("Replica 1", "Replica 2"), each = nrow(abundancias)))
  )
  # Realizar el test de Levene
  resultado_levene <- leveneTest(Abundancia ~ Replica, data = datos)
  # Guardar el p-valor en la tabla de resultados
  resultados_levene <- rbind(resultados_levene,
                             data.frame(Muestra = muestra,
                                           p_valor = resultado_levene$`Pr(>F)`[1]))
  # En caso de no ser MSS, será PD, y repetiremos el proceso para el resto de muestras PD
} else if (col3 %in% colnames(abundance_data) & col4 %in% colnames(abundance_data))
{
  abundancias <- abundance_data[, c(col3, col4)]
  datos <- data.frame(
    Abundancia = c(abundancias[, 1], abundancias[, 2]),
    Replica = factor(rep(c("Replica 1", "Replica 2"), each = nrow(abundancias)))
  )
  resultado_levene <- leveneTest(Abundancia ~ Replica, data = datos)
  resultados_levene <- rbind(resultados_levene,
                             data.frame(Muestra = muestra,
                                           p_valor = resultado_levene$`Pr(>F)`[1]))
}
}

```

Finalmente, mostramos el resultado de los test de Levene.

```
print(resultados_levene)
```

```
##  Muestra  p_valor
## 1      M1 0.6049833
## 2      M5 0.4988294
## 3     T49 0.3808396
## 4     M42 0.4179607
## 5     M43 0.8632354
## 6     M64 0.8864896

```

Vemos entonces, que no hay diferencias significativas en la variabilidad entre las réplicas de las muestras de los dos grupos. Esto implica, que podemos promediar las réplicas, o mantenerlas independientes para los futuros análisis. En este caso, debido a que en los estudios de LC-MS la variabilidad puede mostrar diferencias biológicas subyacentes en las muestras, mantendremos cada réplica como una muestra independiente.

A continuación, mostraremos la distribución de las muestras de la matriz de abundancia.

Al tratar con 12 columnas de datos al mismo tiempo, numéricamente es difícil determinar algún patrón en ellos. Para evaluarlo más gráficamente, realizaremos un boxplot con estos mismos datos.

Primero tendremos que obtener la matriz de abundancias (counts) desde el SE, y convertirla en una matriz de formato largo (usando la función `melt()`) y modificamos los títulos para facilitar su lectura.


```
library(reshape2)
```

```
abundance_data <- assays(se)$counts  
abundance_long <- melt(abundance_data)  
colnames(abundance_long) <- c("Fosfopéptido", "Muestra", "Abundancia")
```

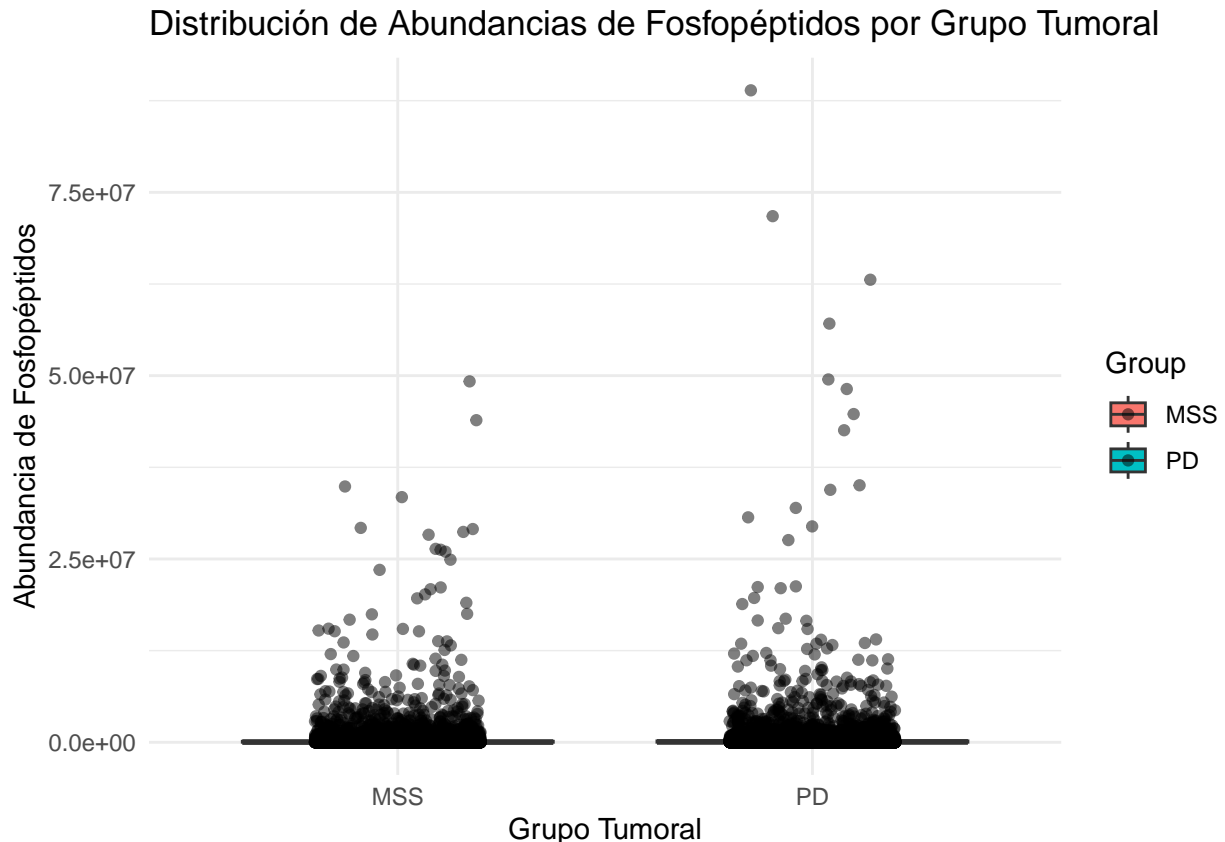
Después, dividiremos los dos grupos tumorales (MSS y PD).

```
grupos <- as.data.frame(colData(se)$Group)  
colnames(grupos) <- "Group"  
abundance_long$Group <- rep(grupos$Group, each = nrow(abundance_data))
```

Finalmente, crearemos un gráfico con los dos grupos

```
library(ggplot2)
```

```
ggplot(abundance_long, aes(x = Group, y = as.numeric(Abundancia), fill = Group)) +  
  geom_boxplot(outlier.shape = NA) + # Oculta los outliers si hay muchos puntos  
  geom_jitter(width = 0.2, alpha = 0.5, color = "black") + # Añade puntos de dispersión  
  labs(title = "Distribución de Abundancias de Fosfopéptidos por Grupo Tumoral",  
        x = "Grupo Tumoral",  
        y = "Abundancia de Fosfopéptidos") +  
  theme_minimal()
```



Vemos que, de entre los dos grupos tumorales, parece ser que el grupo PD tiene una abundancia superior (de hecho, es bastante similar, pero parece que el PD tiene 3 valores mucho más altos).

Para poder ver los datos estadísticos descriptivos de ambos grupos, usaremos el paquete `dplyr`, creando un dataframe con la media, la mediana, la desviación estándar y el ratio intercuartílico.

```
library(dplyr)
```

```
abundance_stats <- abundance_long %>%
  group_by(Group) %>%
  summarise(
    Media = mean(Abundancia, na.rm = TRUE),
    Mediana = median(Abundancia, na.rm = TRUE),
    SD = sd(Abundancia, na.rm = TRUE),
    IQR = IQR(Abundancia, na.rm = TRUE)
  )
print(abundance_stats)
```

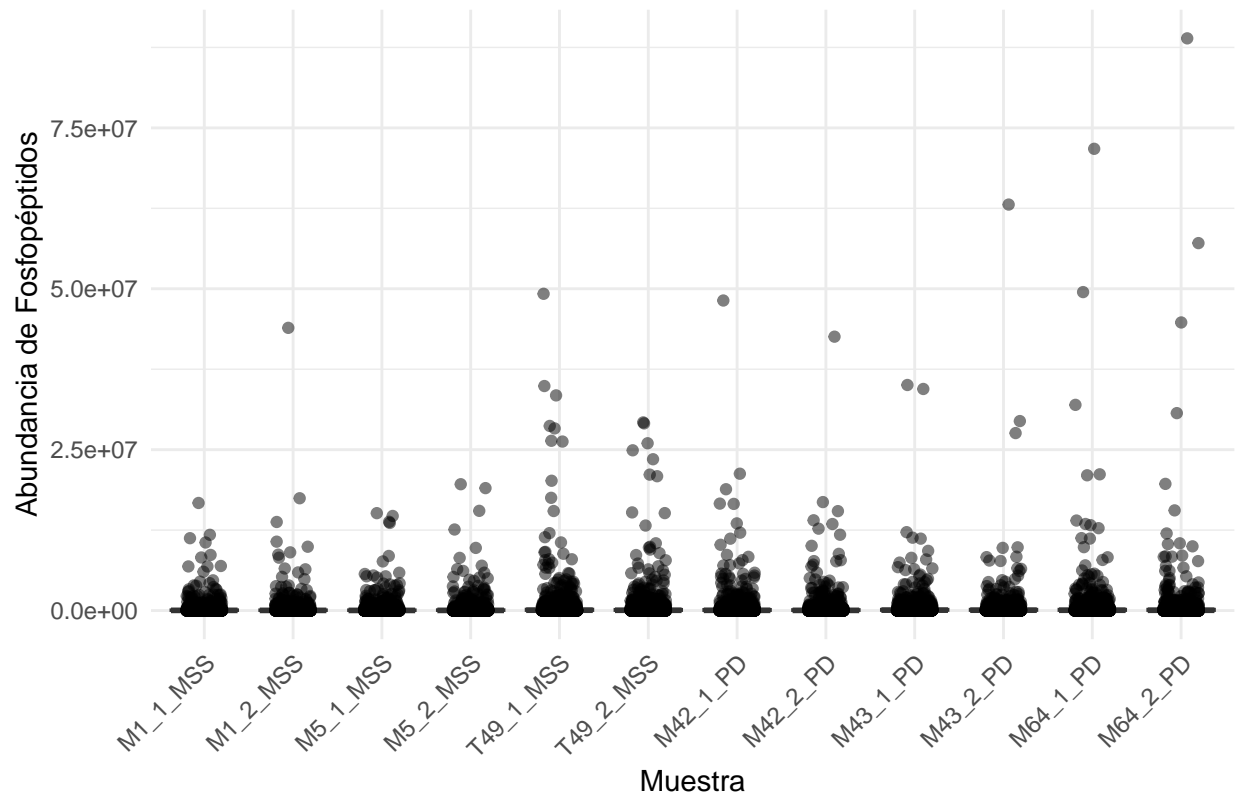
```
## # A tibble: 2 x 5
##   Group   Media Mediana      SD      IQR
##   <chr>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 MSS    330349.   31944. 1670041. 139913.
## 2 PD     397537.   48093. 2308806. 182955.
```

Como ya veíamos en el gráfico original de los grupos tumorales, la media y la mediana de PD es superior a la del grupo MSS. Del mismo modo, la desviación estándar es superior en el grupo PD, por lo que los valores están más dispersos respecto a la media (tiene una mayor variabilidad), y el IQR superior indica que la mitad central de los datos es más dispersa (hay más diferencia entre el Q1 y el Q3).

Como ya tenemos los datos de la muestra en `abundance_long`, solo demos crear el boxplot en base a esos datos, con las muestras (MSS y PD) en base a su abundancia. Además, como tenemos 12 muestras, inclinamos las etiquetas del eje X para facilitar su lectura.

```
ggplot(abundance_long, aes(x = Muestra, y = Abundancia)) +
  geom_boxplot(aes(fill = Group), outlier.shape = NA) + # Oculta los outliers
  geom_jitter(width = 0.2, alpha = 0.5, color = "black") + # Añade puntos
  labs(title = "Distribución de Abundancias de Fosfopéptidos por Muestra",
        x = "Muestra",
        y = "Abundancia de Fosfopéptidos") +
  theme_minimal() +
  theme(axis.text.x = element_text(angle = 45, hjust = 1), # Rotar etiquetas en el eje x
        legend.position = "none") # Eliminar la leyenda
```

Distribución de Abundancias de Fosfopéptidos por Muestra



Como sabemos, las muestras M1, M5 y T49 son del grupo MSS, mientras que las muestras M42, M43 y M64 son del grupo PD. En el gráfico, vemos que las muestras de MSS tienen valores inferiores, excepto la muestra T49, que parece tener unos valores más altos. En cambio, en PD, vemos que todas las muestras parecen tener valores superiores a MSS, lo que explicaría la diferencia de medias.

Para profundizar con el análisis de los datos, evaluaremos los componentes principales (PCA), ya que nos pueden proporcionar datos relevantes, así como permitirnos identificar outliers que desvían el patrón general de los datos. Un paso previo a realizar un análisis de PCA es normalizar los datos. En este caso, obviaremos este proceso, ya que como se explica en la introducción, estas lecturas ya han sido normalizadas previamente.

```
# Extraer los conteos y calcular la varianza de todos los datos de abundancia
counts <- assay(se, "counts")
filtered_counts <- counts[, apply(counts, 2, var) > 0]
fosfo_variances <- apply(filtered_counts, 1, var)

# Filtramos los genes con baja varianza

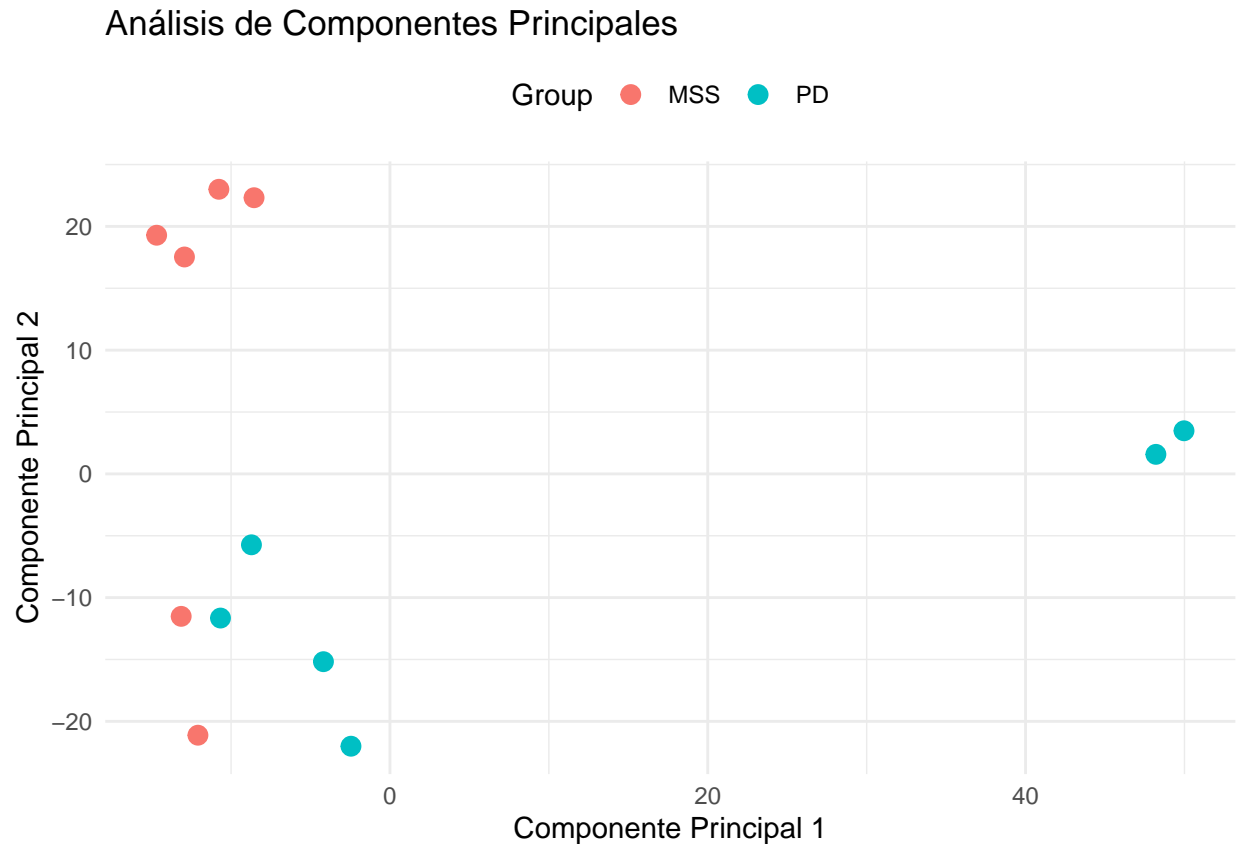
filtered_counts <- filtered_counts[fosfo_variances > 1e-2, ] # Ajustamos el umbral a 0.01

# Realizamos el PCA y creamos el dataframe
pca_result <- prcomp(t(filtered_counts), center = TRUE, scale. = TRUE)
pca_data <- as.data.frame(pca_result$x)

# Agregamos información de grupo
pca_data$Group <- colData(se)$Group

# Visualizamos los resultados de PCA
```

```
ggplot(pca_data, aes(x = PC1, y = PC2, color = Group)) +
  geom_point(size = 3) +
  labs(title = "Análisis de Componentes Principales",
       x = "Componente Principal 1",
       y = "Componente Principal 2") +
  theme_minimal() +
  theme(legend.position = "top")
```



Vemos que, en el caso del grupo MSS, se ve mucho más representado por la componente principal 2, mientras que el grupo PD se ve mucho más representado por la componente principal 1. Aun así, debemos tener en cuenta que cada componente principal representa cierta cantidad de la variabilidad de los datos.

Para determinar que grado de variabilidad se representa con cada PCA, reuniremos los datos en un nuevo dataframe. Como ya calculamos los PCA anteriormente, obtendremos primero las varianzas de cada componente (`variances`), así como la proporción de varianza explicada.

```
# Ver la proporción de varianza explicada
variances <- pca_result$sdev^2 # Varianzas de cada componente
explained_variance <- variances / sum(variances)*100 # Proporción de varianza explicada

# Crear un dataframe para visualización
pca_summary <- data.frame(
  Component = paste0("PC", 1:length(explained_variance)),
  Variance = explained_variance
)
```

En el data frame creado, tenemos una columna con el nombre de las PC y con su varianza, pero no tenemos

el formato correcto (modificaremos el resultado para que se vea en base a 100 y no en notación científica) y no tenemos la varianza acumulada (crearemos la columna CumulativeVariance).

```
#Modificamos el formato de los datos de Variance, así como la nueva columna acumulada.
pca_summary$Variance<-format(pca_summary$Variance, nsmall = 4, scientific = FALSE)
pca_summary$Variance <- round(as.numeric(pca_summary$Variance), 4)
pca_summary$CumulativeVariance <- cumsum(pca_summary$Variance)

# Mostrar la proporción de varianza explicada
print(pca_summary)
```

##	Component	Variance	CumulativeVariance
## 1	PC1	37.5052	37.5052
## 2	PC2	20.1507	57.6559
## 3	PC3	14.6318	72.2877
## 4	PC4	9.2421	81.5298
## 5	PC5	7.2003	88.7301
## 6	PC6	3.8196	92.5497
## 7	PC7	2.2596	94.8093
## 8	PC8	1.6962	96.5055
## 9	PC9	1.4568	97.9623
## 10	PC10	1.0729	99.0352
## 11	PC11	0.9648	100.0000
## 12	PC12	0.0000	100.0000

Vemos, entonces, que las dos componentes principales contienen el 57.6559% de la variabilidad total. Esta variabilidad es importante, y usar únicamente 2 CP podría simplificar el análisis. Aún así, usaremos 3 CP, que contendrán el 72.2877% de la variabilidad total.

De todos modos, aún conociendo la variabilidad asociada a los CP, desconocemos que fosfopéptidos influyen más en cada CP. En **rotation** de PCA, podemos ver que carga supone sobre las CP (nos centraremos en las 3 primeras, pero por separado)

Primero, obtendremos los resultados en **cargas** y crearemos un data frame con estos mismos datos.

Posteriormente, seguiremos un proceso en que extraeremos los 50 valores con cargas más alta, reordenando por cada CP (en cada CP habrá diferentes fosfopéptidos que afecten con mayor intensidad). Finalmente crearemos un gráfico de barras en que veamos como afectan cada uno de esos 50 resultados a la CP.

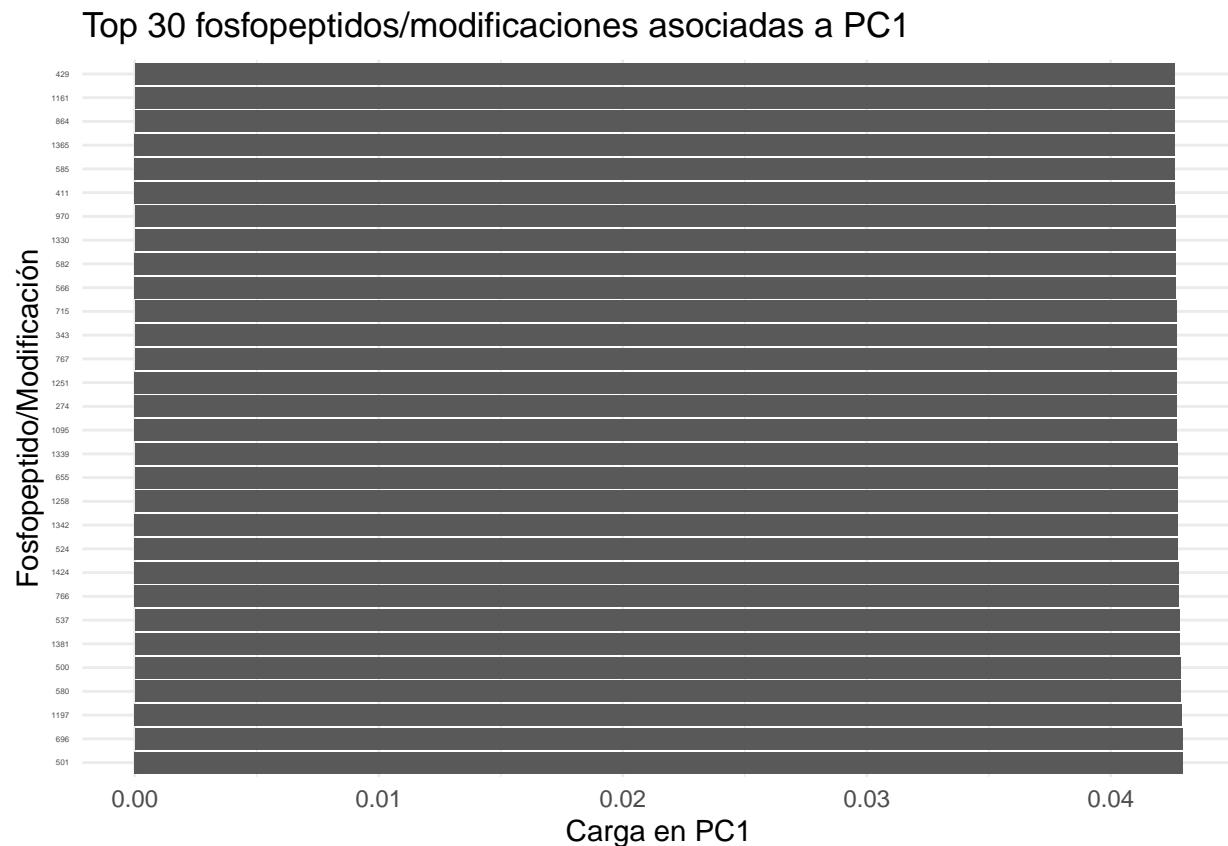
```
library(tibble)
```

```
#Obtenemos los datos de las cargas para las tres primeras CP
cargas <- pca_result$rotation
cargas_df <- as.data.frame(cargas[, c("PC1", "PC2", "PC3")])
cargas_df <- cargas_df %>%
  rownames_to_column(var = "Fosfopeptido")

#Reordenamos los datos en función de CP1

top_cargas_pc1 <- cargas_df %>%
  arrange(desc(PC1)) %>%
  head(30)
top_cargas_pc1$Fosfopeptido <- factor(top_cargas_pc1$Fosfopeptido,
  levels = top_cargas_pc1$Fosfopeptido)
```

```
ggplot(top_cargas_pc1, aes(x = Fosfopeptido, y = PC1)) +
  geom_bar(stat = "identity") +
  coord_flip() +
  labs(title = "Top 30 fosfopeptidos/modificaciones asociadas a PC1",
       x = "Fosfopeptido/Modificación", y = "Carga en PC1") +
  theme_minimal() +
  theme(axis.text.y = element_text(size = 3, angle = 0, hjust = 1)) # Tamaño y ángulo
```

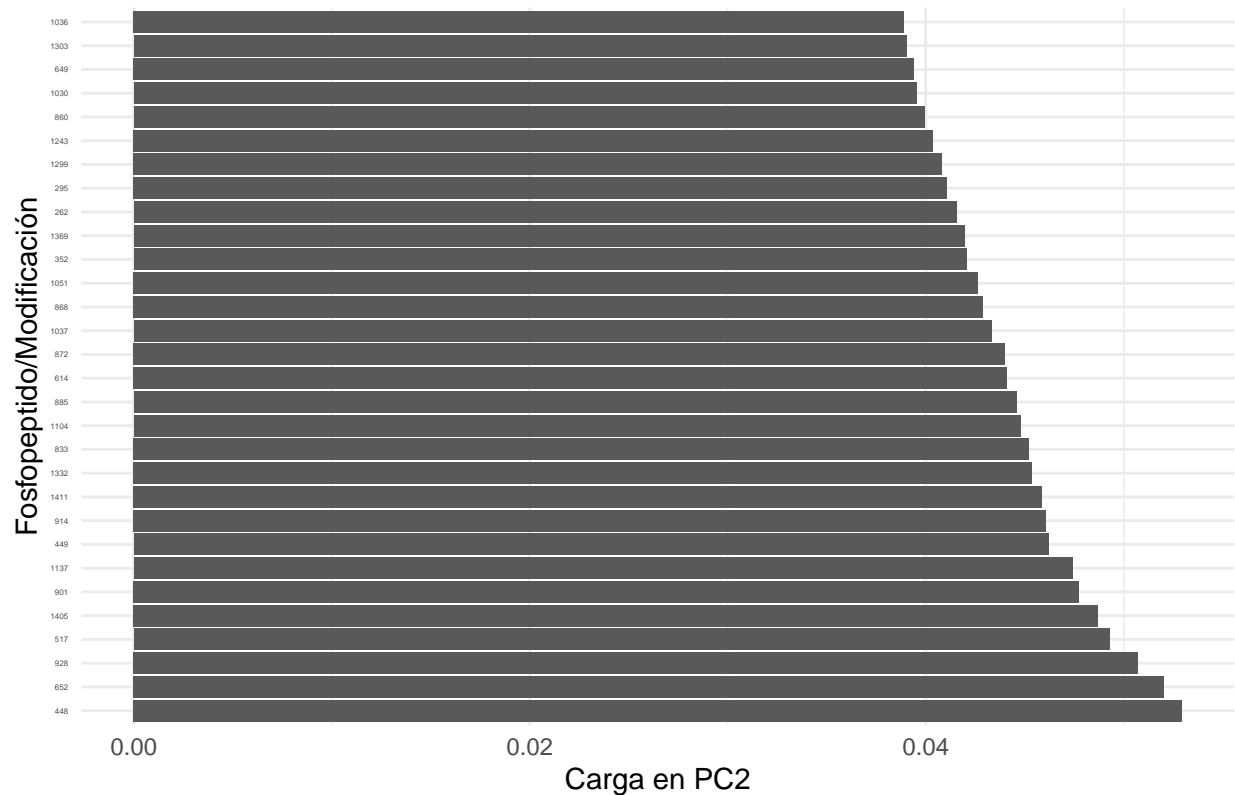


Realizaremos el mismo proceso para CP1:

```
top_cargas_pc2 <- cargas_df %>%
  arrange(desc(PC2)) %>%
  head(30)
top_cargas_pc2$Fosfopeptido <- factor(top_cargas_pc2$Fosfopeptido,
                                       levels = top_cargas_pc2$Fosfopeptido)

ggplot(top_cargas_pc2, aes(x = Fosfopeptido, y = PC2)) +
  geom_bar(stat = "identity") +
  coord_flip() +
  labs(title = "Top 30 fosfopeptidos/modificaciones asociadas a PC2",
       x = "Fosfopeptido/Modificación", y = "Carga en PC2") +
  theme_minimal() + # Tamaño y ángulo
  theme(axis.text.y = element_text(size = 3, angle = 0, hjust = 1))
```

Top 30 fosfopeptidos/modificaciones asociadas a PC2

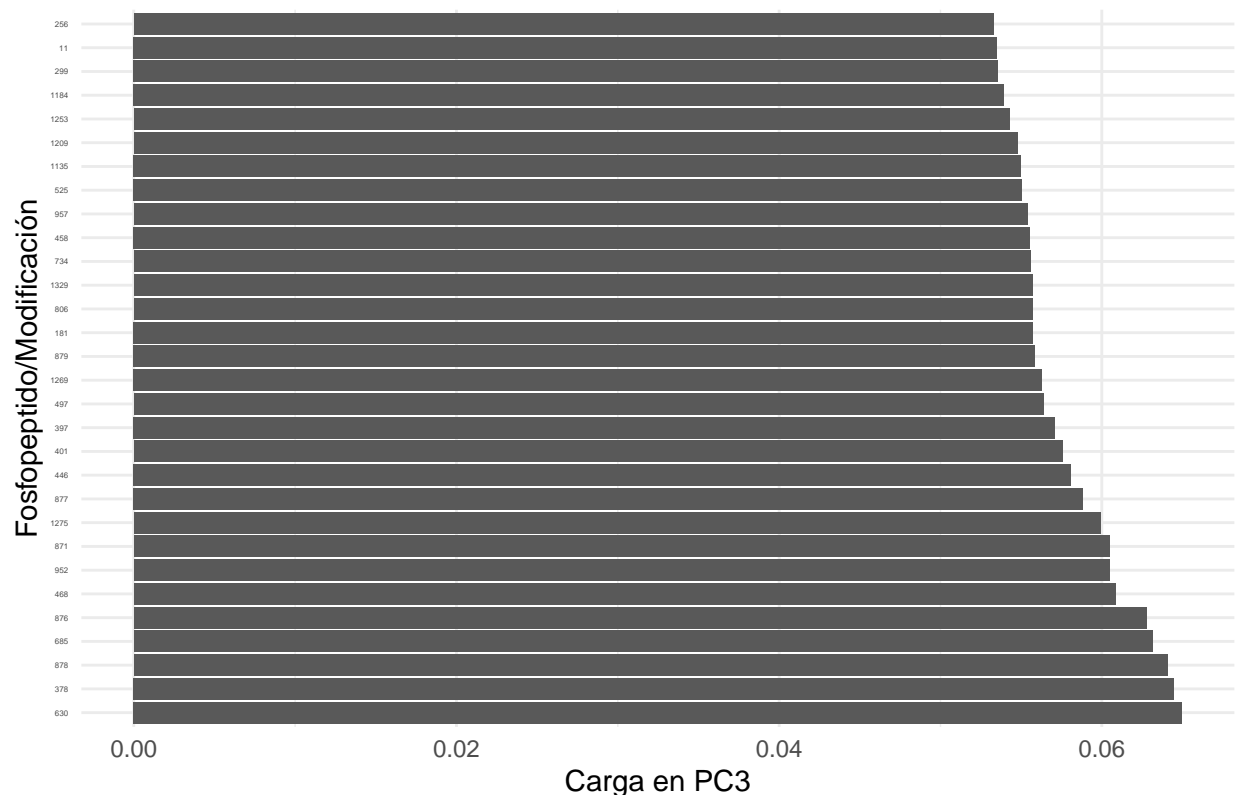


Y para PC3:

```
#Reordenamos los datos en función de CP3
top_cargas_pc3 <- cargas_df %>%
  arrange(desc(PC3)) %>%
  head(30)
top_cargas_pc3$Fosfopeptido <- factor(top_cargas_pc3$Fosfopeptido,
  levels = top_cargas_pc3$Fosfopeptido)

ggplot(top_cargas_pc3, aes(x = Fosfopeptido, y = PC3)) +
  geom_bar(stat = "identity") +
  coord_flip() +
  labs(title = "Top 30 fosfopeptidos/modificaciones asociadas a PC3",
    x = "Fosfopeptido/Modificación", y = "Carga en PC3") +
  theme_minimal() + # Tamaño y ángulo de las etiquetas
  theme(axis.text.y = element_text(size = 3, angle = 0, hjust = 1))
```

Top 30 fosfopeptidos/modificaciones asociadas a PC3



Como podemos ver, cada fosfopéptido tiene una influencia diferente en cada CP, aunque en la CP1 los 50 primeros parece que todos tienen una carga similar.

Para continuar con el análisis, realizaremos un ANOVA de un factor para ver si hay diferencias significativas en las abundancias de fosfopéptidos entre los grupos. Como los datos están estructurados con las variables (fosfopéptidos) en las filas y los grupos (MSS o PD) en las columnas, tendremos que transponer los datos de las abundancias.

```
abundancias <- assay(se)
grupos <- colData(se)$Group

#Para tener las lecturas en horizontal y los fosfopéptidos en vertical
abundancias_df <- as.data.frame(t(abundancias))
abundancias_df$grupo <- grupos
```

Además, debemos tener en cuenta que no podemos aplicar un ANOVA a cada fosfopéptido independientemente, por lo que usaremos sapply para repetir el proceso en todas las columnas (ahora la variable fosfopéptido) y guardaremos el resultado en anova_results.

```
# Función para aplicar ANOVA en cada fosfopéptido
anova_results <- apply(abundancias_df[, -ncol(abundancias_df)], 2, function(y) {
  aov_result <- aov(y ~ abundancias_df$grupo)
  summary(aov_result)[[1]]$Pr(>F)[1] # Extrae el valor p
})
```

Finalmente, crearemos un data frame con los resultados, con la información de los fosfopéptidos y filtraremos solo aquellos que tengan un p-valor significativo. Además,


```

anova_results <- data.frame(Fosfopeptido = names(anova_results),
                             p_value = anova_results)

anova_results$p_adjusted <- p.adjust(anova_results$p_value, method = "fdr")

significativos <- subset(anova_results, p_adjusted < 0.05)
summary(significativos)

```

```

## Fosfopeptido      p_value      p_adjusted
## Length:104      Min.       :7.100e-08    Min.       :0.0001019
## Class :character 1st Qu.:1.698e-04    1st Qu.:0.0091096
## Mode  :character Median :9.309e-04    Median :0.0254426
##                Mean  :1.296e-03    Mean  :0.0260587
##                3rd Qu.:2.360e-03    3rd Qu.:0.0433058
##                Max.  :3.525e-03    Max.  :0.0486788

```

Después del test anova en todos los fosfopéptidos vemos que, de los 1438 fosfopéptidos analizados, solo 104 muestran diferencias significativas entre los grupos MSS y PD.

Toda la documentación se ha recopilado en el siguiente repositorio de github: <https://github.com/dcamacmon/Camacho-Monta-o-Daniel-PEC1>