

# Genómica funcional y análisis de microarrays

## Primera PEC : Selección de genes expresados diferencialmente

Alex Sánchez

30 de mayo de 2012

### Índice

<b>1. Introducción y objetivos</b>	<b>2</b>
1.1. Estructura del documento . . . . .	2
1.2. Objetivos . . . . .	2
<b>2. Materiales y métodos</b>	<b>3</b>
2.1. Etapas del análisis . . . . .	3
2.2. Métodos de análisis . . . . .	3
2.2.1. Herramientas y procedimientos bioinformáticos de análisis	4
<b>3. Obtención y lectura de los datos</b>	<b>5</b>
3.1. Los datos para el análisis . . . . .	5
3.1.1. Localización de los datos . . . . .	5
3.1.2. Selección de muestras para el análisis . . . . .	6
3.1.3. Lectura de los datos . . . . .	6
<b>4. Preprocesado: Exploración, Control de Calidad y Normalización</b>	<b>7</b>
4.1. Exploración y visualización . . . . .	8
4.1.1. Exploración estadística de los datos . . . . .	8
4.1.2. Control de calidad (2): métodos específicos para microarrays	11
4.1.3. Paquetes “all-in-one” para el control de calidad” . . . . .	13
4.1.4. Archivos con los resultados del control de calidad . . . . .	13
4.2. Normalización y Filtraje . . . . .	14
4.2.1. Filtraje . . . . .	16
4.2.2. Archivos de resultados normalizados . . . . .	16
<b>5. Selección de genes diferencialmente expresados</b>	<b>17</b>
5.1. Análisis basado en modelos lineales . . . . .	17
5.1.1. Matrices de diseño y de contrastes . . . . .	17
5.1.2. Estimación del modelo y selección de genes . . . . .	18
5.1.3. Principales genes expresados diferencialmente . . . . .	20

<b>6. Post-procesado de las listas de genes obtenidas</b>	<b>22</b>
6.1. Comparaciones múltiples . . . . .	22
6.2. Anotación de resultados . . . . .	23
6.2.1. Tablas de anotaciones sencillas . . . . .	25
6.2.2. Archivos de resultados con tablas hiperenlazables . . . . .	25
6.3. Visualización de los perfiles de expresión . . . . .	26
6.4. Análisis de la significación biológica . . . . .	28
<b>7. Resumen de resultados y discusión</b>	<b>30</b>
7.1. Discusión . . . . .	30
<b>A. Código R del análisis</b>	<b>32</b>
<b>B. Lista de archivos generados en el estudio</b>	<b>32</b>

## 1. Introducción y objetivos

Este documento presenta una solución al ejercicio de análisis de microarrays propuesto en la asignatura de *Genómica Funcional y Análisis de Microarrays*. Su objetivo no es mostrar como realizar el mejor análisis de microarrays posible sino como hacerlo utilizando los conceptos y ejemplos proporcionados en las unidades dedicadas a este tema.

### 1.1. Estructura del documento

Típicamente un trabajo científico o técnico se organiza en varias partes: *Introducción*, *Material y Métodos*, *Resultados* y *Discusión*. Dado que el objetivo de este estudio es demostrar la aplicación de los métodos de análisis explicados en la asignatura esta estructura se relajará. Después de la introducción se realizará una breve descripción del proceso general y los métodos utilizados en cada paso. A continuación se presentarán los resultados obtenidos en cada fase del análisis. Finalmente se presentará una breve discusión sobre las posibles limitaciones encontradas.

### 1.2. Objetivos

El objetivo de este trabajo es realizar un análisis de datos de microarrays para encontrar genes diferencialmente expresados entre varios tipos de tumores de cancer de mama clasificados en tres grupos: apocrinos (APO), basales (BAS) y luminales (LUMI).

Esta clasificación se basa en la resistencia de los tumores a los receptores de estrógenos y de andrógenos.

- Los tumores clasificados como “APO” son negativos para los receptores de estrógenos (ER-) y positivos para los receptores de andrógenos (AR+).
- los clasificados como “LUMI” son (ER+) y (AR+) y
- los clasificados como “BAS” son (ER-) y (AR-).

Es conveniente destacar que, de hecho, el objetivo de la práctica difiere del objetivo del artículo en que se basa. Mientras que el artículo se busca *descubrir* los tres grupos con los que trabajamos en este ejercicio se pretende caracterizarlos a través de los genes que se expresan de forma distinta entre ellos. No deja de ser un problema del tipo “el huevo y la gallina”, es decir: *¿Puesto que se trata de grupos distintos encontramos genes que se expresan diferencialmente entre ellos o dado que hay genes que se expresan diferencialmente podemos deducir que son grupos distintos?*

## 2. Materiales y métodos

### 2.1. Etapas del análisis

El análisis se ha realizado siguiendo las pautas descritas en los capítulos 6 a 11 de los apuntes de la asignatura <http://ueb.vhir.org/cursos/Microarrays> y resumidos en la figura 1.

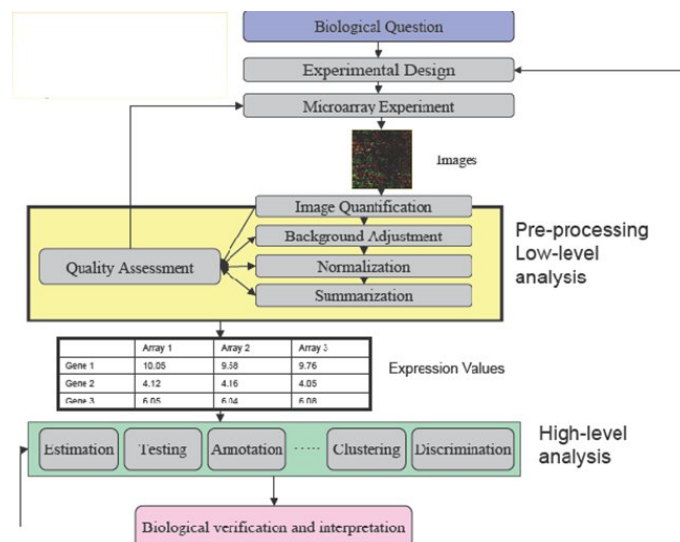


Figura 1: The Microarray Analysis Process

### 2.2. Métodos de análisis

Antes de hibridar los microarrays<sup>1</sup> se comprobó la calidad del cRNA de cada muestra. Sólo las que demostraron una calidad suficiente –presumiblemente un valor de “Bio-analyser”<sup>2</sup> igual o superior a 7– se sometieron a análisis posterior. Las imágenes resultantes del escaneado de los microarrays se procesaron con Microarray Analysis Suite 5.0 (Affymetrix).

<sup>1</sup>Esta primera parte es ficticia en tanto que no hemos tenido acceso al proceso de laboración de los microarrays, por lo que he escrito lo que suele ponerse en esta parte

<sup>2</sup>Sistema para el análisis de la calidad del RNA

Los valores “crudos” de expresión obtenidos directamente de los archivos CEL se preprocesa utilizando el método de RMA ([7]), un proceso de tres pasos que integra la corrección de fondo, la normalización y el resumen de los valores del grupo de sondas en un único valor de expresión “absoluta”. Dichos valores normalizados fueron la base para todos los análisis.

Antes de la selección de genes los valores normalizados se sometieron a un filtraje no específico para eliminar los genes de baja variabilidad (los genes cuyo rango intercuartil entre todas las muestras no superó un umbral mínimo).

La selección de genes diferencialmente expresados entre condiciones experimentales se basó en el método desarrollado por Smyth y otros [11].

Este método extiende el análisis de la varianza clásico utilizando métodos Bayesianos empíricos para combinar la información de cada gen individual con la de todos los genes restantes para obtener mejores estimaciones de error. Esto es de gran utilidad en análisis de microarrays, un contexto en el que el tamaño de las muestras es a menudo pequeño lo que puede dar lugar a estimaciones de los errores erráticos y, en consecuencia, p-valores que no son de fiar.

Los genes más relevantes de cada comparación se resaltaron utilizando “volcano-plots”, que organizan los genes a lo largo de dos dimensiones que podemos considerar de importancia biológica y estadística. El eje horizontal representa el cambio medio de expresión entre los dos grupos (en una escala logarítmica, por lo que la regulación hacia arriba y abajo aparecen simétrica), y el segundo (vertical) representa el “menos logaritmo del p-valor” por lo que los genes cuyo p-valor asociado sea inferior aparecen más arriba. El primer eje indica el impacto biológico del cambio, y el segundo indica la evidencia estadística, o la fiabilidad de dicho cambio.

Con el fin de hacer frente a la problemática derivada del hecho de que muchas pruebas (una por cada gen) se realizan simultáneamente, se realizó un ajuste de p-valores para obtener control sobre la tasa de falsos positivos usando el método de Benjamini y Hochberg ([1]).

Los genes seleccionados como diferencialmente expresados se agruparon para buscar patrones comunes de expresión entre condiciones experimentales. Para ello se utilizaron mapas de colores o “HEatmaps” que realizan una agrupación jerárquica de los genes y/o las muestras y la representan mediante una gama de colores apropiada, de forma que valores altos o bajos se corresponden a colores distintos de la gama escogida.

Las listas de genes diferencialmente expresados se anotaron en diversas bases de datos (Entrez, Unigene, Gene Ontology, KEGG, ...) utilizando los paquetes de anotación para microarrays de affymetrix disponibles en el proyecto Bioconductor (<http://bioconductor.org>).

Para contribuir a la interpretación biológica de los resultados se realizó un análisis de enriquecimiento [5, 3] que busca establecer si las categorías funcionales de los genes seleccionados aparecen entre estos genes con mayor (o menor) frecuencia que entre todos los del genoma. De ser así se indica que la lista de genes se encuentra “enriquecida” en estas funcionalidades, o lo que es lo mismo que los procesos afectados por las diferencias son éstos.

### 2.2.1. Herramientas y procedimientos bioinformáticos de análisis

Los análisis estadísticos se realizaron utilizando el lenguaje estadístico *Ry* las librerías desarrolladas para el análisis de microarray en el proyecto de Biocon-

ductor ([www.bioconductor.org](http://www.bioconductor.org)). Para más detalle sobre los métodos descritos en esta sección puede consultarse [6].

El código siguiente se utilizó para instalar los paquetes de Bioconductor necesarios para el análisis.

```
> installifnot <- function (pkgName){
+   if(!(require(pkgName, character.only=TRUE))){
+     source("http://Bioconductor.org/biocLite.R")
+     biocLite(pkgName)
+   }
+ }
> installifnot("Biobase")
> installifnot("hgu133a.db")
> installifnot("affy")
> installifnot("affyPLM")
> installifnot("arrayQualityMetrics")
> installifnot("genefilter")
> installifnot("limma")
> installifnot("hgu133a.db")
> installifnot("annotate")
> installifnot("annaffy")
> installifnot("hwriter")
> installifnot("gplots")
> installifnot("GOstats")
```

### 3. Obtención y lectura de los datos

#### 3.1. Los datos para el análisis

Los datos en que se basa el estudio se han obtenido a partir de tumores de mama (“advanced/inflammatory breast tumours”) y fueron tomados antes del tratamiento de pacientes enrolados en un estudio clínico (EORTC 10994).

Los microarrays se prepararon a partir de RNA total extraído de secciones de 25 mm de biopsias y amplificados mediante el procedimiento “Eberwine T7 procedure” siguiendo el protocolo indicado por Affymetrix para pequeñas muestras.

Los datos de los microarrays se encuentran en Gene Expression Omnibus (GEO) con el número de serie GSE1561. Puede accederse a ellos en el siguiente enlace <http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE1561>. El artículo de Farmer ([4]) *et al.* (<http://www.ncbi.nlm.nih.gov/pubmed/15897907>) describe el estudio.

##### 3.1.1. Localización de los datos

Como es habitual en este curso supondremos que trabajamos en un directorio escogido por nosotros y cuya localización se asigna a la variable `workingDir`.

Asumiremos que los datos se encuentran en un subdirectorio del anterior, denominado “data” que se almacenará en la variable `dataDir` y que los resultados se almacenarán en un directorio “results” cuyo nombre completo se almacenará en la variable `resultsDir`.

```

> workingDir <- getwd()
> dataDir <- file.path(workingDir, "data")
> resultsDir <- file.path(workingDir, "results")
> celfilesDir <- file.path(workingDir, "celfiles")
> setwd(workingDir)

```

### 3.1.2. Selección de muestras para el análisis

Este análisis se ha basado en un subconjunto de muestras del estudio original. Las muestras se han obtenido de la lista original mediante un pequeño script de *R*.

De la dirección <http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE1561> es posible descargar un archivo comprimido con los 49 archivos “.cel”. Con el fin de simplificar el análisis se decidió seleccionar aleatoriamente 5 arrays de cada grupo utilizando a partir de la información contenida en el archivo: “Asignacion\_de\_Muestras\_A\_Grupos.csv” elaborado a partir de la información sobre el diseño experimental contenida en: <http://www.ncbi.nlm.nih.gov/geo/gds/profileGraph.cgi?gds=1329>

```

> muestras <- read.csv2(file.path(dataDir, "Asignacion_de_muestras_a_grupos.csv"),
+                       head=T)
> misMuestras <- as.character (muestras$Sample)
> paraAnalisis <- c(sample(misMuestras[1:6], 5),
+                   sample(misMuestras[7:20], 5),
+                   sample(misMuestras[23:48], 5))
> alAnalisis <- muestras[muestras$Sample %in% paraAnalisis,]
> write.table(alAnalisis, file=file.path(dataDir, "targets.txt"),
+            sep="\t", row.names=FALSE, quote=FALSE)

```

### 3.1.3. Lectura de los datos

La lectura de datos se lleva a cabo utilizando las clases y métodos definidas en los paquetes *Biobase* y *affy* de Bioconductor.

Una forma cómoda de leer los datos y, al mismo tiempo, asignar a cada muestra los valores de las covariables (por ejemplo el grupo para el analisis) consiste en crear un pequeño archivo de texto, que suele denominarse **targets.txt** y que contiene la identificación de cada archivo con la asignación de cada muestra a cada condición experimental. El archivo targets para este caso tiene el aspecto que se muestra en la tabla 1.

	Sample	Ids	SampleIDs	Group
GSM26878.CEL	GSM26878	PF14 EnPnT2N1G2	PF14	A
GSM26883.CEL	GSM26883	PF19 EnPuT4N0Gu	PF19	A
GSM26886.CEL	GSM26886	PF22 EnPnT2N1G2	PF22	A
GSM26887.CEL	GSM26887	PF23 EnPnT2N0G2	PF23	A
GSM26903.CEL	GSM26903	PF39 EnPuT4N0Gu	PF39	A
GSM26882.CEL	GSM26882	PF18 EnPnT2N1G3	PF18	B
GSM26888.CEL	GSM26888	PF24 EnPnTiN0G3	PF24	B
GSM26893.CEL	GSM26893	PF29 EnPnT3N1G3	PF29	B
GSM26905.CEL	GSM26905	PF41 EnPnT3N1G2	PF41	B
GSM26906.CEL	GSM26906	PF42 EnPnT2N2G3	PF42	B

GSM26804.CEL	GSM26804	PF01	EpPpT2N0G2	PF01	L
GSM26868.CEL	GSM26868	PF03	EpPpT3N0Gu	PF03	L
GSM26877.CEL	GSM26877	PF13	EpPpT2N1G1	PF13	L
GSM26901.CEL	GSM26901	PF37	EpPpT3N1G2	PF37	L
GSM26911.CEL	GSM26911	PF47	EpPpT3N1G3	PF47	L

Cuadro 1: Archivo targets.txt con la asignación a cada muestra de su condición experimental

El contenido del archivo `targets` se utiliza en la lectura de los datos y la creación del objeto `rawData` de la clase `affybatch` que contendrá las intensidades “crudas” de cada archivo .CEL

```
> require(affy)
> sampleInfo <- read.AnnotatedDataFrame(file.path(dataDir,"targets.txt"),
+   header = TRUE, row.names = 1, sep="\t")
> fileNames <- rownames(pData(sampleInfo))
> rawData <- read.affybatch(filename=file.path(celfilesDir,fileNames),
+   phenoData=sampleInfo)
> show(rawData)

AffyBatch object
size of arrays=712x712 features (27 kb)
cdf=HG-U133A (22283 affyids)
number of samples=15
number of genes=22283
annotation=hgu133a
notes=
```

Este objeto es la base para todos los análisis que se realizarán.

## 4. Preprocesado: Exploración, Control de Calidad y Normalización

Los datos procedentes de la lectura de los microarrays se denominan datos “crudos” y deben ser pre-procesados de diversas formas antes de analizarlos.

“Preprocesado” es un término genérico que engloba varios procesos

- Exploración y control de calidad de los datos.
- Normalización y resumen (llamdo “sumarización”) de los valores de las sondas de cada grupo de sondas.
- Filtrado no específico para eliminar el efecto de genes que no se expresan o bien no se expresan de forma distinta entre los grupos.

A su vez la exploración y el control de calidad contempla:

1. Exploraciones estadísticas estándar.
2. Técnicas de control de calidad desarrolladas específicamente para datos de microarrays.

## 4.1. Exploración y visualización

La exploración de los datos suele basarse en técnicas univariantes como los histogramas o los diagramas de caja o en técnicas multivariantes como los análisis de conglomerados (“clusters”), de distancias o de análisis de componentes principales.

### 4.1.1. Exploración estadística de los datos

Un gráfico de densidad –mal llamado en este contexto, histograma– permite hacerse una idea de si las distribuciones de los distintos arrays son similares en forma y posición.

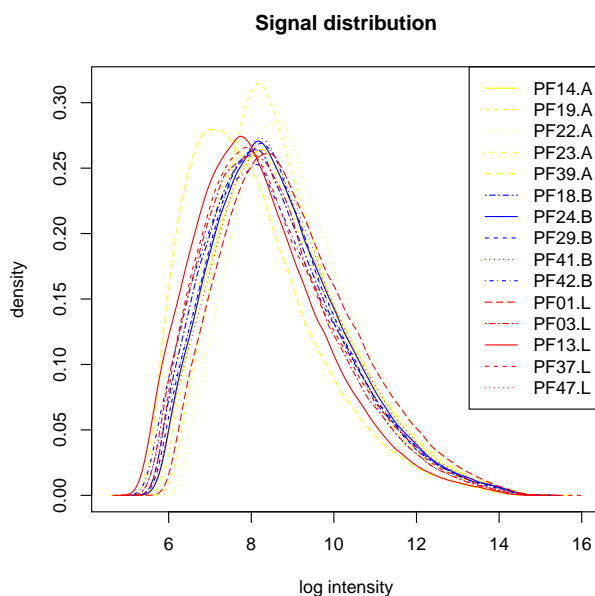


Figura 2: Distribución de las expresiones de cada array. Aunque los gráficos muestran una cierta heterogeneidad entre muestras el aspecto de todas ellas es compatible con la "normalidad".<sup>es</sup> decir no sugieren que pueda haber algún problema en los datos.

Un diagrama de cajas muestra la misma información facilitando más la comparación entre distribuciones.

Una representación multivariante como la ofrecida por los gráficos de las componentes principales suele resultar muy informativa de posibles problemas como efectos batch o presencia de valores raros". El gráfico de componentes principales muestra los datos en dimensión reducida de forma que cada componente explica una dimensión de mayor variabilidad e independiente de la siguiente. Representando el porcentaje de variabilidad explicada para cada componente tenemos una medida de la importancia de los grupos que se puedan visualizar. Si la suma de los porcentajes es alta, por ejemplo superior al 50 % las conclusiones que de ellos se deriven serán más fiables que con valores bajos, por ejemplo inferiores al 30 % de varianza explicada.



```
> boxplot(rawData, cex.axis=0.6, col=colores, las=2, names=sampleNames,
+         main="Signal distribution for selected chips")
```

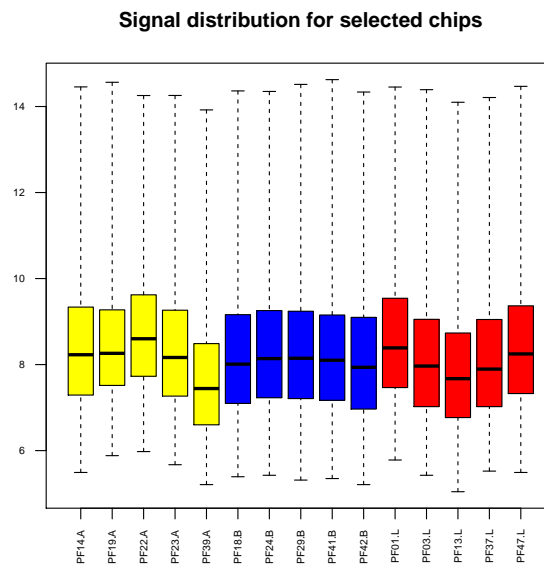


Figura 3: El diagrama de cajas muestra, de forma similar al histograma, cómo las distribuciones de los datos son relativamente similares. Hay heterogeneidad pero no se muestra ninguna diferencia sistemática, lo que sugiere que es conveniente normalizar pero no parece que hayan arrays problemáticos

```

> plotPCA <- function ( X, labels=NULL, colors=NULL, dataDesc="", scale=FALSE)
+ {
+   pcX<-prcomp(t(X), scale=scale) # o prcomp(t(X))
+   loads<- round(pcX$sdev^2/sum(pcX$sdev^2)*100,1)
+   xlab<-c(paste("PC1",loads[1],"%"))
+   ylab<-c(paste("PC2",loads[2],"%"))
+   if (is.null(colors)) colors=1
+   plot(pcX$x[,1:2],xlab=xlab,ylab=ylab, col=colors,
+        xlim=c(min(pcX$x[,1])-10, max(pcX$x[,1])+10))
+   text(pcX$x[,1],pcX$x[,2], labels, pos=3, cex=0.8)
+   title(paste("Plot of first 2 PCs for expressions in", dataDesc, sep=" "), cex=0.8)
+ }

> plotPCA(exprs(rawData), labels=sampleNames, dataDesc="selected samples")

```

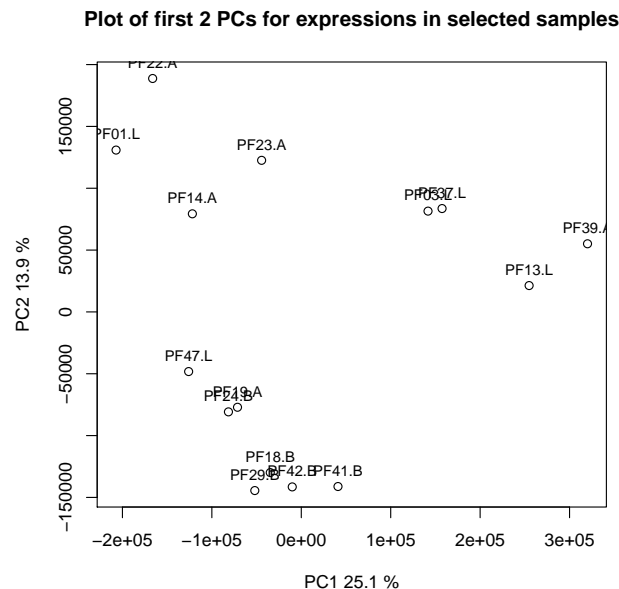


Figura 4: Representación de las dos primeras componentes de un análisis de componentes principales. Los grupos aparecen algo separados pero el porcentaje de variabilidad explicado por las dos primeras componentes sugiere que puede no ser muy explicativas

Un enfoque alternativo al del PCA aunque relacionado es realizar un análisis basado en distancias. Podemos hacerlo calculando y visualizando la matriz de distancias mediante un mapa de colores o mediante un agrupamiento jerárquico seguido de un dendrograma.

Un cluster jerárquico seguido de un dendrograma nos puede ayudar a hacernos una idea de si las muestras se agrupan por condiciones experimentales.

Si lo hacen es bueno, pero si no, no es necesariamente indicador de problemas, puesto que es un gráfico basado en todo los datos.

```
> manDist <- dist(t(exprs(rawData)))
> heatmap (as.matrix(manDist), col=heat.colors(16))
```

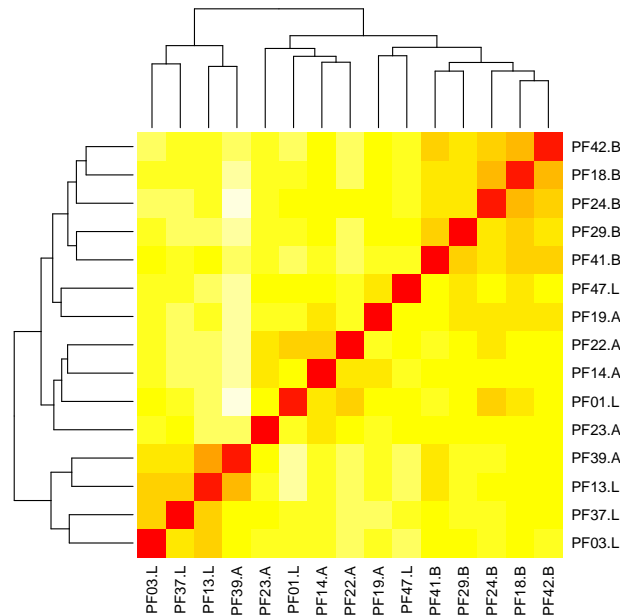


Figura 5: Mapa de colores de las distancias (euclídeas) entre los distintos arrays

#### 4.1.2. Control de calidad (2): métodos específicos para microarrays

Las exploraciones anteriores nos proporcionan una idea general acerca de la distribución de los datos y la posible presencia de grupos naturales (e.g. por tratamientos) o artificiales (ej. días de procesad).

Existen también métodos para el control de calidad específicos de microarrays, como por ejemplo:

- Los controles de calidad estándar de Affymetrix, descritos en el paquete `simpleaffy`.
- Controles basados en modelos a nivel de sondas, descritos en el paquete `affyPLM`.

De entre los controles de calidad “estándar” de affymetrix nos fijaremos en el gráfico de degradación que permite hacerse una idea de como ha sido el proceso de hibridación de las muestras. Si las líneas que forman el gráfico son más o menos paralelas sugiere una calidad similar en todos los arrays.

El “probe-level-model analysis” ajusta un model a la intensidad de las sondas y analiza los residuos de este ajuste a través de dos graficos: el RLE -de expresiones relativas- y el NUSE -de residuos centrados y normalizados- (vease [6]). Si los datos son de calidad ambos graficos deben ser centrados y relativamente simetricos. Cambios en esta situacion sugieren problemas en los arrays que no las verifiquen. Estos métodos se encuentran implementados en el paquete `affyPLM`.

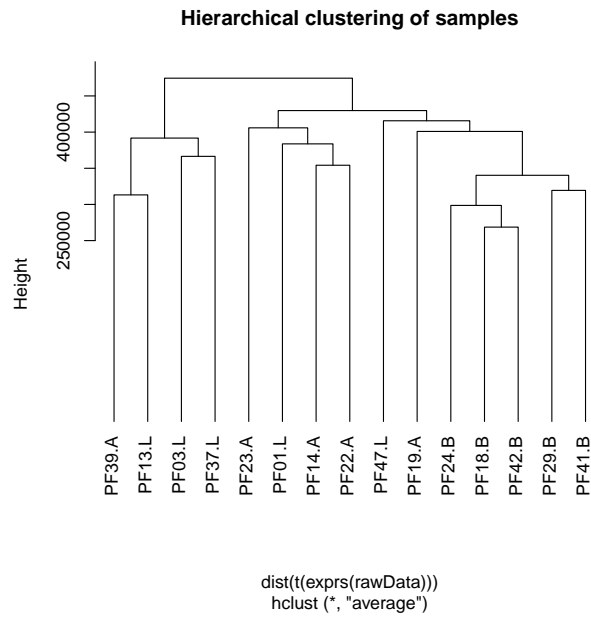


Figura 6: Dendrograma resultante de un agrupamiento jerárquico entre las muestras. Como en otras representaciones (Heatmaps, PCA) se observa una cierta agrupación de muestras de tipos similares aunque también una cierta heterogeneidad, atribuible probablemente a la diferencia entre pacientes

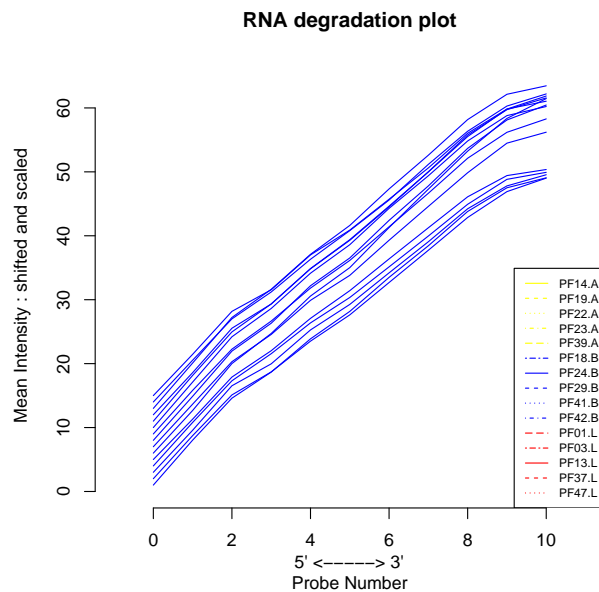


Figura 7: Gráfico de degradación. Las líneas paralelas muestran que el nivel de degradación del RNA es similar en todos los chips.

```

> stopifnot(require(affyPLM))
> if (!(exists("Pset")))
+   Pset<- fitPLM(rawData)

```

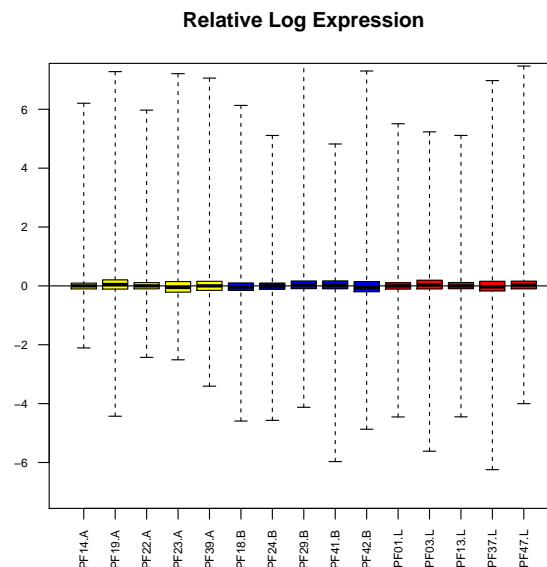


Figura 8: Gráfico de expresiones relativas (RLE) resultante del análisis basado en modelos (PLM). Las muestras son simétricas y similares lo que sugiera una calidad aceptable de los datos

#### 4.1.3. Paquetes “all-in-one” para el control de calidad”

El paquete `arrayQualityMetrics` encapsula todos los análisis anteriores, y alguno más, facilitando su ejecución e incluso su interpretación. La instrucción `arrayQualityMetrics` lleva a cabo todos los análisis de golpe y genera un informe de resultados con ayudas a la interpretación y a la detección de arrays problemáticos.

```

> stopifnot(require(arrayQualityMetrics))
> arrayQualityMetrics(rawData, outdir = file.path(resultsDir, "arrayQuality"),
+                     force=TRUE)

```

#### 4.1.4. Archivos con los resultados del control de calidad

Los resultados del control de calidad realizado con `arrayQualityMetrics` se encuentran accesibles a través del archivo `index.html` contenido en el subdirectorío creado al invocarlo -en este caso denominado `arrayQuality`.

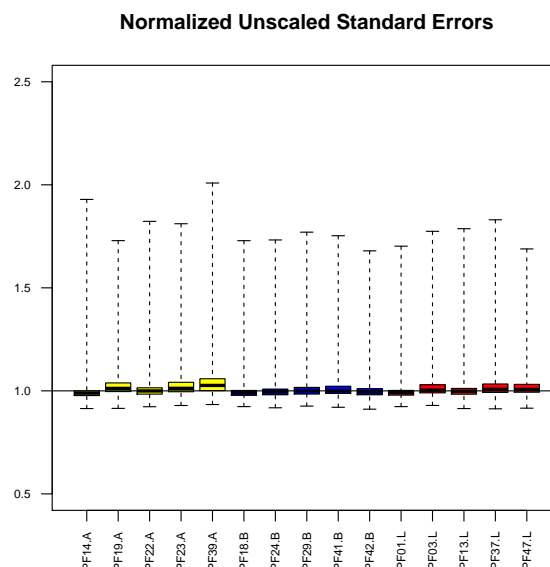


Figura 9: Gráfico de errores estandarizados normalizados provenientes de los residuos del análisis basado en modelos (PLM). No aparece ninguna caja *claramente* diferenciada de las otras lo que sugiera una calidad aceptable de los datos

## 4.2. Normalizacion y Filtraje

Una vez realizado el control de calidad se procede a normalizar los datos y sumarizarlos.

La normalización puede hacerse por distintos métodos (MAS5, VSN, RMA, GCRMA, ...) pero en este caso se utilizará el método RMA que es sin duda el más utilizado entre arrays de affymetrix.

El procesado mediante RMA implica un proceso en tres etapas:

- Corrección de fondo (el RMA hace precisamente esto).
- Normalización para hacer los valores de los arrays comparables.
- Summarización de las diversas sondas asociadas a cada grupo de sondas para dar un único valor.

```
> stopifnot(require(affy))
> eset_rma <- rma(rawData)
```

```
Background correcting
Normalizing
Calculating Expression
```

Un boxplot de los valores normalizados muestra que los valores han quedado claramente en una escala común en donde se pueden comparar.

```
> boxplot(eset_rma,main="RMA", names=sampleNames, cex.axis=0.7, col=colores,las=2)
```

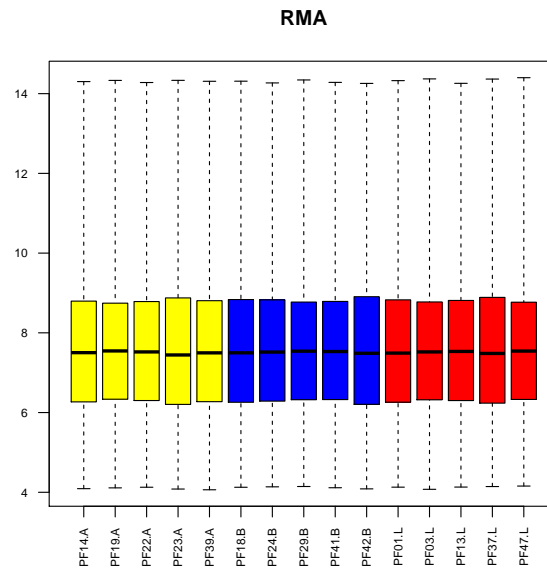


Figura 10: El diagrama de cajas después de normalizar muestra como dicho proceso ha puesto los valores en una escala completamente comparable. *Esto no significa que si hubieran errores los habría arreglado* por lo que dicho proceso debe hacerse tras descartar la ausencia de arrays problemáticos

### 4.2.1. Filtraje

El filtraje *no específico* permite eliminar los genes que varían poco entre condiciones o que deseamos quitar por otras razones como por ejemplo que no disponemos de anotación para ellos. La función `nsFilter` permite eliminar los genes que, o bien varían poco, o bien no se dispone de anotación para ellos.

Si al filtrar deseamos usar las anotaciones, o la falta de ellas, como criterio de filtraje debemos disponer del correspondiente paquete de anotaciones.

```
> require(genefilter)
> filtered <- nsFilter(eset_rma, require.entrez=TRUE,
+   remove.dupEntrez=TRUE, var.func=IQR,
+   var.cutoff=0.5, var.filter=TRUE,
+   filterByQuantile=TRUE, feature.exclude=~AFFX")
```

La función `nsFilter` devuelve los valores filtrados en un objeto `expressionSet` y un informe de los resultados del filtraje.

```
> names(filtered)

[1] "eset"          "filter.log"

> class(filtered$eset)

[1] "ExpressionSet"
attr(,"package")
[1] "Biobase"

> print(filtered$filter.log)

$numDupsRemoved
[1] 7659

$numLowVar
[1] 6352

$numRemoved.ENTREZID
[1] 1910

$feature.exclude
[1] 10

> eset_filtered <-filtered$eset
```

Podemos grabar el objeto `eset_rma` y los datos filtrados para su posterior uso.

### 4.2.2. Archivos de resultados normalizados

El resultado de la normalización y el filtraje es un objeto `expressionSet` almacenado en un archivo binario `datos.normalizados.Rda`, que será la base para todos los estudios posteriores.



```
> save(eset_rma, eset_filtered, file=file.path(resultsDir, "datos.normalizados.Rda"))
```

Es muy común también que los investigadores deseen ver los datos normalizados por lo que también se pueden guardar en un archivo de texto para facilitar su posterior inspección. El archivo `Datos.Normalizados.Filtrados.csv2` contiene los datos normalizados y filtrados separados por “;” para no generar problemas de formato con el punto decimal en inglés o español.

```
> write.csv2(exprs(eset_rma), file.path(resultsDir, "Datos.Normalizados.csv2"))
```

## 5. Selección de genes diferencialmente expresados

Como en las etapas anteriores la selección de genes diferencialmente expresados (GDE) puede basarse en distintas aproximaciones, desde la  $t$  de Student al programa SAM pasando por multitud de variantes.

En este ejemplo, dado que se realizaran tres comparaciones que luego deseamos comparar entre ellas, se aplicará la aproximación presentada por Smyth *et al.* (2004) basado en la utilización del *modelo lineal general* combinada con un método para obtener una estimación mejorada de la varianza.

### 5.1. Análisis basado en modelos lineales

#### 5.1.1. Matrices de diseño y de contrastes

El primer paso para el análisis basado en modelos lineales es crear la matriz de diseño. Básicamente se trata de una tabla que describe la asignación de cada muestra a un grupo. Tiene tantas filas como muestras y tantas columnas como grupos (si solo se considera un factor). Cada fila contiene un uno en la columna del grupo al que pertenece la muestra y un cero en las restantes.

La matriz de contrastes se utiliza para describir las comparaciones entre grupos. Consta de tantas columnas como comparaciones y tantas filas como grupos (es decir como columnas de la matriz de diseño). Una comparación entre grupos –llamada “contraste”– se representa con un “1” y un “-1” en las filas de los grupos a comparar y ceros en las restantes. Si varios grupos intervinieran en la comparación se tendría tantos coeficientes como grupos con la única restricción de que su suma sería cero.

La matriz de diseño puede definirse manualmente o a partir de un factor creado específicamente para ello.

Manualmente, sería:

```
> design<-matrix(
+           c(1,1,1,1,1,0,0,0,0,0,0,0,0,0,0,
+           0,0,0,0,0,1,1,1,1,1,0,0,0,0,0,
+           0,0,0,0,0,0,0,0,0,0,1,1,1,1,1),
+           nrow=15,byrow=F)
> colnames(design)<-c("A", "B", "L")
> rownames(design) <- sampleNames
> print(design)
```

```

      A B L
PF14.A 1 0 0
PF19.A 1 0 0
PF22.A 1 0 0
PF23.A 1 0 0
PF39.A 1 0 0
PF18.B 0 1 0
PF24.B 0 1 0
PF29.B 0 1 0
PF41.B 0 1 0
PF42.B 0 1 0
PF01.L 0 0 1
PF03.L 0 0 1
PF13.L 0 0 1
PF37.L 0 0 1
PF47.L 0 0 1

```

Las comparaciones que nos interesan son las diferencias, dos a dos entre cada tipo de tumor lo que puede hacerse con los contrastes siguientes:

```

> require(limma)
> cont.matrix <- makeContrasts (
+     AvsB = B-A,
+     AvsL = L-A,
+     BvsL = L-B,
+     levels=design)
> print(cont.matrix)

```

```

      Contrasts
Levels AvsB AvsL BvsL
      A   -1   -1    0
      B    1    0   -1
      L    0    1    1

```

### 5.1.2. Estimación del modelo y selección de genes

Una vez definida la matriz de diseño y los contrastes podemos pasar a estimar el modelo, estimar los contrastes y realizar las pruebas de significación que nos indiquen, para cada gen y cada comparación, si puede considerarse diferencialmente expresado.

El método implementado en `limma` amplía el análisis tradicional utilizando modelos de Bayes empíricos para combinar la información de toda la matriz de datos y de cada gen individual y obtener estimaciones de error mejoradas.

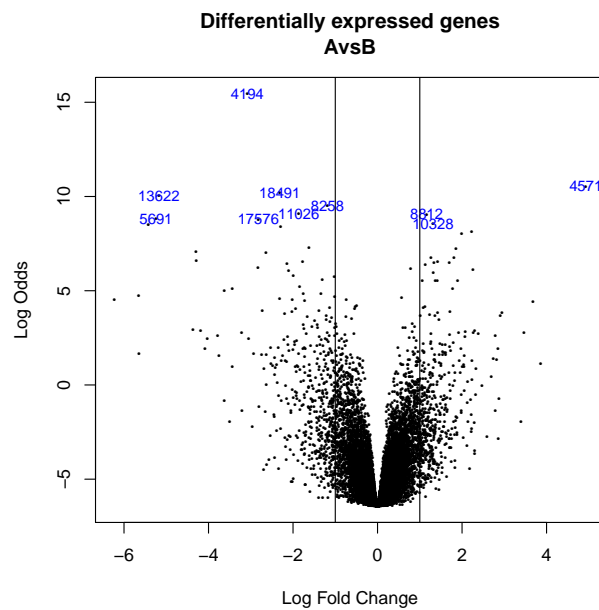
El análisis proporciona los estadísticos de test habituales como **Fold-change** *t*-moderados o *p*-valores ajustados que se utilizan para ordenar los genes de mas a menos diferencialmente expresados.

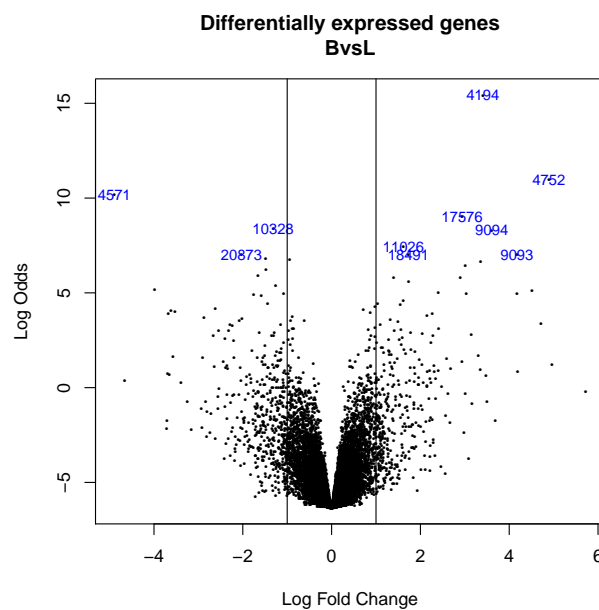
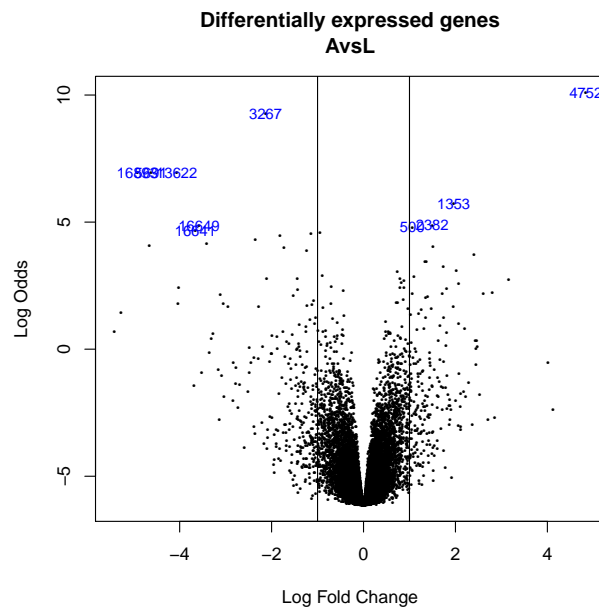
A fin de controlar el porcentaje de falsos positivos que puedan resultar del alto numero de contrastes realizados simultaneamente los *p*-valores se ajustan de forma que tengamos control sobre la tasa de falsos positivos utilizando el metodo de Benjamini y Hochberg ([1]).

La función `topTable` genera para cada contraste una lista de genes ordenados de mas a menos diferencialmente expresados.

```
> topTab_AvsB <- topTable (fit.main, number=nrow(fit.main), coef="AvsB", adjust="fdr")
> topTab_AvsL <- topTable (fit.main, number=nrow(fit.main), coef="AvsL", adjust="fdr")
> topTab_BvsL <- topTable (fit.main, number=nrow(fit.main) , coef="BvsL", adjust="fdr")
```

Los volcano-plots de la figuras ?? permiten visualizar si hay muchos o pocos genes con un gran fold-change y significativamente expresados o si este número es bajo. Estos gráficos representa en abscisas los cambios de expresión en escala logarítmica y en ordenadas el “menos logaritmo” del p-valor o alternativamente el estadístico  $B$  (“log-odds”).





### 5.1.3. Principales genes expresados diferencialmente

Si consideráramos genes expresados diferencialmente aquellos que tienen un p-valor ajustado inferior a 0.05 o 0.01 el número de genes diferencialmente expresado en cada caso sería:

Numero de genes con un p-valor inferior a 0.05 en cada comparación:

En la comparación 'A vs B': 804

En la comparación 'A vs L': 74

En la comparación 'B vs L': 348

Numero de genes con un p--valor inferior a 0.01 en cada comparación:

En la comparación 'A vs B': 162

En la comparación 'A vs L': 20

En la comparación 'B vs L': 78

Los 10 genes más expresados diferencialmente en cada comparación se muestran en las tablas siguientes:

	ID	logFC	AveExpr	t	P.Value	adj.P.Val	B
	4194	204667_at	-3.08	8.83	-21.93	0.00	15.44
	4571	205044_at	4.93	7.75	12.74	0.00	10.54
	18491	219127_at	-2.31	8.86	-12.34	0.00	10.22
	13622	214243_s_at	-5.17	9.22	-12.12	0.00	10.03
	8258	208763_s_at	-1.19	10.49	-11.53	0.00	9.51
	11026	211621_at	-1.86	7.55	-11.07	0.00	9.07
	8812	209318_x_at	1.16	6.17	11.05	0.00	9.05
	5691	206165_s_at	-5.24	7.72	-10.80	0.00	8.81
	17576	218211_s_at	-2.82	10.85	-10.79	0.00	8.80
	10328	210873_x_at	1.31	6.06	10.53	0.00	8.54

Cuadro 2: 10 genes más expresados diferencialmente en la comparación "A vs B"

	ID	logFC	AveExpr	t	P.Value	adj.P.Val	B
	4752	205225_at	4.84	9.30	13.57	0.00	10.09
	3267	203740_at	-2.13	10.81	-12.30	0.00	9.26
	16893	217528_at	-4.92	6.57	-9.54	0.00	6.95
	5691	206165_s_at	-4.62	7.72	-9.52	0.00	6.94
	13622	214243_s_at	-4.06	9.22	-9.52	0.00	6.93
	1353	201825_s_at	1.96	9.08	8.37	0.00	5.71
	2382	202854_at	1.50	10.33	7.68	0.00	4.87
	16649	217284_x_at	-3.58	9.20	-7.67	0.00	4.87
	500	200972_at	1.06	10.57	7.62	0.00	4.80
	16641	217276_x_at	-3.66	8.91	-7.51	0.00	4.66

Cuadro 3: 10 genes más expresados diferencialmente en la comparación "A vs L"

	ID	logFC	AveExpr	t	P.Value	adj.P.Val	B
	4194	204667_at	3.41	8.83	24.26	0.00	15.42
	4752	205225_at	4.89	9.30	13.73	0.00	10.97
	4571	205044_at	-4.89	7.75	-12.64	0.00	10.18

17576	218211_s_at	2.94	10.85	11.25	0.00	0.00	9.02
10328	210873_x_at	-1.32	6.06	-10.54	0.00	0.00	8.37
9094	209604_s_at	3.61	10.80	10.48	0.00	0.00	8.31
11026	211621_at	1.62	7.55	9.61	0.00	0.00	7.42
20873	221510_s_at	-2.03	8.82	-9.26	0.00	0.00	7.03
18491	219127_at	1.73	8.86	9.24	0.00	0.00	7.01
9093	209603_at	4.17	8.62	9.23	0.00	0.00	6.99

Cuadro 4: 10 genes más expresados diferencialmente en la comparación “B vs L”

## 6. Post-procesado de las listas de genes obtenidas

Una vez obtenidas las listas de genes diferencialmente expresados pueden llevarse a cabo todo tipo de análisis sobre ellas, generalmente encaminados a facilitar la interpretación de los resultados.

Entre estas exploraciones –que podemos llamar genericamente “post-procesado de las listas” se encuentra

- La anotación de las listas de genes en diversas bases de datos.
- La comparación entre las listas para determinar que genes cambian simultáneamente -o no- en varias comparaciones (o cuales cambian en una comparación pero no en otra).
- La visualización de todos los genes seleccionados en varias comparaciones para, de forma similar a lo anterior, detectar grupos de genes con patrones de cambio similares –o distintos– entre distintas comparaciones.
- EL análisis de significación biológica de las listas mediante análisis de enriquecimiento o mediante “gene set analysis” para detectar si las listas se encuentran enriquecidas en genes asociados a funciones o procesos biológicos determinados.

### 6.1. Comparaciones múltiples

Cuando se realizan varias comparaciones a la vez puede resultar importante ver que genes cambian simultáneamente en más de una comparación. Si el número de comparaciones es alto también puede ser necesario realizar un ajuste de p-valores entre las comparaciones, distinto del realizado entre genes.

La función `decidetests` permite realizar ambas cosas. En este ejemplo no se ajustaran los p-valores entre comparaciones. Tan solo se seleccionaran los genes que cambian en una o más condiciones.

EL resultado del análisis es una tabla `res` que para cada gen y cada comparación contiene un 1 (si el gen esta sobre-expresado o “up” en esta condicion), un 0 (si no hay cambio significativo) o un -1 (si esta “down”-regulado).

Para resumir dicho análisis podemos contar qué filas tienen como mínimo una celda distinta de cero:

```
> vennDiagram (res.selected[,1:3], main="Genes in common #1", cex=0.9)
```

**Genes in common #1**

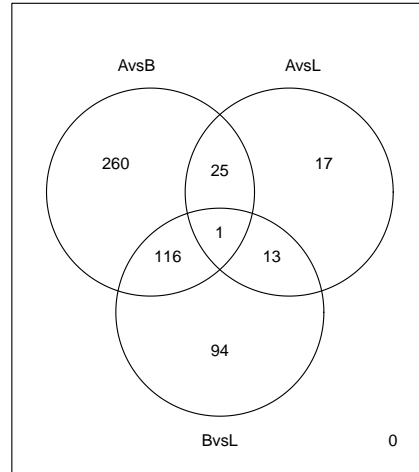


Figura 11: Número de genes seleccionado en cada comparación

```
> sum.res.rows<-apply(abs(res),1,sum)
> res.selected<-res[sum.res.rows!=0,]
> print(summary(res))
```

	AvsB	AvsL	BvsL
-1	81	13	30
0	22151	22264	22220
1	51	6	33

En vista de estos valores podemos aplicar otros criterios de selección, por ejemplo genes con un p-valor ajustado inferior a 0.05 y **log Fold change** mayor o igual a 1. *Este criterio combina a la vez la significación estadística y la significación biológica por lo que, en un estudio real sería probablemente el escogido.*

	AvsB	AvsL	BvsL
-1	228	29	113
0	21881	22227	22059
1	174	27	111

Un diagrama de Venn permite visualizar la tabla anterior sin diferenciar entre genes “up” o “down” regulados.

## 6.2. Anotación de resultados

La identificación de los genes seleccionados puede resultar más sencilla para el especialista en un campo si se utilizan nombres estándar como el símbolo del

gen o “gene symbol”. Con esta finalidad cada paquete de anotaciones tiene tablas de correspondencia entre los distintos tipos de identificadores, principalmente entre los del array y los de otras bases de datos.

Para saber que anotaciones estan disponibles debe cargarse el paquete y llamar la función del mismo nombre

```
> require(hgu133a.db)
> hgu133a()
```

Quality control information for hgu133a:

This package has the following mappings:

```
hgu133aACCNUM has 22283 mapped keys (of 22283 keys)
hgu133aALIAS2PROBE has 43839 mapped keys (of 92418 keys)
hgu133aCHR has 20366 mapped keys (of 22283 keys)
hgu133aCHRLNGTHS has 93 mapped keys (of 93 keys)
hgu133aCHRLOC has 20152 mapped keys (of 22283 keys)
hgu133aCHRLOCEND has 20152 mapped keys (of 22283 keys)
hgu133aENSEMBL has 19785 mapped keys (of 22283 keys)
hgu133aENSEMBL2PROBE has 13042 mapped keys (of 19985 keys)
hgu133aENTREZID has 20373 mapped keys (of 22283 keys)
hgu133aENZYME has 3109 mapped keys (of 22283 keys)
hgu133aENZYME2PROBE has 905 mapped keys (of 975 keys)
hgu133aGENENAME has 20373 mapped keys (of 22283 keys)
hgu133aGO has 19409 mapped keys (of 22283 keys)
hgu133aGO2ALLPROBES has 15392 mapped keys (of 15942 keys)
hgu133aGO2PROBE has 11827 mapped keys (of 12445 keys)
hgu133aMAP has 20312 mapped keys (of 22283 keys)
hgu133aOMIM has 17061 mapped keys (of 22283 keys)
hgu133aPATH has 8125 mapped keys (of 22283 keys)
hgu133aPATH2PROBE has 229 mapped keys (of 229 keys)
hgu133aPFAM has 20227 mapped keys (of 22283 keys)
hgu133aPMID has 20311 mapped keys (of 22283 keys)
hgu133aPMID2PROBE has 313900 mapped keys (of 334784 keys)
hgu133aPROSITE has 20227 mapped keys (of 22283 keys)
hgu133aREFSEQ has 20275 mapped keys (of 22283 keys)
hgu133aSYMBOL has 20373 mapped keys (of 22283 keys)
hgu133aUNIGENE has 20328 mapped keys (of 22283 keys)
hgu133aUNIPROT has 19818 mapped keys (of 22283 keys)
```

Additional Information about this package:

```
DB schema: HUMANCHIP_DB
DB schema version: 2.1
Organism: Homo sapiens
Date for NCBI data: 2012-Mar7
Date for GO data: 20120303
```



Date for KEGG data: 2011-Mar15  
Date for Golden Path data: 2010-Mar22  
Date for IPI data: 2011-Sep27  
Date for Ensembl data: 2012-Feb24

Cada tabla de asociación puede consultarse de diversas formas,

- Con las funciones `get` o `mget`.
- Convirtiéndola en una tabla y extrayendo valores
- En algunos casos utilizando funciones específicas como `getSYMBOL` o `getEG` (por “Entrez Gene”) cuando exitan.

Bioconductor dispone de algunos paquetes que permiten aprovechar esta funcionalidad anterior para obtener las anotaciones de cada gen y generar una tabla HTML con enlaces a algunas bases de datos.

De forma sencilla es posible obtener tablas con las anotaciones correspondientes a los genes seleccionados. Si se desea ser más ambicioso es posible generar tablas en las que se combinen hiperenlaces a las anotaciones con los resultados de la selección de genes.

Veamos las distintas posibilidades, en orden de complejidad creciente

#### 6.2.1. Tablas de anotaciones sencillas

El paquete `annafy` permite de forma muy simple generar una tabla de anotaciones con hiperenlaces a las bases de datos para cada anotación seleccionada.

La instrucción siguiente crea una tabla con las anotaciones disponibles para los genes seleccionados en la sección de comparaciones múltiples.

```
> require(annafy)
> genesSelected <- rownames(res.selected)
> at <- aafTableAnn(genesSelected, "hgu133a.db")
> saveHTML(at, file.path(resultsDir, "anotations.html"),
  "Annotations for selected genes")
```

#### 6.2.2. Archivos de resultados con tablas hiperenlazables

La función `htmlpage` permite generar un archivo html a partir de un `data.frame` con una funcionalidad muy interesante: podemos hacer que una o varias de sus columnas contengan hiperenlaces a alguna base de datos permitiendo acceder a gran cantidad de información sobre los genes de la lista.

Para aprender más sobre su uso se puede leer la viñeta `How to get pretty output` disponible en Bioconductor.

Mediante un bucle apropiado creamos una tabla anotada para cada comparación.

```
> listOfTables <- list(AvsB = topTab_AvsB, AvsL = topTab_AvsL, BvsL = topTab_BvsL)
> for (i in 1:length(listOfTables)){
+   # Seleccionamos la "topTable"
+   topTab <- listOfTables[[i]]
+   # Escogemos los grupos de sondas a incluir en la tabla
```

```

+   whichGenes<-topTab["P.Value"]<0.05
+   selectedIDs <- topTab$ID[whichGenes]
+   # Los convertimos a identificadores Entrez ("EG") y a Gene Symbols
+   genes<- getEG(selectedIDs, "hgu133a.db")
+   symbols <-getSYMBOL(selectedIDs, "hgu133a.db")
+   # Haremos la columna de Entrez sea hiperenlazable
+   paraEnlace <- list (misgenes=genes)
+   # Preparamos el data.frame con el que se creará el archivo de resultados
+   otherNames = data.frame(selectedIDs, symbols, topTab[whichGenes,-1])
+   names(otherNames) = c("Affy ID", "Gene Symbol", colnames(topTab)[-1])
+   # Invocamos la función "htmlpage"
+   comparison <- names(listOfTables)[i]
+   htmlpage(paraEnlace,
+             filename =file.path(resultsDir,
+                                   paste("Selected Genes in comparison ",comparison,".html", sep="")) ,
+             title = paste("Diff. expressed genes in comparison ", comparison, sep=""),
+             othertnames = otherNames,
+             table.head = c("Entrez IDs", names(otherNames)),
+             table.center = TRUE,
+             repository=list("en"))
+ }

```

### 6.3. Visualización de los perfiles de expresión

Tras seleccionar los genes diferencialmente expresados podemos visualizar las expresiones de cada gen agrupándolas para destacar los genes que se encuentran up o down regulados simultáneamente constituyendo *perfiles de expresión*.

Hay distintas formas de visualización pero aquí tan sólo se presenta el uso de mapas de color o Heatmaps.

En primer lugar seleccionamos los genes a visualizar: Se toman todos aquellos que han resultado diferencialmente expresados en alguna de las tres comparaciones.

```

> probeNames<-rownames(res)
> probeNames.selected<-probeNames[sum.res.rows!=0]
> exprs2cluster <-exprs(eset_rma)[probeNames.selected,]
> colnames(exprs2cluster)<-sampleNames
> color.map <- function(grupo) {
+   if (grupo=="A"){
+     c<- "yellow"
+   }else{
+     if (grupo=="B"){
+       c<- "red"
+     }else{
+       c<- "blue"
+     }
+   }
+ }
+ return(c)}

```

Para representar el Heatmap tan sólo necesitamos la matriz de datos resultante.

```
> grupColors <- unlist(lapply(pData(eset_rma)$Group, color.map))
> heatmap(exprs2cluster, col=rainbow(100), ColSideColors=grupColors, cexCol=0.9)
```

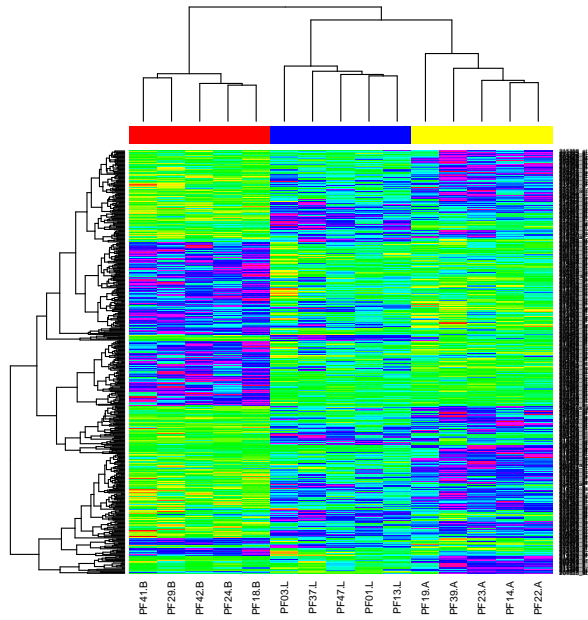


Figura 12: Mapa de colores basado en los genes seleccionados por estar diferencialmente expresados. Como puede verse los tumores Apocrinos y Luminales tienen perfiles de expresión más parecidos entre ellos que cada uno con los de tipo Basal

Si se desea realizar mapas de color más sofisticados puede utilizarse el paquete `gplots` que implementa una versión mejorada en la función `heatmap.2`

```
> grupColors <- unlist(lapply(pData(eset_rma)$Group, color.map))
> require("gplots")
> heatmap.2(exprs2cluster,
+           col=bluered(75), scale="row",
+           ColSideColors=grupColors, key=TRUE, symkey=FALSE,
+           density.info="none", trace="none", cexCol=1)
```

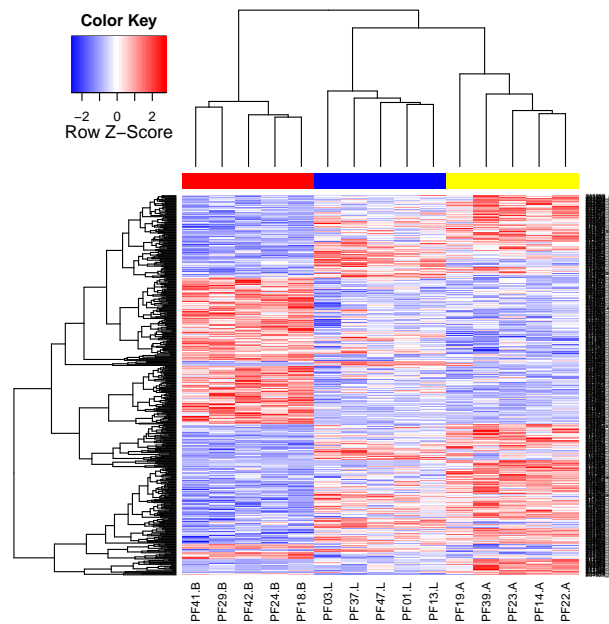


Figura 13: Mapa de colores mejorado basado en los genes seleccionados por estar diferencialmente expresados. Como puede verse los tumores Apocrinos y Luminales tienen perfiles de expresión más parecidos entre ellos que cada uno con los de tipo Basal

#### 6.4. Análisis de la significación biológica

El análisis de significación biológica busca establecer si, dada una lista de genes seleccionados por estar diferencialmente expresados entre dos condiciones, las funciones y procesos biológicos que los caracterizan aparecen en esta lista con mayor frecuencia que entre el resto de genes analizados.

Se han desarrollado multitud de variantes de estos tipos de análisis ([8]) pero aquí utilizaremos el análisis básico de enriquecimiento tal como se describe en los trabajos de Gentleman ([5]) implementados en el paquete `G0stats` de Bioconductor.

El análisis se realiza sobre dos bases de datos de anotaciones, la "Gene Ontology" o la "Kyoto Encyclopedia of Genes and Genomes".

Los análisis de este tipo necesitan un número mínimo de genes para resultar fiables por lo que se incluyan en todos los genes con p-valores ajustados inferiores a 0.05 (sin filtrar por mínimo “fold-change”).

```
> require(GOstats)
> listOfTables <- list(AvsB = topTab_AvsB, AvsL = topTab_AvsL, BvsL = topTab_BvsL)
> for (i in 1:length(listOfTables)){
+   # Seleccionamos la "topTable"
+   topTab <- listOfTables[[i]]
+   # Definimos el universo de genes: todos los que se han incluido en el análisis
+   # EL programa trabaja con identificadores "entrez" y no admite duplicados
+
+   entrezUniverse = unique(getEG(as.character(topTab$ID), "hgu133a.db"))
+
+   # Escogemos los grupos de sondas a incluir en el análisis
+   # Este análisis trabaja bien con varios centenares de genes
+   # por lo que es habitual basarse en p-valores sin ajustar para incluirlos
+
+   whichGenes<-topTab["adj.P.Val"]<0.05
+   geneIds <- unique(getEG(as.character(topTab$ID[whichGenes]),"hgu133a.db"))
+
+   # Creamos los "hiperparámetros" en que se basa el análisis
+   GOparams = new("GOHyperGParams",
+     geneIds=geneIds, universeGeneIds=entrezUniverse,
+     annotation="org.Hs.eg.db", ontology="BP",
+     pvalueCutoff=0.001, conditional=FALSE,
+     testDirection="over")
+   KEGGparams = new("KEGGHyperGParams",
+     geneIds=geneIds, universeGeneIds=entrezUniverse,
+     annotation="org.Hs.eg.db",
+     pvalueCutoff=0.01, testDirection="over")
+
+   # Ejecutamos los análisis
+
+   GOhyper = hyperGTest(GOparams)
+   KEGGhyper = hyperGTest(KEGGparams)
+
+   # Creamos un informe html con los resultados
+   comparison = names(listOfTables)[i]
+   GOfilename =file.path(resultsDir,
+     paste("GOResults.",comparison,".html", sep=""))
+   KEGGfilename =file.path(resultsDir,
+     paste("KEGGResults.",comparison,".html", sep=""))
+   htmlReport(GOhyper, file = GOfilename, summary.args=list("htmlLinks"=TRUE))
+   htmlReport(KEGGhyper, file = KEGGfilename, summary.args=list("htmlLinks"=TRUE))
+ }
```

## 7. Resumen de resultados y discusión

Los controles de calidad llevados a cabo en la primera parte de este estudio han permitido establecer que los datos con los que se ha trabajado eran de buena calidad, a pesar de presentar una cierta heterogeneidad que queda ilustrada en la incompleta agrupación de las muestras segun sus tipos de tumores.

Los análisis llevados a cabo sobre los datos normalizados y filtrados han permitido detectar un cierto número de genes diferencialmente expresados (ver tabla 5) que se prestan a estudios posteriores acerca de su significación biológica.

Cuadro 5: Genes diferencialmente expresados basados en un cirte definido por un p-valor ajustado (FDR) inferior a 0.05 y un cambio de expresión de cómo mínimo 1.

	A vs B	A vs L	B vs L
Genes “down-regulados”	228	29	113
Genes “sobre-expresados”	174	27	111

La lista completas de genes seleccionados (ordenados de mayor a menor p-valor) pueden verse en los archivos “Selected genes in comparison AvsB.html”, “Selected genes in comparison AvsL.html” y “Selected genes in comparison BvsL.html”.

El análisis de significación biológica ha permitido detectar una serie de funciones y procesos biológicos que caracterizan las listas de genes seleccionadas en las bases de datos de anotaciones más populares, la Gene Ontology (GO) y la “Kyoto Encyclopedia of Genes and Genomes”. Estas listas pueden verse en las tablas “GOResults.XXvsYY.html” y “KEGGResults.XXvsYY.html” (XXvsYY es AvsB, AvsL o BvsL) pero no se discuten aquí.

### 7.1. Discusión

El estudio que se ha presentado en este documento es un análisis estándar de microarrays y como tal adolece de sus ventajas e inconvenientes.

Como inconveniente podemos destacar dos cosas

- El tamaño de las muestras utilizadas es bastante limitado lo que determina que el estudio tenga poca potencia por lo que *probablemente* habrá menos reproducibilidad y más falsos negativos de los que sería deseable si se utilizara un mayor número de muestras.
- En cada paso del proceso se han tomado decisiones relativamente arbitrarias acerca de los métodos a seguir para la normalización, filtradom selección de genes, etc. La decisión de si estos métodos son los más adecuados o no es probablemente subjetiva (véase por ejemplo los artículos [2, 12]) por lo que sería interesante saber como cambian los resultados si se tomaran otras decisiones.

Los problemas anteriores no son, sin embargo, problemas de este estudio concreto, sino en general de los estudios basados en microarrays por lo que, limitaciones aparte, el estudio aportará probablemente información valiosa que permitirá un seguimiento posterior del problema.

Finalmente debe de tenerse en cuenta que cualquier gen que se acabe considerando realmente expresado diferencialmente se tendrá que verificar mediante

otras técnicas como RT-qPCR por lo que este estudio debe de considerarse como un paso hacia el descubrimiento de genes candidatos, no como una fase definitiva.

## Referencias

- [1] Y. Benjamini and Y. Hochberg. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society B*, 57:289–300, 1995.
- [2] Sung Choe, Michael Boutros, Alan Michelson, George Church, and Marc Halfon. Preferred analysis methods for affymetrix genechips revealed by a wholly defined control dataset. *Genome Biology*, 6(2):R16, 2005.
- [3] S. Falcon and R. Gentleman. Using gostats to test gene lists for go term association. *Bioinformatics*, 23(2):257–258, 2007.
- [4] Pierre Farmer, Herve Bonnefoi, Veronique Becette, Michele Tubiana-Hulin, Pierre Fumoleau, Denis Larsimont, Gaetan Macgrogan, Jonas Bergh, David Cameron, Darlene Goldstein, Stephan Duss, Anne-Laure Nicoulaz, Cathrin Briskin, Maryse Fiche, Mauro Delorenzi, and Richard Iggo. Identification of molecular apocrine breast tumours by microarray analysis. *Oncogene*, 24(29):4660–4671, July 2005. PMID: 15897907.
- [5] R. Gentleman. Using go for statistical analysis. *Bioconductor's Compendiums* ([www.bioconductor.org](http://www.bioconductor.org)), 2004.
- [6] Robert Gentleman, Vince Carey, Wolfgang Huber, Rafael Irizarry, and Sandrine Dudoit. *Bioinformatics and Computational Biology Solutions using R and Bioconductor*. Springer, New York, 2005.
- [7] R. A. Irizarry, B. Hobbs, F. Collin, Y. D. Beazer-Barclay, K. J. Antonellis, U. Scherf, and T. P. Speed. Exploration, normalization, and summaries of high density oligonucleotide array probe level data. *Biostatistics*, 4:249–264, 2003.
- [8] P. Khatrı and S. Drăghici. Ontological analysis of gene expression data: current tools, limitations, and problems. *Bioinformatics*, 18:3587–3595, 2005.
- [9] Friedrich Leisch. Sweave: Dynamic generation of statistical reports using literate data analysis. In Wolfgang Härdle and Bernd Rönz, editors, *Compstat 2002 — Proceedings in Computational Statistics*, pages 575–580. Physica Verlag, Heidelberg, 2002. ISBN 3-7908-1517-9.
- [10] Friedrich Leisch. Sweave, part I: Mixing R and L<sup>A</sup>T<sub>E</sub>X. *R News*, 2(3):28–31, December 2002.
- [11] Gordon K Smyth. Linear models and empirical bayes methods for assessing differential expression in microarray experiments. *Stat Appl Genet Mol Biol*, 3:Article3, 2004.

- [12] Qianqian Zhu, Jeffrey Miecznikowski, and Marc Halfon. Preferred analysis methods for affymetrix genechips. ii. an expanded, balanced, wholly-defined spike-in dataset. *BMC Bioinformatics*, 11(1):285, 2010.

## A. Codigo R del análisis

El archivo `Solucion_PEC_Microarrays.R` contiene el código R utilizado para realizar todos los análisis presentados en este informe.

El análisis se ha llevado a cabo utilizando el sistema **Sweave** ([10, 9]) de análisis reproducibles que integra el sistema  $\text{\LaTeX}$  el lenguaje *R* para generar a la vez los análisis y el informe de resultados por lo que el código puede contener alguna instrucción inusual cuando su objetivo es generar una tabla de  $\text{\LaTeX}$ . Aparte de ésto es un código estándar en el que se ha procurado no usar funciones avanzadas sino básicamente lo que se ha discutido en el curso.

## B. Lista de archivos generados en el estudio

La tabla 6 contiene una lista de los archivos generados en este estudio. Una versión mejorada de dicha tabla incluiría una descripción de los mismos así como hiperenlaces que permitan acceder a cada uno.

	Archivo	Descripción
1	anotations.html	Anotaciones de los genes seleccionados
2	arrayQuality/index.html	Archivo de navegación entre los archivos de control de calidad
3	Datos.Normalizados.csv2	Datos normalizados y filtrados
4	GOResults.AvsB.html	Análisis de enriquecimiento (GO) para la comparación A vs B
5	GOResults.AvsL.html	Análisis de enriquecimiento (GO) para la comparación A vs L
6	GOResults.BvsL.html	Análisis de enriquecimiento (GO) para la comparación B vs L
7	KEGGResults.AvsB.html	Análisis de enriquecimiento (KEGG) para la comparación A vs B
8	KEGGResults.BvsL.html	Análisis de enriquecimiento (KEGG) para la comparación B vs L
9	Selected Genes in comparison AvsB.html	Lista de genes en la comparación A vs B
10	Selected Genes in comparison AvsL.html	Lista de genes en la comparación A vs L
11	Selected Genes in comparison BvsL.html	Lista de genes en la comparación B vs L
12	Solucion_PEC_Microarrays.R	Código R utilizado en el análisis

Cuadro 6: Lista de archivos generados en este análisis

Utilizando el paquete **hwriter** es posible generar de forma sencilla una tabla html con hiperenlaces que permita acceder a los distintos archivos generados en el análisis.

```
> listaArchivos <-read.table(file.path(resultsDir,"listaArchivos.txt"),
head=TRUE, sep="\t")
> require(hwriter)
> hwrite(listaArchivos,file.path(resultsDir, "listaArchivos.html"))
```