# Análisis transcriptómicos de la expresión génica

## Máster Universitario en Bioinformática

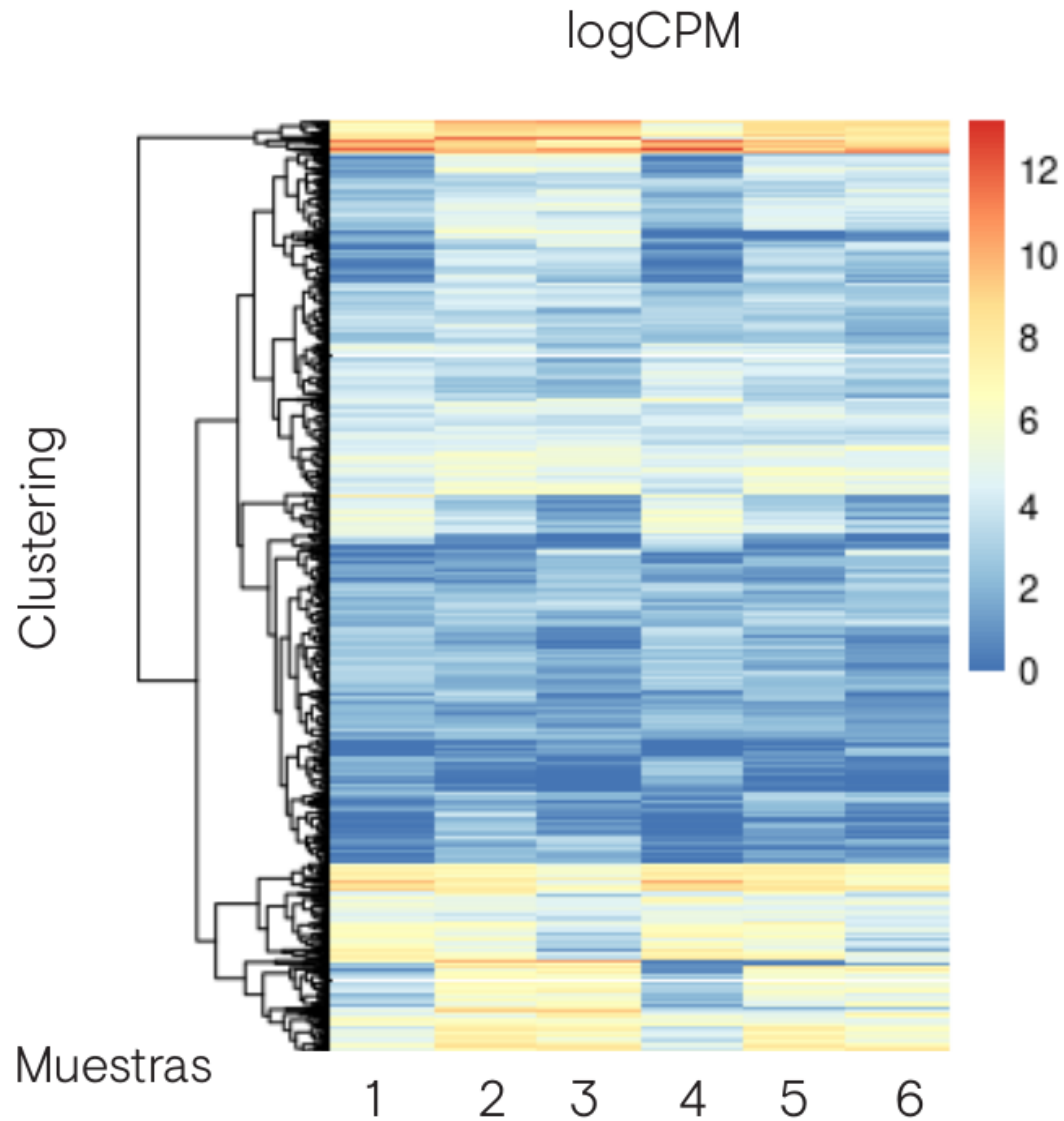**Sesión 12**

Dra. Paula Soler Vila
paula.solerv@professor.universidadviu.com

De:
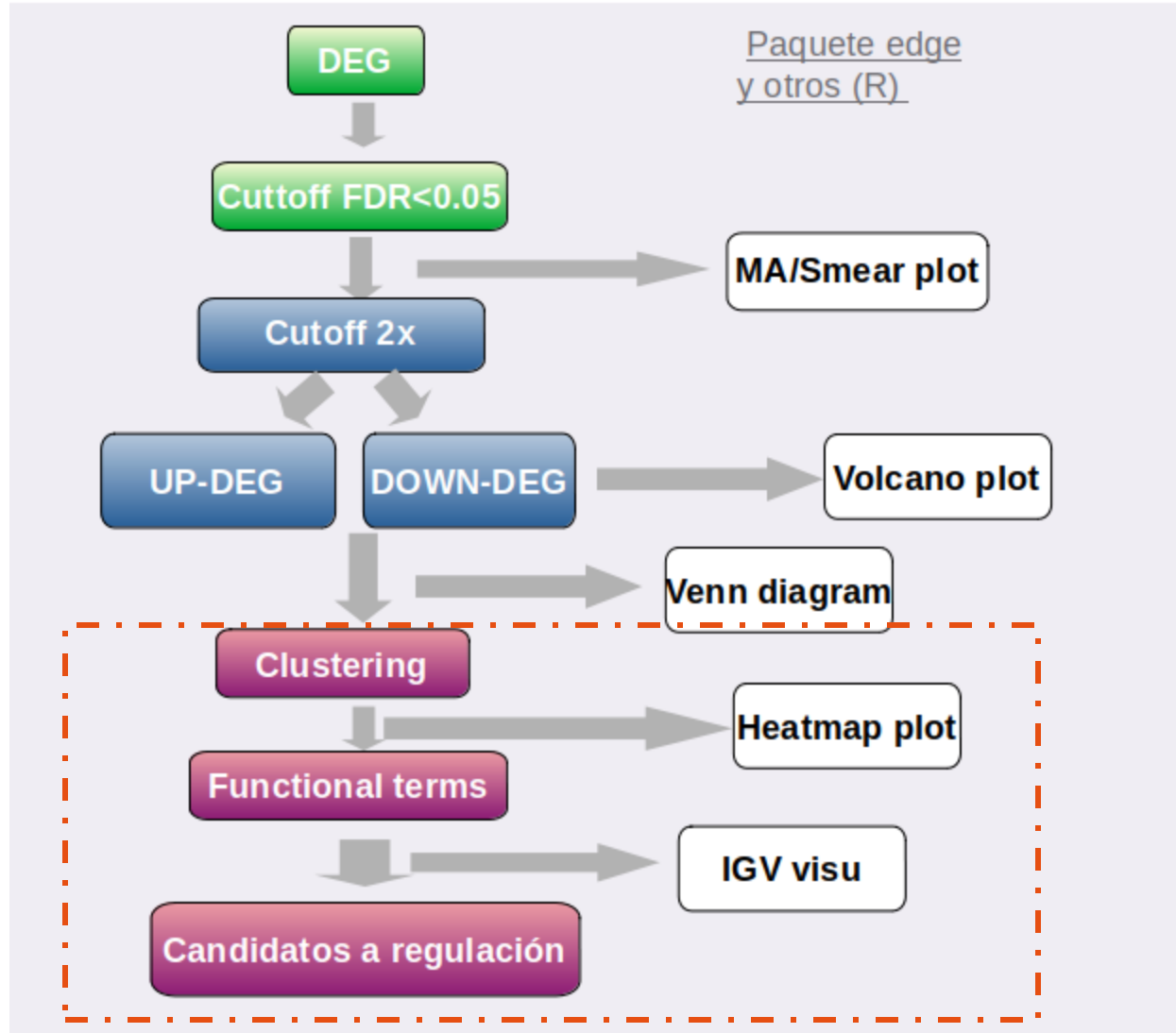Planeta Formación y Universidades

# Bloque V: Exploración y visualización de resultados

# Objetivos

**1** Conocer las posibilidades del "**agrupamiento jerárquico**" en la exploración de datos de RNA-seq y en la visualización de patrones de expresión de genes DE.

**2** Generar gráficos de **heatmap** a partir de recuentos normalizados ajustando los datos a diferentes **escalas** y métodos de cálculo de **distancias**.

**3** Conocer y analizar la **ontología génica** y su aplicación en datos de expresión.

**4** Calcular el **enriquecimiento funcional** de listas de genes DE.

# Flujo de trabajo para la exploración y visualización de resultados

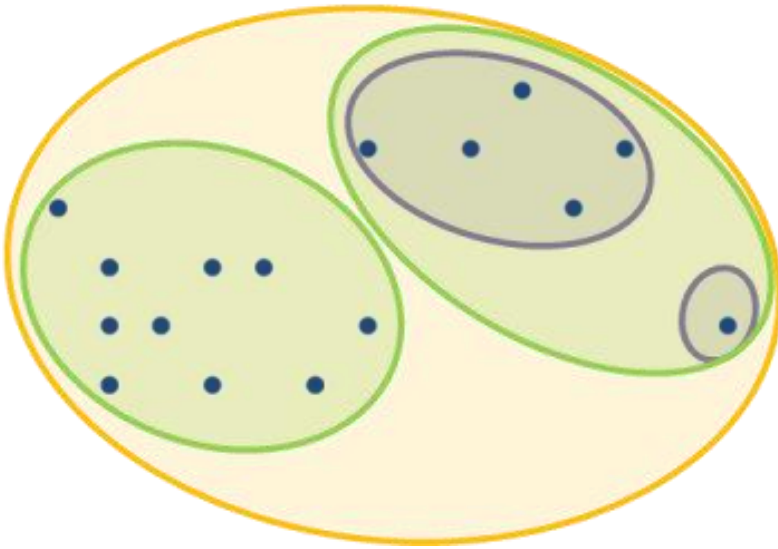# ¿A qué hace referencia el término *clustering*?

Técnicas *no supervisadas* cuya finalidad

es encontrar patrones o grupos (*clusters*)

dentro de un conjunto de observaciones

que comparten patrones comunes

Las observaciones que están dentro de un mismo grupo son similares entre ellas y distintas a las observaciones de otros grupos.

My super advanced intelligence solves **another impossible task!**

# ¿A qué hace referencia el término *clustering*?



**Hierarchical Clustering**

**Partitional Clustering**

***Agrupamiento jerárquico
(Hierarchical Clustering)***

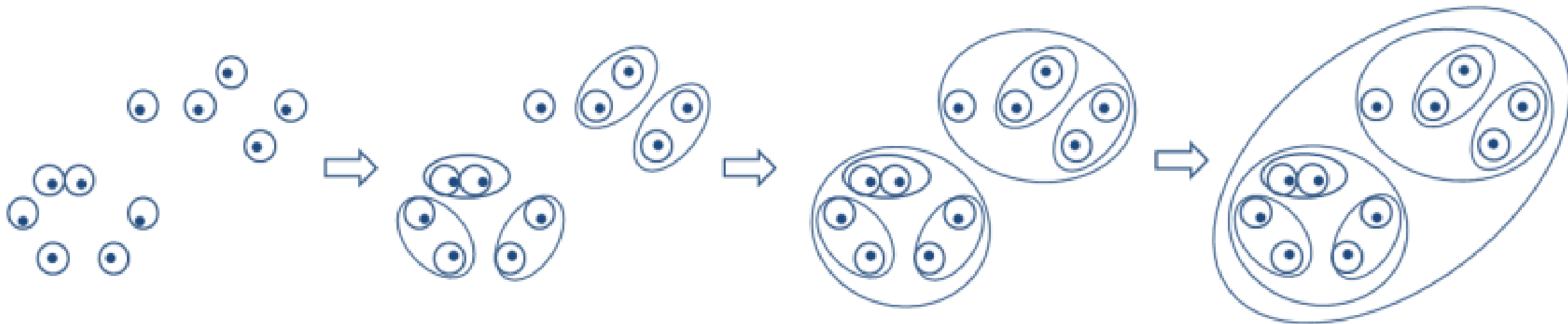- Cada observación pertenece a múltiples grupos
- Hay solapamiento

***Agrupamiento particional
(Partitional Clustering)***

- Cada observación pertenece únicamente a un grupo
- No hay solapamiento
- **K-means clustering**

# ¿A qué hace referencia el término *clustering*?



Agglomerative Hierarchical Clustering

Se van agrupando atendiendo a su cercanía

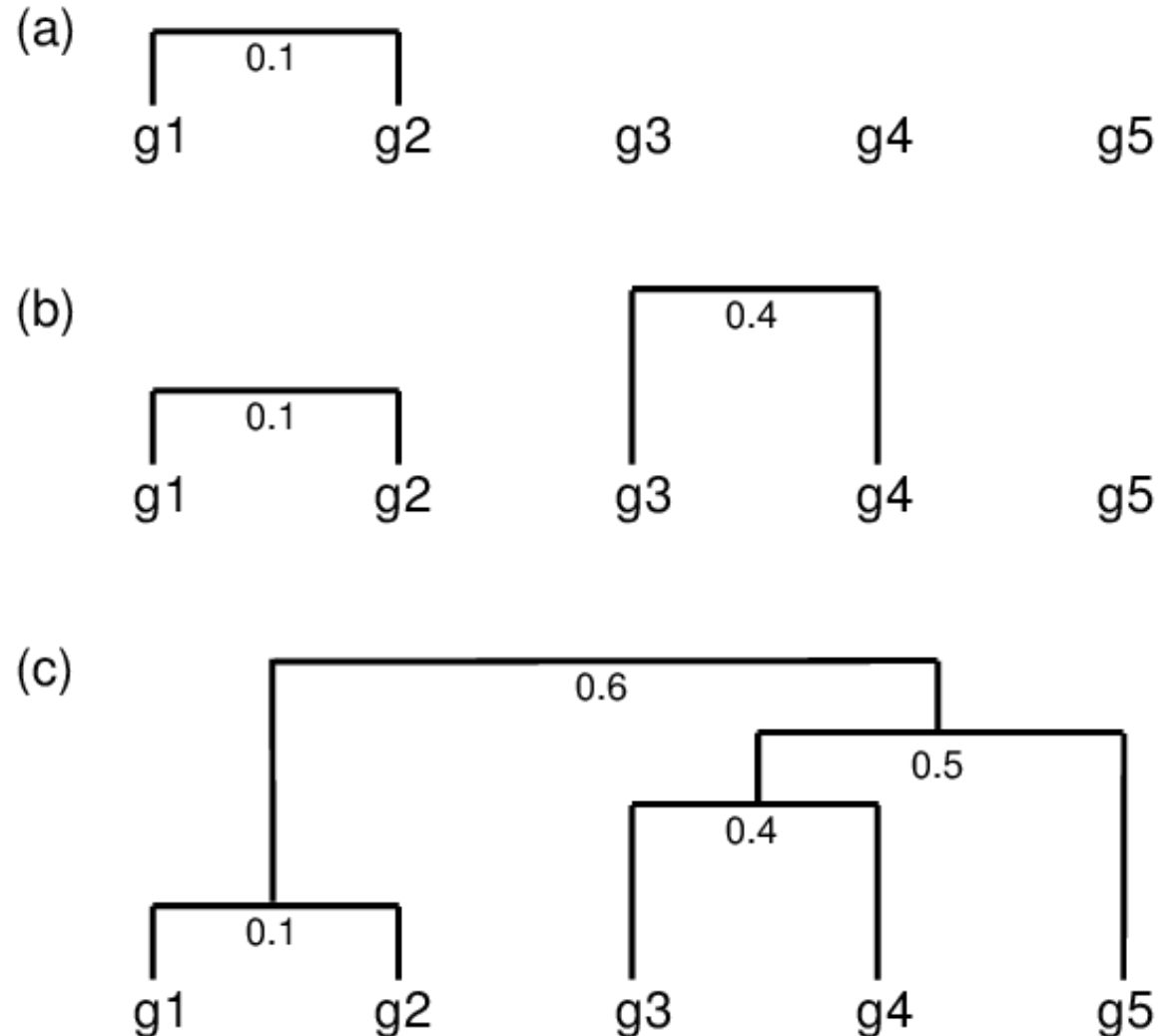# Agrupación jerárquica aglomerativa



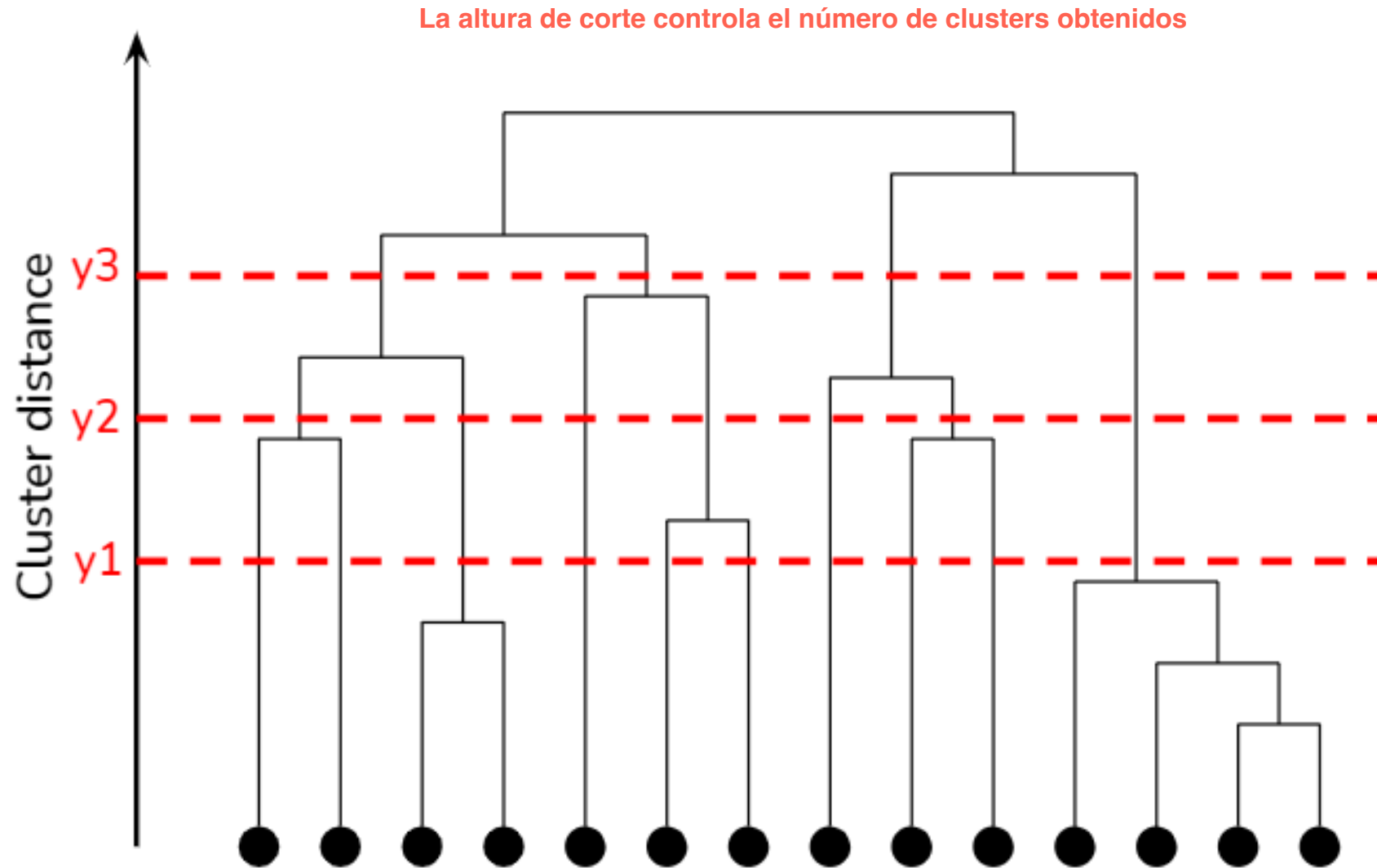1. Inicialización

2. Cálculo de distancias

Medidas de distancias

3. Fusión de clusters
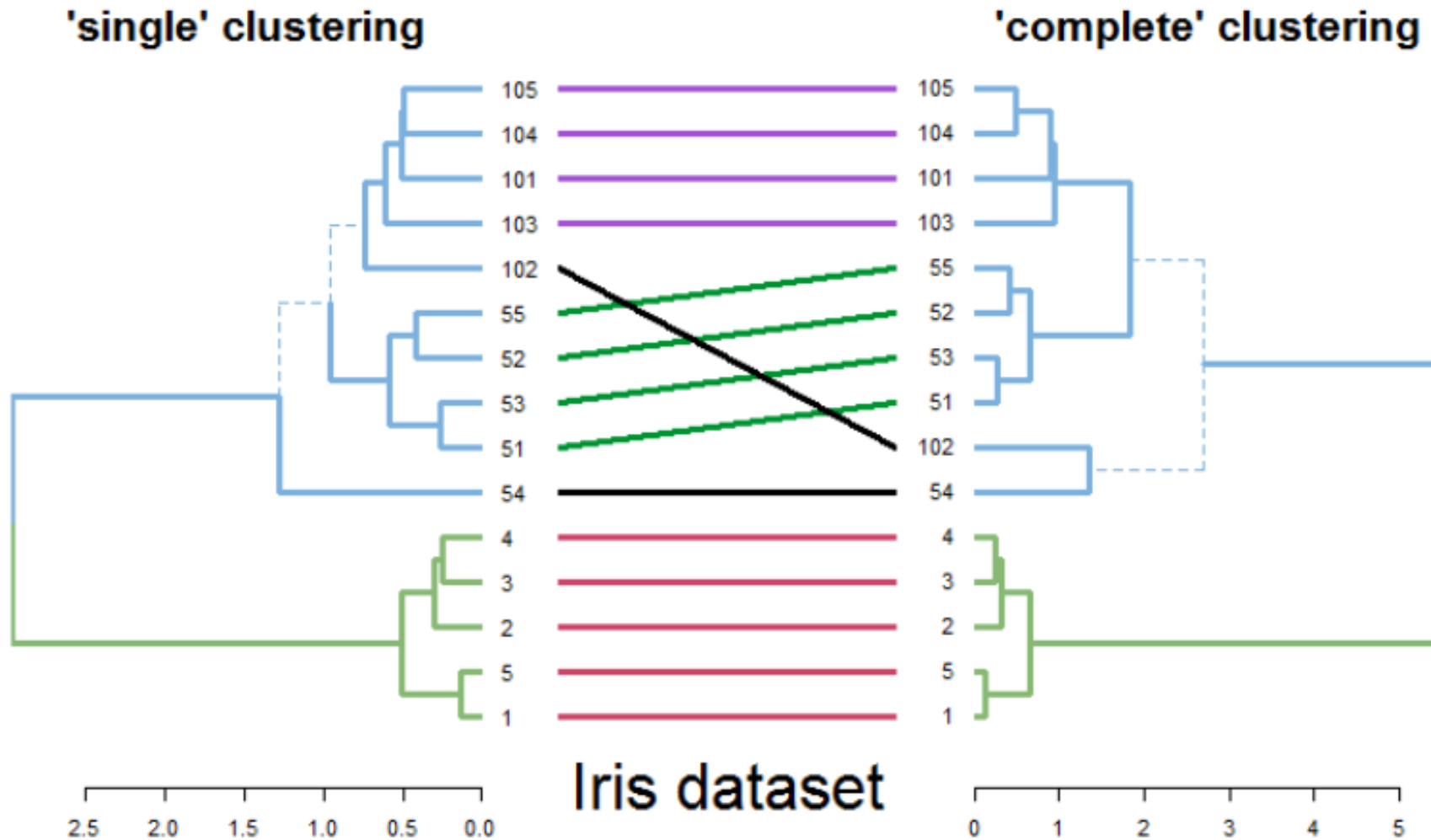
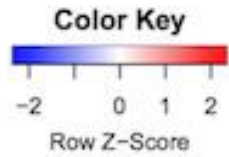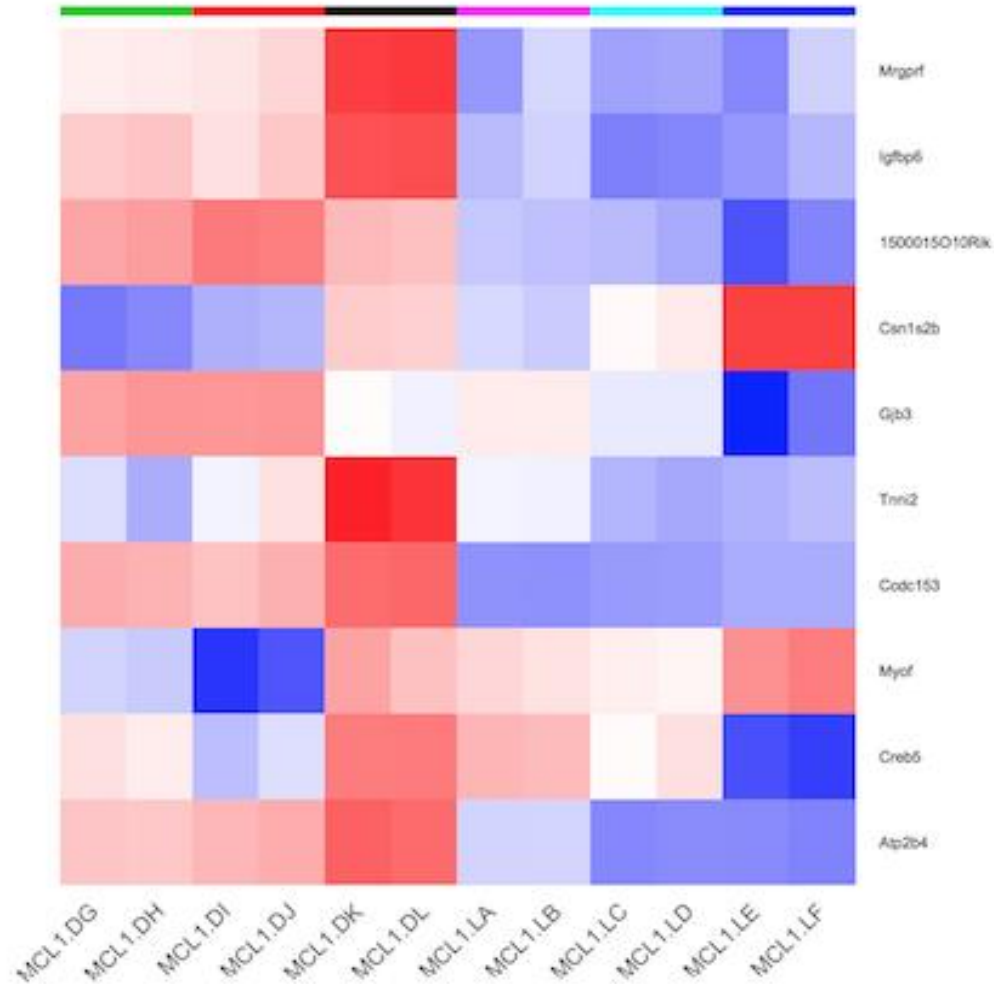Métodos de aglomeramiento (*linkage*)

4. **Dendrograma** final

# Cálculo de las distancias entre observaciones

**Definir y cuantificar la similitud/distancia entre las observaciones**

Table 1: Expresiones genéticas de pacientes.

| Muestra | IRX4 | OCT4 | PAX6 |
|---|---|---|---|
| paciente 1 | 11 | 10 | 1 |
| paciente 2 | 13 | 13 | 3 |
| paciente 3 | 12 | 4 | 10 |
| paciente 4 | 1 | 3 | 9 |



Figure 1: Valores de expresión génica para diferentes pacientes
Paciente 1 y paciente 2 tienen valores de expresión genética similares.

# Cálculo de las **distancias** entre observaciones

Definir y cuantificar la similitud/distancia entre las observaciones

- **Distancia euclídea**

- Distancia de Manhattan

- Índice de Jaccard

- Correlación de Pearson

- Correlación de Spearman ...



Cluster Dendrogram

d
hclust (*, "complete")

# Linkage: métodos de aglomeración

Estrategias que pueden ser empleadas a la hora de unir los *clusters* en las diversas etapas o niveles de un procedimiento jerárquico



**Single linkage: la distancia de similitud viene dada por la mínima distancia entre sus componentes**



**No hay ningún método bueno o malo, no hay ningún método que sea el más óptimo, por tanto lo mejor es hacerlo con varios y comparar la robustez**

**Complete linkage: viene determinado por la similitud de sus elementos más dispares, aquellos genes que se encuentran a una mayor distancia**

# Agrupación jerárquica aglomerativa

1. Inicialización

2. Cálculo de distancias

Medidas de distancias

3. Fusión de clusters

Métodos de aglomeramiento (*linkage*)

4. **Dendrograma** final

# Dendrograma



La altura de corte controla el número de clusters obtenidos

# Comparación de dendrogramas



'single' clustering     'complete' clustering

Iris dataset

*Función: tanglegram*

# Heatmap: Visualización de la expresión de DEG



**Contrast: basalpregnant−basallactate**
**Top 10 genes by adj.P.Val**

**Color Key**
Row Z-Score

1. Extracción de los **conteos normalizados** para estos genes

Colores frios son genes que estarán menos expresados, mientras que los colores cálidos son genes que están mas expresados

2. **Escalado**
Calcular el **Z-score** de los conteos normalizados

# Heatmap: Visualización de la expresión de DEG



① Extracción de los **conteos** **normalizados** para estos genes

logCPM

Cuando escalamos los datos todos los genes tendrán la misma escala, y por tanto, ya son comparables. Si no, genes que están muy expresados no dejan ver otro tipo de relaciones

# Heatmap: Visualización de la expresión de DEG



**Escalado**

② Calcular el **Z-score** de los conteos normalizados

**Scale** *function*

El concepto de escalado es transformar el logCPM en Zscore, que va a calcular la media y la desviación estándar de un gen a lo largo de todas las muestras

# Heatmap: Visualización de la expresión de DEG

## Usage

```
pheatmap(mat, color = colorRampPalette(rev(brewer.pal(n = 7, name =
  "RdYlBu")))(100), kmeans_k = NA, breaks = NA, border_color = "grey60",
  cellwidth = NA, cellheight = NA, scale = "none", cluster_rows = TRUE,
  cluster_cols = TRUE, clustering_distance_rows = "euclidean",
  clustering_distance_cols = "euclidean", clustering_method = "complete",
  clustering_callback = identity2, cutree_rows = NA, cutree_cols = NA,
  treeheight_row = ifelse((class(cluster_rows) == "hclust") || cluster_rows,
  50, 0), treeheight_col = ifelse((class(cluster_cols) == "hclust") ||
  cluster_cols, 50, 0), legend = TRUE, legend_breaks = NA,
  legend_labels = NA, annotation_row = NA, annotation_col = NA,
  annotation = NA, annotation_colors = NA, annotation_legend = TRUE,
  annotation_names_row = TRUE, annotation_names_col = TRUE,
  drop_levels = TRUE, show_rownames = T, show_colnames = T, main = NA,
  fontsize = 10, fontsize_row = fontsize, fontsize_col = fontsize,
  angle_col = c("270", "0", "45", "90", "315"), display_numbers = F,
  number_format = "%.2f", number_color = "grey30", fontsize_number = 0.8
  * fontsize, gaps_row = NULL, gaps_col = NULL, labels_row = NULL,
  labels_col = NULL, filename = NA, width = NA, height = NA,
  silent = FALSE, na_col = "#DDDDDD", ...)
```

# PRACTIQUEMOS

amazon WorkSpaces

Clustering (heatmap)

# Heatmap: Visualización de la expresión de DEG

```r
# 1. Transform the data
logcpm <- cpm(y, log = TRUE)
rownames(logcpm) <- y$genes$SYMBOL
colnames(logcpm) <- paste(y$samples$group, 1:2, sep="-")
head(logcpm)

# 2. Selection of the top genes
DEG <- res_corrected$table[res_corrected$table$FDR <= 0.01 &
                abs(res_corrected$table$logFC) >= 3.5 ,]

DEG_selection <- logcpm[na.omit(DEG$SYMBOL),]

# 3. Creation of pheatmap.
pheatmap(DEG_selection, scale = "row",
     cluster_rows = T,
     cluster_cols = T,
     clustering_distance_rows = "euclidean",
     clustering_distance_cols = "euclidean",
     clustering_method = "ward.D2",
     cutree_cols = 1, cutree_rows = 2,
     display_numbers = T, fontsize_number = 6, fontsize_row = 7, border_color = NA)
```

# Ontología Génica

# Ontología génica

La **ontología génica** (**GO**, por sus siglas en inglés) es el conjunto de términos o de vocabulario estructurado creado para **describir y categorizar los genes y sus productos**.

Actualmente, la mayoría de las principales bases de datos de plantas, animales y microorganismos forman parte del proyecto, y **un total de 45.003 términos** —aplicables a una amplia variedad de organismos biológicos— aparecen registrados en su base de datos.

# Ontología génica se divide en:

**Ontología**

§ La función molecular de los productos génicos.

§ Su participación en los procesos biológicos.

§ Su localización celular.

### Number of GO terms by aspect



**# BP**    **# MF**    **# CC**

# Ontología génica-> Ejemplo (gen de la lactasa)

# Ontología génica-> ejemplo

Gene Ontology    Provided by GOA

| Function | Evidence Code | Pubs |
|---|---|---|
| enables beta-glucosidase activity | IDA | PubMed |
| enables beta-glucosidase activity | ISS | |
| enables cellobiose glucosidase activity | IEA | |
| enables galactosylceramidase activity | ISS | |
| enables glucosylceramidase activity | ISS | |
| enables glycosylceramidase activity | IEA | |
| enables lactase activity | IBA | PubMed |
| enables lactase activity | IDA | PubMed |
| enables lactase activity | IMP | PubMed |
| enables phlorizin hydrolase activity | IDA | PubMed |
| enables protein homodimerization activity | IDA | PubMed |

| Process | Evidence Code | Pubs |
|---|---|---|
| involved_in cellobiose catabolic process | IDA | PubMed |
| involved_in glycosylceramide catabolic process | ISS | |
| involved_in lactose catabolic process | IDA | PubMed |
| involved_in quercetin catabolic process | IDA | PubMed |

| Component | Evidence Code | Pubs |
|---|---|---|
| located_in external side of apical plasma membrane | IDA | PubMed |
| located_in integral component of plasma membrane | TAS | PubMed |
| located_in plasma membrane | TAS | |

# Ontología génica-> ejemplo

# Ontología génica-> ejemplo

# Ontología génica-> ejemplo

# Ontología génica *anotación/análisis*



http://geneontology.org/

10/09/2024

# Análisis de enriquecimiento de Ontología Genética (GO)



*https://github.com/jzsh2000/goana*

# PRACTIQUEMOS

goana

viu

**Universidad**
Internacional
de Valencia

De:

Planeta Formación y Universidades