

PEC 1

Joan Canseco i Ribas

2024-10-30

Contents

Abstract PENDIENTE	1
Objetivos del estudio	2
Materiales y Métodos	2
Origen y naturaleza de los datos	2
Herramientas informáticas	2
Creación de un contenedor del tipo SummarizedExperiment	3
Exploración del dataset PENDIENTE	3
Creación del repositorio Github	3
Resultados	3
Creación del contenedor del tipo SummarizedExperiment	3
Exploración del dataset PENDIENTE	7
Discusión y limitaciones y conclusiones del estudio PENDIENTE	11

Abstract PENDIENTE

En este proyecto, se ha llevado a cabo un análisis estructurado de datos metabolómicos para simular un proceso de análisis ómico. El objetivo ha sido familiarizarse con la organización y manejo de datos ómicos, aplicando diversas herramientas y métodos que se han ido mostrando durante la asignatura de “Análisis de datos ómicos”. Para ello, se ha seleccionado un conjunto de datos de un estudio en el que se investiga los metabolitos en respuesta a la cirugía bariátrica en 4 distintos timepoints en 39 pacientes con más de 10 años de obesidad. Estos 39 pacientes están clasificados en dos grupos, un grupo de pacientes que se les aplica cirugía by pass (n=26) y otro que se les aplica cirugía tubular (n=13). Estos datos han sido organizados en un contenedor del tipo SummarizedExperiment para facilitar el trabajo con este dataset. Además, se han realizado análisis exploratorios en el que se han reflejado las características más relevantes del proyecto. También se han hecho ##### para verificar cuáles son los metabolitos que más cambian al aplicar una u otra cirugía. Finalmente, se ha creado un repositorio del proyecto en Github donde se han subido todos los documentos del proyecto.

Objetivos del estudio

Los objetivos generales de esta actividad son planificar y ejecutar una versión simplificada del proceso de análisis de datos ómicos, utilizando algunas de las herramientas y métodos que se han trabajado en la unidad. En concreto, se deberá:

- Seleccionar un dataset de un repositorio de github de metabolómica <https://github.com/nutrimetabolomics/metaboData/>.
- Crear un contenedor del tipo *SummarizedExperiment* que contenga los datos y metadatos del dataset.
- Llevar a cabo una exploración del dataset que dé una visión general del mismo.
- Elaborar un informe que describa el proceso realizado (este documento).
- Crear un repositorio en Github y subir el informe, el objeto contenedor con los datos y metadatos, el código R para la exploración de los datos, los datos y los metadatos.

Materiales y Métodos

Esta actividad se ha desarrollado desde un ordenador Windows.

Origen y naturaleza de los datos

Los datos analizados han sido extraídos del siguiente repositorio online: <https://github.com/nutrimetabolomics/metaboData/> Una vez accedemos a dicho repositorio, clicamos el botón verde “<> Code”, y clicamos la opción “Download ZIP”. Una vez se ha descargado el repositorio en nuestro ordenador, clicamos botón derecho y seleccionamos “Extraer todo”. Exploramos los distintos datasets y seleccionamos el dataset “2018-MetabotypingPaper” porque tiene los documentos bien estructurados y parece que se puede trabajar bien con él. Cogemos los documentos de dicho dataset y los llevamos a nuestro directorio de trabajo de la PEC1. Se trata de un dataset de un estudio de 39 pacientes con más de 10 años de obesidad a los que se les ha realizado una cirugía bariátrica. Se clasifica el dataset en 2 grupos según la cirugía bariátrica recibida: grupo “by pass” compuesto por 26 pacientes y grupo “tubular” compuesto por 13 pacientes. También se recoge su estado metabólico, en cuanto a si los pacientes estaban sanos o enfermos metabólicamente en la variable Group. En estos pacientes se pretende estudiar como evolucionan ciertos metabolitos en 4 distintos timepoints tras la cirugía según qué cirugía hayan tomado.

Herramientas informáticas

Se han usado principalmente 3 herramientas informáticas y bioinformáticas para la resolución de la PEC1:

- R y Rstudio: Se ha usado para realizar la exploración del dataset para obtener una visión general del mismo. Además, se ha usado para la redacción del informe, usando Rmarkdown.
- Bioconductor: Aunque es un paquete de R, lo consideraremos como una herramienta independiente para hacer esta explicación. Bioconductor se ha usado para crear un contenedor del tipo *SummarizedExperiment*.
- Github: Se ha usado para la descarga de los datos (ya se ha indicado en la sección previa el procedimiento realizado), y se ha usado para presentar los resultados complementarios de la PEC1.

Creación de un contenedor del tipo SummarizedExperiment

Para la creación del contenedor se ha usado Bioconductor dentro de R (como se ha especificado anteriormente). Los pasos se han especificado en el apartado de Resultados.

Exploración del dataset PENDIENTE

Para la exploración de los datos se usará principalmente R y Rstudio. Se usará en este proyecto dos documentos:

- “DataValues_S013.csv”: Se trata del documento donde hay almacenados los datos experimentales de los pacientes. Es un data frame con 39 filas (de 39 pacientes) y 696 variables. De las 696 variables, las 6 primeras son para identificar a los pacientes (X.1, SUBJECTS, SURGERY, AGE, GENDER y Group, en la que X.1 es un duplicado de SUBJECTS, por lo que realmente se podría considerar solo las 5 primeras columnas), y el resto de variables son los metabolitos tomados en 4 timepoints distintos (T0, T2, T4 y T5). Se denota a qué timepoint hace referencia cada variable porque al final del nombre de cada variable hay especificado el timepoint al que pertenecen con “_TX” (donde X es 0, 2, 4 o 5).
- “DataInfo_S013.csv”: Se trata de un documento donde hay la metadata del documento anterior y está organizado de manera distinta. Tiene solo 4 columnas (X, es un duplicado de VarName, VarName, contiene el nombre de cada una de las variables del documento anterior, varTpe, especifica el tipo de variable que es, y Description, especifica qué información se recoge en cada variable) y 695 filas, correspondientes a las 695 variables del documento anterior (recordemos que eran 696 pero la primera columna era un duplicado de la segunda columna, por lo que son 695 efectivas). Una exploración más exhaustiva será llevada a cabo en el apartado de Resultados. En esta exploración se realizará un análisis inicial de los datos con dim, head y summary para entender el tipo de datos que hay. Además, se realizarán análisis de componentes principales (PCA)

Creación del repositorio Github

Primeramente, se ha creado un nuevo repositorio online llamado “Canseco-i-Ribas-Joan-PEC1” desde el usuario “CansecoRibas”. Dentro del directorio local del proyecto (donde hay todos los documentos que se han ido mencionando en la parte de “Materiales y Métodos”), se ha creado primeramente un proyecto de R, para facilitar el trabajo con estos documentos. Posteriormente, y desde la Terminal de Rstudio, se ha creado un repositorio local Git usando el comando “git init”. Después se ha establecido conexión con el repositorio online con “git remote add origin <https://github.com/CansecoRibas/Canseco-i-Ribas-Joan-PEC1.git>”. Se han añadido todos los documentos al área de ensayo con “git add .” y se ha hecho el primer commit con “git commit -m ‘Primer commit del proyecto’”. Seguidamente, se han enviado todos los documentos del repositorio local al repositorio online con “git push origin master”. Una vez se han vinculados los repositorios, se han ido haciendo commits a medida que se ha ido avanzando en el proyecto. Se puede acceder al repositorio online desde: <https://github.com/CansecoRibas/Canseco-i-Ribas-Joan-PEC1>

Resultados

En este apartado se volcarán los resultados del proyecto.

Creación del contenedor del tipo SummarizedExperiment

El primer paso será asegurarnos de que “BiocManager” está instalado:

```
if (!requireNamespace("BiocManager", quietly = TRUE)) install.packages("BiocManager")
```

Posteriormente instalamos desde Bioconductor el paquete “SummarizedExperiment” que es el paquete específico para crear el contenedor.

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

Y lo llamamos con library.

```
library("SummarizedExperiment")
```

Ahora cargamos los datos (“DataValues_S013.csv”) y los metadatos (“DataInfo_S013.csv”) del dataset que están dentro de la carpeta “2018-MetabotypingPaper”.

```
data_values <- read.csv("2018-MetabotypingPaper/DataValues_S013.csv")
metadata <- read.csv("2018-MetabotypingPaper/DataInfo_S013.csv")
```

Observamos estos dos archivos. Posteriormente ya se hará una mejor exploración de los datos en próximos apartados. Dejamos este código en forma de comentarios porque sinó el informe se llena de datos densos.

```
#head(data_values)
#head(metadata)
```

En estos dos documentos hay dos variables iniciales que sobran, ya que son duplicaciones de las columnas adyacentes. Estas columnas son X.1 en “DataValues_S013.csv” y X en “DataInfo_S013.csv”, por lo que las vamos a eliminar para que no nos molesten.

```
data_values_mod <- data_values[, 2:696]
metadata_mod <- metadata[, 2:4]
```

Vemos que hay algunas variables identificativas (las 5 primeras) en el documento “DataValues_S013.csv” y el resto son variables experimentales. En el documento “DataInfo_S013.csv” se recoge información de las variables del documento “DataValues_S013.csv”. Algunas de las variables son las mismas pero en distintos timepoints por lo que vamos a crear una nueva variable que especifique el timepoint. De esta manera, conseguiremos obtener muestras de manera individual para las variables. Primero nos aseguramos que los paquetes dplyr y tidyr estén instalados.

```
if (!requireNamespace("tidyr", quietly = TRUE)) install.packages("tidyr")
if (!requireNamespace("dplyr", quietly = TRUE)) install.packages("dplyr")
```

Los llamamos con library.

```
library("tidyr")
library("dplyr")
```

Creamos una nueva variable que se llame Timepoint haciendo pivot_longer para así obtener las muestras de manera individual que es lo que nos interesa.

```
data_values_samples <- data_values_mod %>%
  # Hacemos pivot_longer para transformar el df.
  pivot_longer(
    cols = 6:695, # Seleccionamos las columnas de metabolitos.
    names_to = c(".value", "Timepoint"), # Extraemos los metabolitos como
    # columnas y establecemos timepoints como nueva columna.
    names_sep = "_T" # Establecemos el separador entre el nombre del metabolito
    # y el timepoint.
  ) %>%
  # Especificamos como queremos que se muestren los timepoints.
  mutate(Timepoint = paste0("T", Timepoint)) %>%
  filter(!Timepoint == "TNA") # Eliminamos las filas que se nos han creado por
  # error que están completamente vacías.
```

Ahora modificamos el metadata para adecuarlo a las nuevas variables. Primero añadimos la nueva fila de timepoints creando un nuevo df.

```
metadata_timepoints <- data.frame(
  VarName = "Timepoint",
  varTpe = "character",
  Description= "Punto temporal de recogida de muestra del paciente (0, 2, 4 o 5)"
)
```

Modificamos las variables del metadata eliminando “_TX” del final y eliminamos variables repetidas.

```
metadata_metabolites <- metadata_mod %>%
  mutate(VarName = gsub("_T[0-5]$", "", VarName)) %>% # Buscamos y eliminamos el "_TX" del final.
  distinct(VarName, .keep_all = TRUE) # Quitamos las variables repetidas en Varname.
```

Finalmente unimos los dos dfs con rbind.

```
# Finalmente unimos los dos dfs con rbind.
metadata_samples <- rbind(metadata_timepoints, metadata_metabolites)
```

Vemos que metadata_samples tiene las mismas filas que columnas tiene data_values_samples que es lo que estábamos buscando.

```
dim(metadata_samples)
```

```
## [1] 180 3
```

```
dim(data_values_samples)
```

```
## [1] 156 180
```

Una vez modificados ambos documentos, nos quedan los datos separados por muestras individuales cogidas en diferentes timepoints. Ahora ya podemos crear el objeto “SummarizedExperiment”. Para ello, seleccionamos las variables experimentales (a partir de la variable 7 hasta el final) y creamos una matriz transpuesta llamada “metabolite_data”, ya que las filas deben ser las variables y las columnas deben ser las distintas muestras.

```
metabolite_data <- t(as.matrix(data_values_samples[, 7:ncol(data_values_samples)]))
```

Creamos rowData a partir de la información de la metadata para especificar la información de las variables experimentales.

```
rowData <- metadata_samples[7:nrow(metadata_samples), ]
```

Creamos ColData seleccionando las variables de identificación de los pacientes.

```
colData <- data_values_samples[, 1:6]
```

Creamos el objeto *SummarizedExperiment* con los 3 componentes que hemos mencionado.

```
sum_exp_object <- SummarizedExperiment(  
  assays = list(counts = metabolite_data),  
  rowData = rowData,  
  colData = colData  
)
```

Vemos qué pinta tiene.

```
sum_exp_object
```

```
## class: SummarizedExperiment  
## dim: 174 156  
## metadata(0):  
## assays(1): counts  
## rownames(174): MEDDM MEDCOL ... X lysoPC.a.C14.0  
## rowData names(3): VarName varTpe Description  
## colnames: NULL  
## colData names(6): SUBJECTS SURGERY ... Group Timepoint
```

Vemos que es un objeto con 174 filas (variables experimentales de metabolitos) y 156 columnas, que son las muestras (4 para cada paciente, ya que son 4 timepoints). Se puede estudiar un poco más el objeto con las siguientes funciones, pero las dejaremos como comentarios para que no se muestren unos tochos terribles en el informe.

```
#assay(sum_exp_object) # Devuelve la matriz principal de datos.  
#rowData(sum_exp_object) # Devuelve un df con información sobre los metabolitos  
#(filas de la matriz).  
#colData(sum_exp_object) # Devuelve un df con información sobre las muestras  
#(columnas de la matriz).
```

Solo queda guardar el objeto tal y como nos pide el enunciado.

```
save(sum_exp_object, file = "sum_exp_object.Rda")
```

Exploración del dataset PENDIENTE

Empezamos con el análisis de ambos documentos, “DataValues_S013.csv” y “DataInfo_S013.csv”. Como se ha indicado anteriormente, el primero contiene los datos en relación a los metabolitos recogidos en muestras de 4 distintos timepoints de 39 pacientes, y el segundo contiene la explicación de las características de las variables y su información. Por lo que nos fijaremos para hacer los análisis en el primer documento (“DataValues_S013.csv”).

Primero haremos una exploración inicial. Cogemos el documento data_values_samples, en el que hemos transformado los datos para que las filas sean las muestras recogidas (y no los pacientes como antes). Verificaremos las dimensiones del documento para ver si lo que se cumple es cierto.

```
dim(data_values_samples)
```

```
## [1] 156 180
```

Vemos que es un documento con 156 filas (correspondiente a los 4 timepoints de los 39 pacientes) y 180 columnas (correspondientes a las variables de estudio). Hacemos ahora un resumen estadístico muy básico de las 180 variables para ver qué valores toman y qué tipos de variables tenemos. Es un poco complicado trabajar con este número tan elevado de variables ya que los outputs son enormes, por lo que en este caso tampoco vamos a mostrar el resultado.

```
#summary(data_values_samples)
```

También es importante verificar los missings de los datos. Vamos a mirar cuantos missings hay por cada variable.

```
colSums(is.na(data_values_samples))
```

##	SUBJECTS	SURGERY	AGE	GENDER	Group
##	0	0	0	0	0
##	Timepoint	MEDDM	MEDCOL	MEDINF	MEDHTA
##	0	13	13	15	13
##	GLU	INS	HOMA	HBA1C	HBA1C.mmol.mol
##	11	14	14	91	91
##	PESO	bmi	CC	CINT	CAD
##	10	10	24	14	14
##	TAD	TAS	TG	COL	LDL
##	25	25	11	11	18
##	HDL	VLDL	PCR	LEP	ADIPO
##	13	49	68	68	71
##	GOT	GPT	GGT	URICO	CREAT
##	12	12	15	16	15
##	UREA	HIERRO	TRANSF	FERR	Ile
##	15	13	25	16	18
##	Leu	Val	Ala	Pro	Gly
##	18	18	18	18	18
##	Ser	Trp	Phe	Met	Orn
##	18	18	18	18	18
##	Arg	His	Asn	Asp	Glu
##	18	18	18	18	18
##	Gln	Cit	Tyr	Thr	Lys

##	18	18	18	18	18
##	Creatinine	Kynurenine	Putrescine	Sarcosine	Serotonin
##	18	18	21	18	18
##	Taurine	SDMA	C0	C2	C3.OH
##	18	18	18	18	18
##	C6..C4.1.DC.	C5.DC..C6.OH.	C7.DC	C8	C10
##	18	18	18	18	18
##	C10.1	C10.2	C14.1	C14.2	C16.1
##	18	18	18	18	18
##	C16.2	C16.2.OH	C18.1	C18.1.OH	C18.2
##	18	18	18	18	18
##	lysoPC.a.C16.0	lysoPC.a.C16.1	lysoPC.a.C17.0	lysoPC.a.C18.0	lysoPC.a.C18.1
##	18	18	18	18	18
##	lysoPC.a.C18.2	lysoPC.a.C20.3	lysoPC.a.C20.4	lysoPC.a.C24.0	lysoPC.a.C26.0
##	18	18	18	18	18
##	lysoPC.a.C26.1	lysoPC.a.C28.0	lysoPC.a.C28.1	PC.aa.C24.0	PC.aa.C28.1
##	18	18	18	18	18
##	PC.aa.C30.0	PC.aa.C32.0	PC.aa.C32.1	PC.aa.C32.3	PC.aa.C34.1
##	18	18	18	18	18
##	PC.aa.C34.2	PC.aa.C34.3	PC.aa.C34.4	PC.aa.C36.0	PC.aa.C36.1
##	18	18	18	18	18
##	PC.aa.C36.2	PC.aa.C36.3	PC.aa.C36.4	PC.aa.C36.5	PC.aa.C38.0
##	18	18	18	18	18
##	PC.aa.C38.1	PC.aa.C38.3	PC.aa.C38.4	PC.aa.C38.5	PC.aa.C38.6
##	18	18	18	18	18
##	PC.aa.C40.1	PC.aa.C40.2	PC.aa.C40.3	PC.aa.C40.4	PC.aa.C40.5
##	18	18	18	18	18
##	PC.aa.C40.6	PC.aa.C42.0	PC.aa.C42.1	PC.aa.C42.2	PC.aa.C42.4
##	18	18	18	18	18
##	PC.aa.C42.5	PC.aa.C42.6	PC.aa.C30.0	PC.aa.C32.1	PC.aa.C32.2
##	18	18	18	18	18
##	PC.aa.C34.0	PC.aa.C34.1	PC.aa.C34.2	PC.aa.C34.3	PC.aa.C36.0
##	18	18	18	18	18
##	PC.aa.C36.1	PC.aa.C36.2	PC.aa.C36.3	PC.aa.C36.4	PC.aa.C36.5
##	18	18	18	18	18
##	PC.aa.C38.0	PC.aa.C38.2	PC.aa.C38.3	PC.aa.C38.4	PC.aa.C38.5
##	18	18	18	18	18
##	PC.aa.C38.6	PC.aa.C40.1	PC.aa.C40.2	PC.aa.C40.3	PC.aa.C40.4
##	18	18	18	18	18
##	PC.aa.C40.5	PC.aa.C40.6	PC.aa.C42.1	PC.aa.C42.2	PC.aa.C42.3
##	18	18	18	18	18
##	PC.aa.C42.4	PC.aa.C42.5	PC.aa.C44.3	PC.aa.C44.4	PC.aa.C44.5
##	18	18	18	18	18
##	PC.aa.C44.6	SM..OH..C14.1	SM..OH..C16.1	SM..OH..C22.1	SM..OH..C22.2
##	18	18	18	18	18
##	SM..OH..C24.1	SM.C16.0	SM.C16.1	SM.C18.0	SM.C18.1
##	18	18	18	18	18
##	SM.C20.2	SM.C24.0	SM.C24.1	X lysoPC.a.C14.0	
##	18	18	18	156	118

Vemos que las variables identificativas están completas, mientras que todas las variables de metabolitos tienen algún missing, aunque en general son pocos (la mayoría tienen menos de 19 missings de 156 valores totales). Además, hay una variable “X” que son todo missings.

Como hay muchas variables de metabolitos, haremos una exploración inicial de las variables identificativas

de los pacientes, que, como hemos dicho son las 5 primeras variables (excluyendo la variable Timepoint que en este caso no tiene mucho sentido mantenerla porque hace referencia a las muestras).

```
data_values_patients <- data_values_samples %>%  
  select("SUBJECTS", "SURGERY", "AGE", "GENDER", "Group") %>% # Seleccionamos columnas  
  distinct() # Eliminamos duplicados de pacientes
```

Verificamos que tenemos 39 pacientes en este df.

```
nrow(data_values_patients)
```

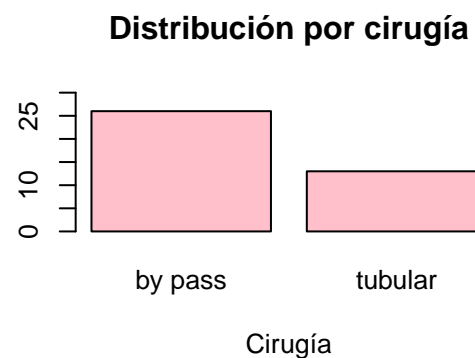
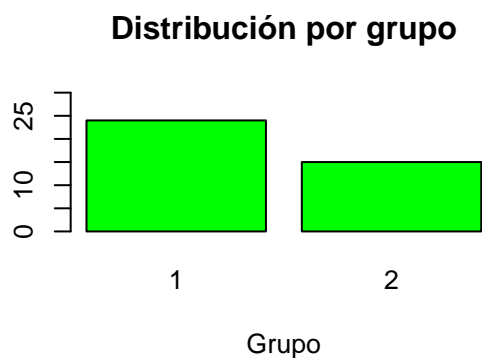
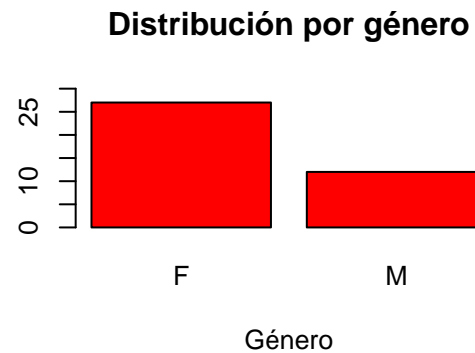
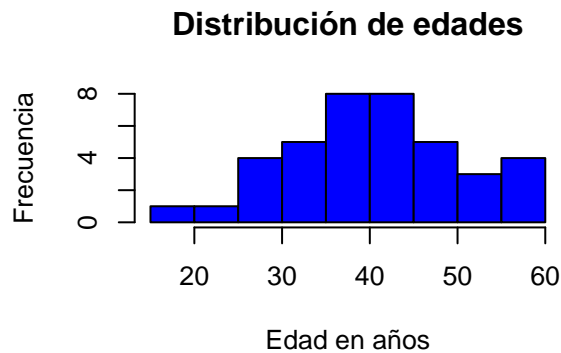
```
## [1] 39
```

Verificamos las variables que deban ser factors realmente sean factors.

```
data_values_patients$SURGERY <- factor(data_values_patients$SURGERY)  
data_values_patients$GENDER <- factor(data_values_patients$GENDER)  
data_values_patients$Group <- factor(data_values_patients$Group)
```

Ahora podemos hacer estadística básica con estos datos de los pacientes. Los datos realmente importantes de este df se encuentran en las variables “SURGERY”, “AGE”, “GENDER” y “Group”, ya que “SUBJECTS” solo hace referencia al número de paciente. La variable “AGE” la veremos con un histograma mientras que los factores los veremos con diagramas de barras.

```
par(mfrow = c(2,2)) # Así podemos ver los 4 gráficos a la vez.  
# Histograma de la edad.  
hist(data_values_patients$AGE,  
      main="Distribución de edades",  
      col="blue",  
      xlab="Edad en años",  
      ylab="Frecuencia")  
  
# Diagrama de barras para el género.  
barplot(table(data_values_patients$GENDER),  
        main="Distribución por género",  
        col="red",  
        xlab="Género",  
        ylim=c(0,30))  
  
# Diagrama de barras para el grupo.  
barplot(table(data_values_patients$Group),  
        main="Distribución por grupo",  
        col="green",  
        xlab="Grupo",  
        ylim=c(0,30))  
  
# Diagrama de barras para la cirugía.  
barplot(table(data_values_patients$SURGERY),  
        main="Distribución por cirugía",  
        col="pink",  
        xlab="Cirugía",  
        ylim=c(0,30))
```



Estos gráficos muestran que la mayoría de edades de los participantes están entre los 30 y los 50 años, que hay el doble de mujeres que de hombres, que hay más componentes del grupo 1 que del grupo 2 y que hay más pacientes con cirugía by pass que tubular. Ahora miraremos la relación entre la variable SURGERY (la variable que indica el tipo de cirugía que ha tomado cada paciente) y las demás, ya que la variable SURGERY es la variable clave del estudio.

```
# Separamos los datos según la variable SURGERY.
split_data <- split(data_values_patients, data_values_patients$SURGERY)

# Redistribuimos el espacio gráfico.
par(mfrow = c(2, 3))

# Hacemos loop.
for (surgery_status in names(split_data)) {
  data_subset <- split_data[[surgery_status]]

  # Histograma de la edad.
  hist(data_subset$AGE,
        main=paste("Edades (", surgery_status, ")"),
        col="blue",
        xlab="Edad en años",
        ylab="Frecuencia")

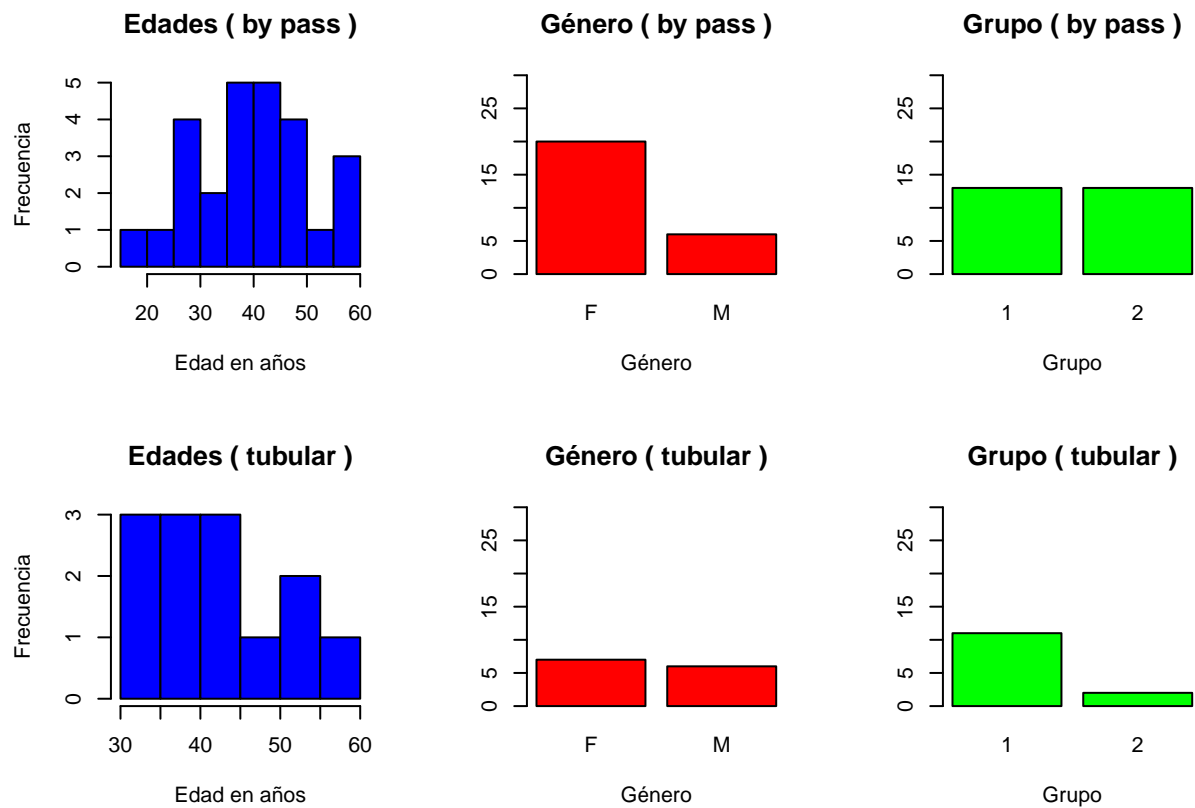
  # Diagrama de barras para el género.
  barplot(table(data_subset$GENDER),
           main=paste("Género (", surgery_status, ")"),
```

```

col="red",
xlab="Género",
ylim=c(0,30))

# Diagrama de barras para el grupo.
barplot(table(data_subset$Group),
main=paste("Grupo (", surgery_status, ")"),
col="green",
xlab="Grupo",
ylim=c(0,30))
}

```



Discusión y limitaciones y conclusiones del estudio PENDIENTE