# Advanced data processing

**3**

---

## Combining data from multiple sources
### *Gene clustering example*

- R has powerful functions to combine heterogeneous data into a single data set
- Gene clustering example data:
  - five sets of differentially expressed genes from various experimental conditions
  - file with names of experimentally verified genes
- Gene clustering exercise:
  1. combine this dataset into a single table and cluster to see which conditions are similar
  2. repeat the clustering but only on a subset of experimentally verified genes

---

## Combining gene tables

- input files have two columns: gene names and fold change
- we want to combine all five tables into a single table, with 0 for missing values

| Gene | Value |   | Gene | Value |   | Gene | Value |   | Gene | Value |   | Gene | Value |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| LpR2 | 3.5795 |   | Psa | 3.8529 |   | lola | 3.0121 |   | lola | 3.3019 |   | brat | 5.2812 |
| fs(1)h | 3.1376 |   | vnd | 3.6457 |   | CG31368 | 2.8063 |   | CG6919 | 2.9965 |   | ct | 4.828 |
| CG6954 | 2.7492 |   | ct | 3.201 |   | Kr-h1 | 2.7262 |   | CG31368 | 2.617 |   | CG31163 | 4.3345 |
| Psa | 2.7012 |   | fs(1)h | 3.1489 |   | svp | 2.7055 |   | CG5149 | 2.7675 |   | LpR2 | 3.6882 |
| zfh2 | 2.6247 |   | btd | 3.1229 |   | mub | 2.6475 |   | Kr-h1 | 2.7647 |   | vnd | 3.6866 |
| Fur1 | 2.4413 |   | zfh2 | 2.8421 |   | CG5149 | 2.5248 |   | TER94 | 2.6286 |   | zfh2 | 3.5314 |
| ct | 2.3804 |   | RhoBTB | 2.6022 |   | run | 2.4759 |   | tna | 2.5748 |   | pros | 3.4307 |
| S | 2.3674 |   | pros | 2.5679 |   | tna | 2.4302 |   | CG11153 | 2.4795 |   | Psa | 3.3998 |
| rux | 2.3574 |   | CG1124 | 2.5475 |   | CG6954 | 2.4235 |   | run | 2.3831 |   | fs(1)h | 3.3869 |
| RhoBTB | 2.26 |   | S | 2.5424 |   | CG11153 | 2.3045 |   | CG14888 | 2.0938 |   | CG31241 | 2.9973 |
| CG14889 | 2.1735 | **+** | oc | 2.5111 | **+** | Awd | 2.2295 | **+** | S | -2.0243 | **+** | HmgZ | 2.9226 |
| oc | 2.1421 |   | Fur1 | 2.43 |   | S | 2.1324 |   | rux | -2.0668 |   | Fur1 | 2.7469 |
| pros | 2.0882 |   | PHDP | 2.304 |   | CG14888 | 2.067 |   | oc | -2.3437 |   | oc | 2.6543 |
| Kr-h1 | -2.0447 |   | CG31241 | 2.2802 |   | Psa | -2.0276 |   | corto | -2.5556 |   | Toll-7 | 2.6161 |
| CG5149 | -2.1521 |   | nux | 2.2232 |   | nux | -2.093 |   | fs(1)h | -2.6211 |   | nux | 2.5975 |
| tna | -2.2102 |   | CG14889 | 2.1752 |   | fs(1)h | -2.141 |   | brat | -2.9904 |   | CG14889 | 2.3054 |
| CG14888 | -2.4346 |   | CG31163 | 2.1606 |   | CG1124 | -2.155 |   | ct | -3.3404 |   | S | 2.2324 |
| CG31368 | -2.4793 |   | HmgZ | 2.0795 |   | Fur1 | -2.1588 |   | zfh2 | -4.4947 |   | CG1124 | 2.0216 |
| Trim9 | -2.616 |   | svp | -2.0404 |   | S | -2.2539 |   | CG6954 | -4.7244 |   | Kr-h1 | -2.1439 |
| Awd | -3.0595 |   | TER94 | -2.1807 |   | corto | -2.2618 |   |   |   |   | tna | -2.1793 |
|   |   |   | corto | -2.3481 |   | oc | -2.3017 |   |   |   |   | CG5149 | -2.1892 |
|   |   |   | olf413 | -2.4404 |   | CG14889 | -2.4393 |   |   |   |   | run | -2.2194 |
|   |   |   | brat | -2.7256 |   | zfh2 | -2.5884 |   |   |   |   | Trim9 | -2.251 |
|   |   |   | CG31368 | -2.7293 |   | HmgZ | -3.6328 |   |   |   |   | olf413 | -2.3621 |
|   |   |   | mub | -2.9555 |   | btd | -3.7627 |   |   |   |   | btd | -3.0293 |
|   |   |   | Awd | -3.1413 |   | brat | -3.7716 |   |   |   |   | CG6919 | -3.3719 |
|   |   |   | lola | -3.8882 |   |   |   |   |   |   |   |   |   |

## Gene clustering
### Script walkthrough 1

- To make the big table we first need to find out all the genes present in at least one of the files
- Make sure not to use factors in read.delim()

```
# start with en empty collection of genes
genes <- c()
for( fileNum in 1:5 ){
    # load in files 13_DiffGenes1.tsv ...
    t <- read.delim(paste("13_DiffGenes", fileNum, ".tsv", sep=""),
                as.is=TRUE, header=FALSE)
    # label the input columns to help code readability
    names(t) <- c("gene", "expression")
    genes <- union(genes, t$gene)
}

# for tidiness order our genes by name
genes <- sort(genes)

genes # show all genes
```

> when loading in character data use **as.is=T** to prevent it being converted to factors!

> union() is a set operation, combines two vectors by eliminating duplicates. There are also intersect() and setdiff()

> Example code:
> 13_geneClustering.R

---

## Gene clustering
### Script walkthrough 2

- Using the complete list of genes, we can create the big table and fill in the values:

```
# make the destination table [rows = unique genes, cols = file numbers]
values <- matrix(0, nrow=length(genes), ncol=5)
rownames(values) <- genes # name the rows with the complete gene names

for(fileNum in 1:5){
    # read in the file again
    t <- read.delim(paste("13_DiffGenes", fileNum, ".tsv", sep=""),
            as.is=T, header=F)
    names(t) <- c("gene", "expression")

    # match the names of the genes to the rows in our big table
    index <- match(t$gene, rownames(values))
    # copy the expression levels
    values[index,fileNum] <- t$expression
}
```

> match() returns the index of first argument in the second, i.e. index of input file genes in the big table

> we use index to pick the rows in such way that they match the gene order in the input file

---

## Gene clustering
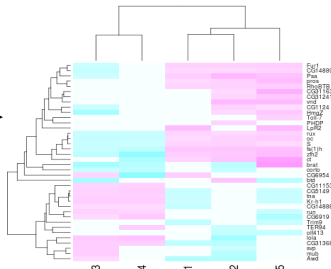### Script walkthrough 3

- Now we can do hierarchical clustering:

```
heatmap(values, scale="none", col = cm.colors(256))
```

> Values from the matrix are colour-coded. Rows and columns are re-arranged according to similarity

# Gene clustering
Script walkthrough 4

• In a second part of our analysis, we want to produce the same heatmap but only based on a list of experimentally verified genes

• The problem is data is not formatted in the most convenient way:

| genes | citation |
|---|---|
| oc,run,RhoBTB,CG5149,CG11153,S,Fur1 | Segal et al, Development 2001 |
| tna,Kr-h1,rux | Krejci et al, Development 2002 |

---

# Gene clustering
Script walkthrough 5

• We load in this table, and only extract the gene names, then we use them to select a subset of values matrix

```
# load in the tab-delimited file with genes and citations
t.exp <- read.delim("13_ExperimentalGenes.tsv", as.is=T)
# split all gene names by "," and then flatten it out into a single vector
experim.genes <- unlist( strsplit(t.exp$genes, ",") )
```

unlist() flattens out a nested list into a single vector

strsplit() splits a vector of strings by a custom split character (","), the results is a list of split values for each element of input vector

```
# redo the heatmap by using just the genes in the experimentally verified set
is.experimental <- rownames(values) %in% experim.genes
heatmap(values[ is.experimental, ], scale="none", col = cm.colors(256))
```

---

# Gene clustering review

• We load in the five tables twice - first to collect gene names, then to load expression values
• Based on expression table (values) we construct a clustered heatmap first on the whole set of genes, then on a selected subset

• Go through the code, try it out it and understand it

• Try answering the following questions:
  • what is rownames(values) ?
  • why is rownames(values)[index] and t$gene giving the same output?
  • what is a difference between rownames(values) %in% experim.genes and experim.genes %in% rownames(values)

Example code:
13_geneClustering.R