# Babelomics

# Microarray Data Analysis
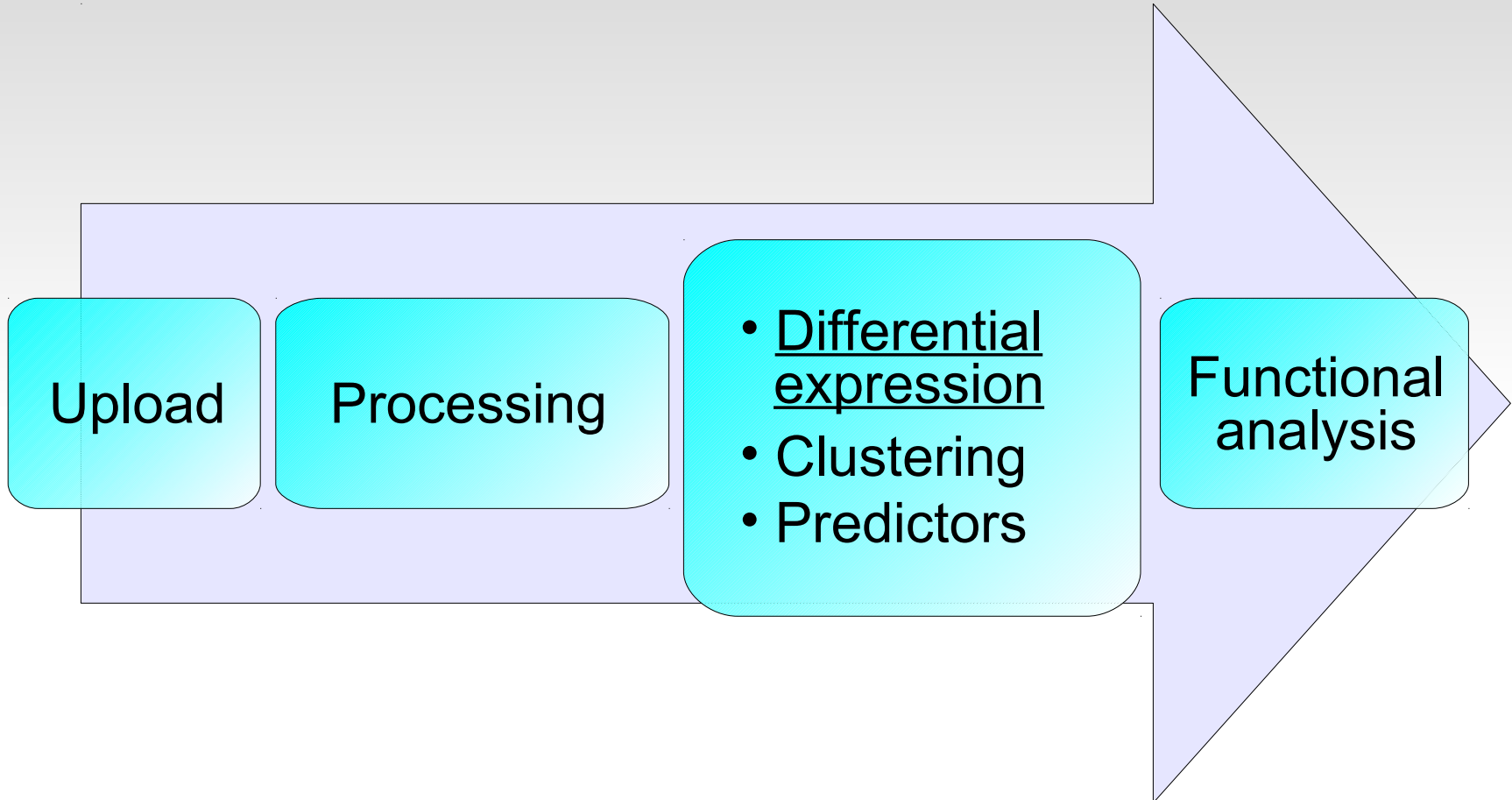# Differential Gene Expression

Valencia, March 2011

Francisco García,  fgarcia@cipf.es
Bioinformatics and Genomics Department
Centro de Investigacion Principe Felipe (CIPF)
(Valencia, Spain)

INB

PRINCIPE FELIPE
CENTRO DE INVESTIGACION

er cïberer
CENTRO DE INVESTIGACIÓN BIOMEDICA EN RED
DE ENFERMEDADES RARAS

# Differential Expression

Class comparison

Correlation

Survival

Time / dosage series

# Input



1. After normalization, we will add information relative to experimental design

2. Assigning values of variables to each array

# Input

Array names        Arrays        Tab separated file

genes

| #NAMES | col1 | col2 | col3 | col4 | col5 | col6 | col7 |
|---|---|---|---|---|---|---|---|
| YGR138C | -1.23 | -0.81 | 1.79 | 0.78 | -0.42 | -0.69 | 0.58 |
| YPR156C | -1.76 | -0.94 | 1.16 | 0.36 | 0.41 | -0.35 | 1.12 |
| YOR230W | -2.19 | 0.13 | 0.65 | -0.51 | 0.52 | 1.04 | 0.36 |
| YAL018C | -1.22 | -0.98 | 0.79 | -0.76 | -0.29 | 1.54 | 0.93 |
| YBR287W | -1.47 | -0.83 | 0.85 | 0.07 | -0.81 | 1.53 | 0.65 |
| YCL075W | -1.04 | -1.11 | 0.87 | -0.14 | -0.80 | 1.74 | 0.48 |
| YDR055w | -1.57 | -1.17 | 1.29 | 0.23 | -0.20 | 1.17 | 0.26 |
| YOR358W | -1.53 | -1.25 | 0.59 | -0.30 | 0.32 | 1.41 | 0.77 |
| YBR006W | -1.76 | -0.72 | 0.13 | -0.01 | -0.23 | 1.30 | 1.28 |
| YBR241C | -1.39 | -0.42 | -0.08 | -0.29 | -0.65 | 1.85 | 0.98 |
| YCR021c | -1.52 | -0.99 | 0.26 | 0.04 | -0.42 | 1.43 | 1.19 |
| YCR061W | -1.57 | -0.39 | 0.33 | -0.54 | -0.51 | 1.59 | 1.09 |
| YDL024c | -1.27 | -1.14 | 0.57 | -0.30 | -0.47 | 1.46 | 1.14 |
| YDR298C | -1.49 | -0.87 | 0.41 | -0.47 | -0.25 | 1.38 | 1.29 |
| YER141w | -1.69 | -0.60 | 0.00 | 0.41 | -0.62 | 1.45 | 1.05 |

......

# Results



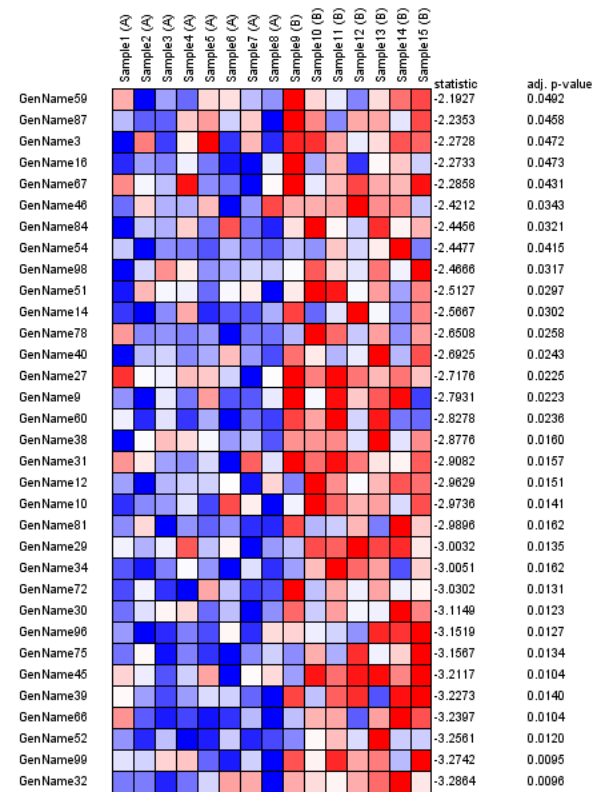| #NAMES | statistic | p-value adj. | p-value |
|---|---|---|---|
| 200067_x_at | 5.538204986516149 | 4.974604961551534E-6 | 2.4375564311602516E-4 |
| 200052_s_at | 5.2110706353314535 | 1.4519552681244469E-5 | 4.743053875873193E-4 |
| 200054_at | 5.102797363044661 | 4.263454480923057E-5 | 0.0010445463478261499 |
| 200009_at | 4.209329258084638 | 1.9598818661190837E-4 | 0.0027557324122247896 |
| 200017_at | 4.0805286865632855 | 2.2495774793860376E-4 | 0.0027557324122247896 |
| 1053_at | 3.9460740578057503 | 6.082189146003286E-4 | 0.005960545363083221 |
| 200013_at | 3.767033234598989 | 7.042746674112254E-4 | 0.006274447036936371 |
| 200071_at | 3.5180398564848283 | 0.0014872364080140634 | 0.012145763998781516 |
| 200076_s_at | 3.137574787036864 | 0.003912733450155826 | 0.0247030394179239 |
| 177_at | 3.0053355520231624 | 0.0061374669029413305 | 0.030073587824412523 |

| name | statistic | p-value | adj. p-value |
|---|---|---|---|
| 200067_x_at | 5.5382 | 0.0000049746 | 0.00024376 |
| 200052_s_at | 5.2111 | 0.00001452 | 0.00047431 |
| 200054_at | 5.1028 | 0.000042635 | 0.0010445 |
| 200009_at | 4.2093 | 0.00019599 | 0.0027557 |
| 200017_at | 4.0805 | 0.00022496 | 0.0027557 |
| 1053_at | 3.9461 | 0.00060822 | 0.0059605 |
| 200013_at | 3.767 | 0.00070427 | 0.0062744 |
| 200071_at | 3.518 | 0.0014872 | 0.012146 |
| 200076_s_at | 3.1376 | 0.0039127 | 0.024703 |
| 177_at | 3.0053 | 0.0061375 | 0.030074 |

prev  next  page 1 of 10

Search the term 200067_x_at

General databases

e! Ensembl
novo|seek

Other info

# Results

# Different experimental designs

# Class comparison



1. data

2. classes

3. methods

4. adj. p-value

5. job name

# Two-classes *form*



Methods:

Limma, t-test:

$$H_0: \mu_1 = \mu_2$$
$$H_a: \mu_1 \neq \mu_2$$

Fold-change:

$$\text{Log}_2 \left( \overline{y}_1 \, / \, \overline{y}_2 \right)$$

$$\overline{y}_1 - \overline{y}_2$$

# t – test for a gene expression

For each gene, we check if its mean expression is equal or different across the **two** classes

$H_0$: $\quad \mu_1 = \mu_2$

**Null** hypothesis: the mean expression is **equal** n both groups.

$H_a$: $\quad \mu_1 \neq \mu_2$

**Alternative** hypothesis: the mean expression is **different** between the groups.

Mean in group 1

Mean in group 2

Test Statistic: $T = \dfrac{\bar{Y}_1 - \bar{Y}_2}{\sqrt{s_1^2/N_1 + s_2^2/N_2}}$

Estimation of the variability of the differences

# Two-classes *results*

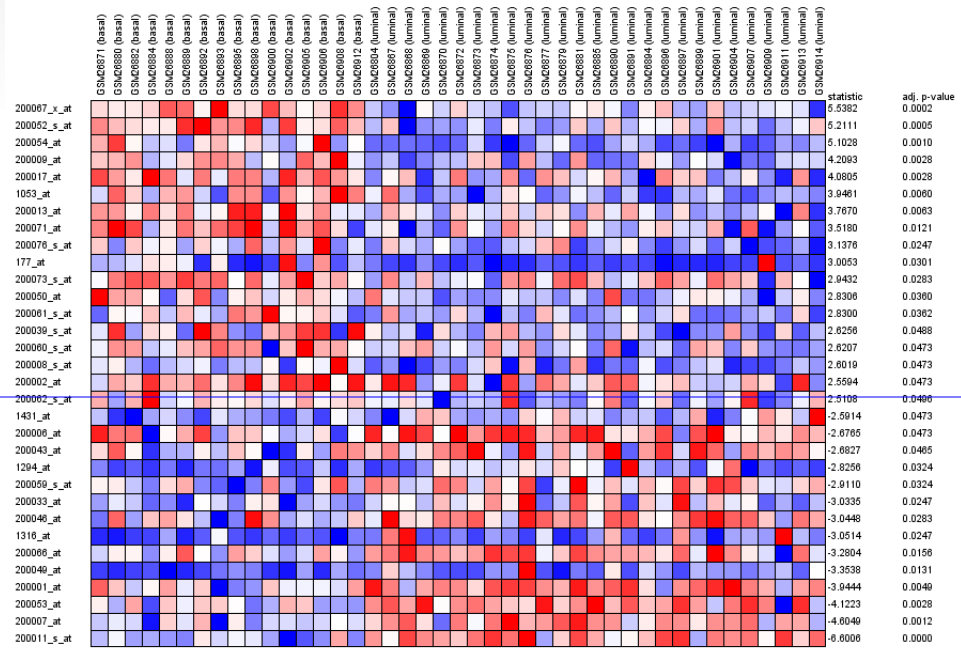Limma, t-test



```
#NAMES          statistic.            p-value.              adj. p-value
200067_x_at.    5.538204986516149.    4.974604961551534E-6. 2.4375564311602516E-4
200052_s_at.    5.211070635331453.    1.4519552681244469E-5. 4.743053875873193E-4
200054_at.      5.102797363044661.    4.263454480923057E-5. 0.001044546347826149
200009_at.      4.209329258084638.    1.9598818661190837E-4. 0.002755732412247896
200017_at.      4.0805286865632855.   2.2495774793860376E-4. 0.002755732412247896
1053_at.        3.9460740578057503.   6.082189146003286E-4. 0.005960545363083221
200013_at.      3.767033234598989.    7.042746674112254E-4. 0.006274447036936371
200071_at.      3.5180398564848283.   0.0014872364080140634. 0.012145763998781516
200076_s_at.    3.137574787036864.    0.003912733450155826. 0.024703039417792398
177_at.         3.0053355520231624.   0.0061374669029413305. 0.030073587824412523
200073_s_at.    2.9431619702299616.   0.005421086687530431. 0.028266292930286572
```

Fold-change

```
#NAMES          log                     diff
AFFX-BioB-5_at  0.023216278176490163    0.07865499999999948
AFFX-BioB-M_at  -0.01743916189488063    -0.06712000000000007
AFFX-BioB-3_at  0.011365357483625202    0.03912499999999941
AFFX-BioC-5_at  -0.014803025848131719   -0.06477500000000092
AFFX-BioC-3_at  0.012163222943743255    0.05631999999999948
AFFX-BioDn-5_at -0.0319905025238944     -0.1625899999999998
AFFX-BioDn-3_at -0.04557267653715912    -0.2699649999999991
AFFX-CreX-5_at  0.005800532734088386    0.042909999999999116
AFFX-CreX-3_at  0.005947534544836981    0.04539000000000115
```

# Two-classes *results*

# One-class



$$H_o: \quad \mu = 0$$
$$H_a: \quad \mu \neq 0$$

Genes ordered by statistic

# Multi-classes



$$H_0: \quad \mu_1 = \mu_2 = \dots = \mu_n$$
$$H_a: \quad \text{not } H_0$$

# Gene expression related to a continuous variable, *form*



Select your data

browse server    no data selected.
Or go to Upload Data form:  Upload [datamatrix]

Select the class to analyse

Class name:  No classes available

Select test

- Pearson's correlation
- Spearman correlation
- Regression

Select multiple-test correction

- Benjamini and Hochberg (BH), FDR
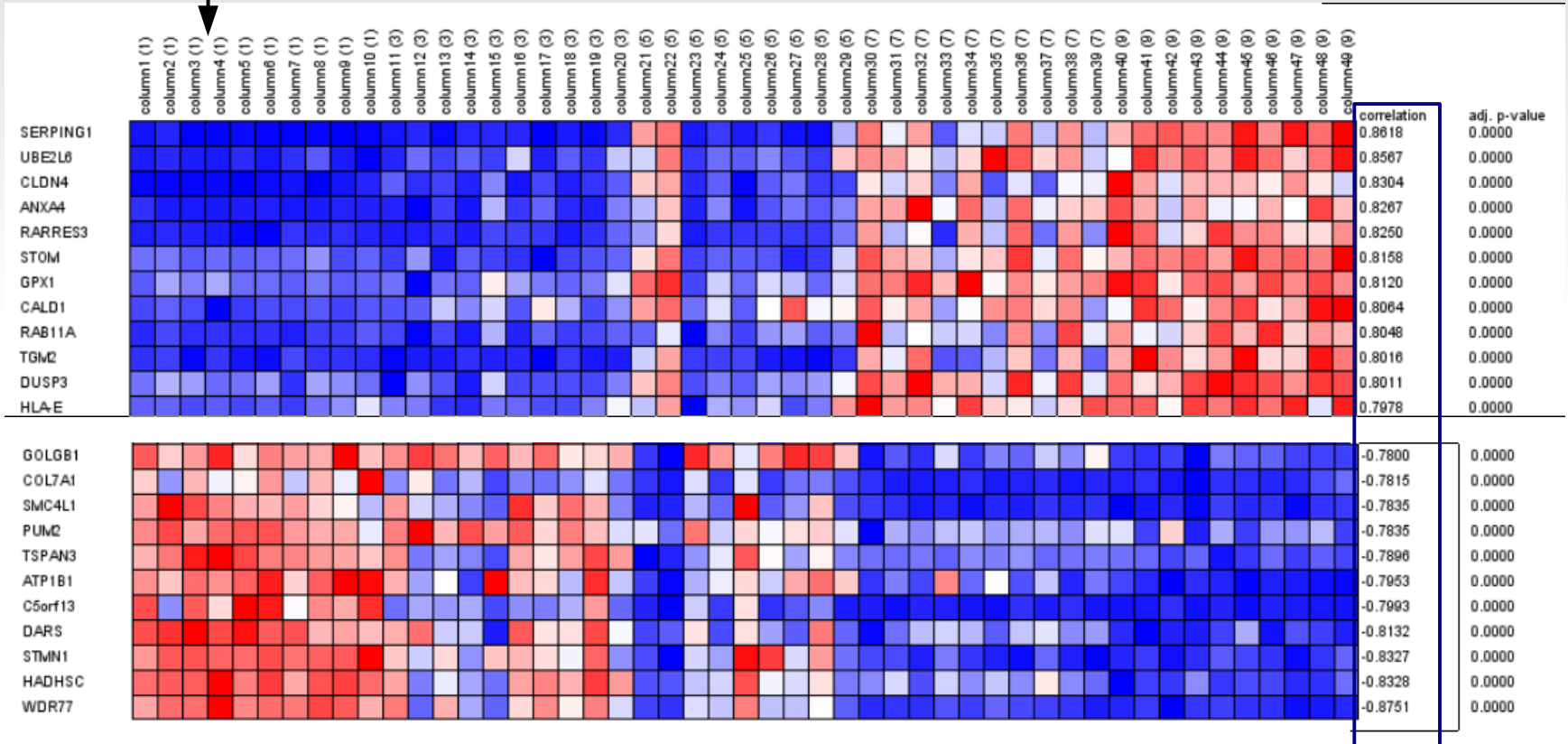- Benjamini and Yekutieli (BY)
- Hochberg
- Holm
- Bonferroni

methods

$$\rho = \frac{\sum_i (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_i (x_i - \bar{x})^2 \sum_i (y_i - \bar{y})^2}}.$$
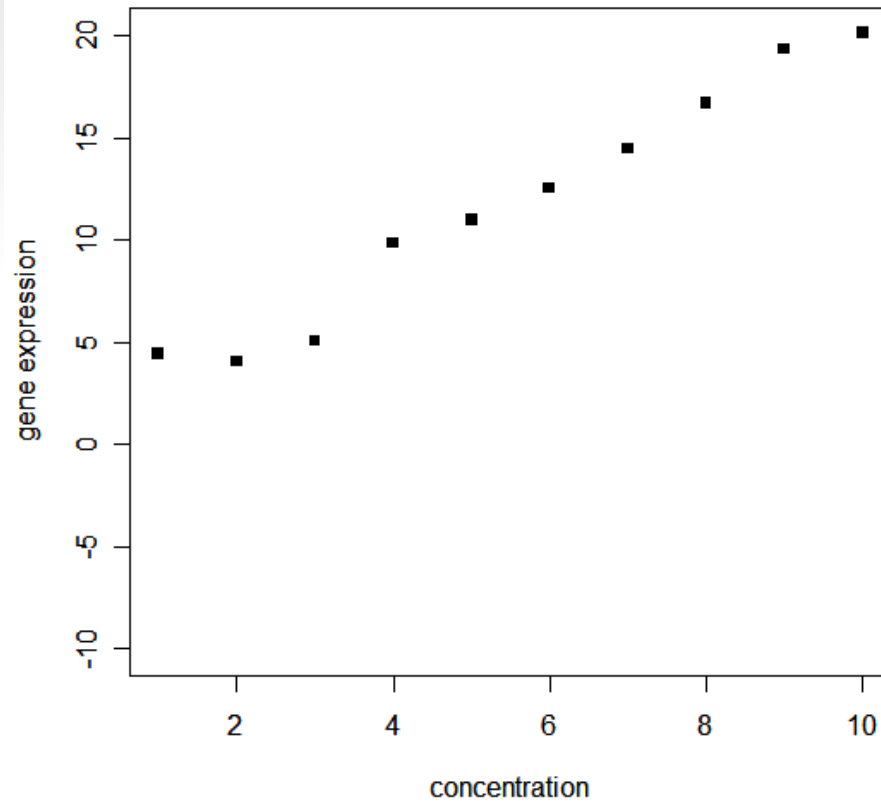
$$\rho = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}.$$

# Correlation *results*
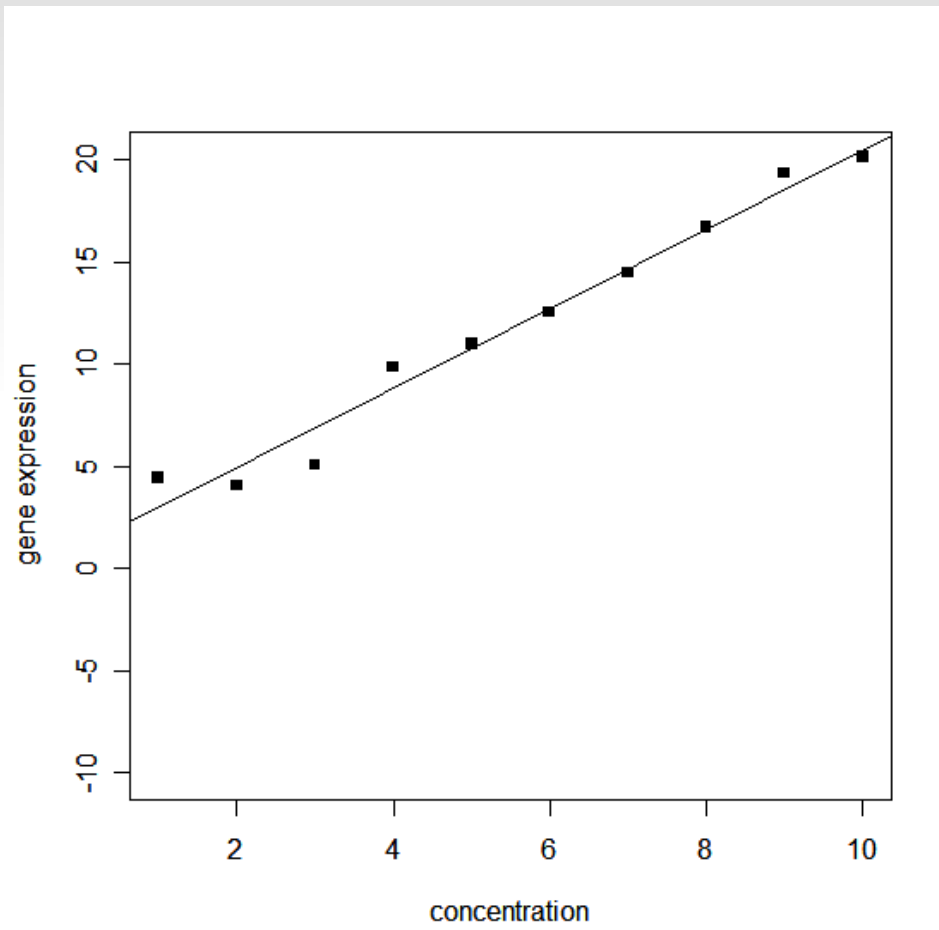
Arrays ranked according to the independent variable



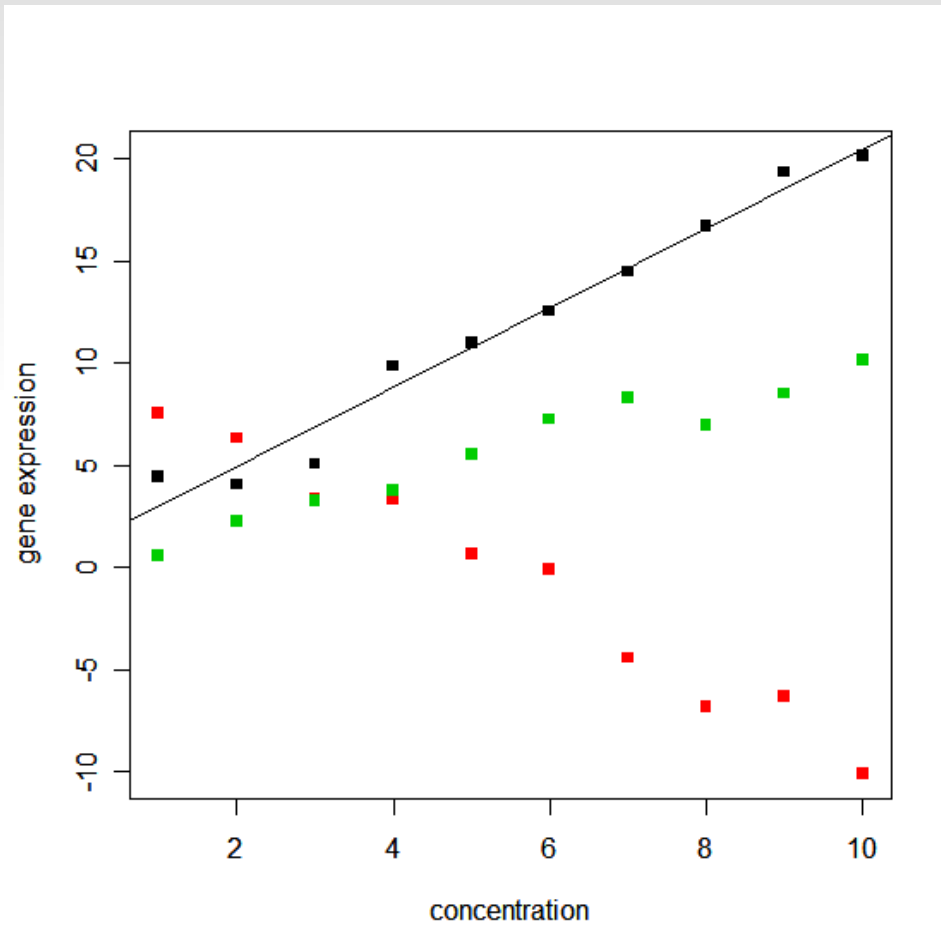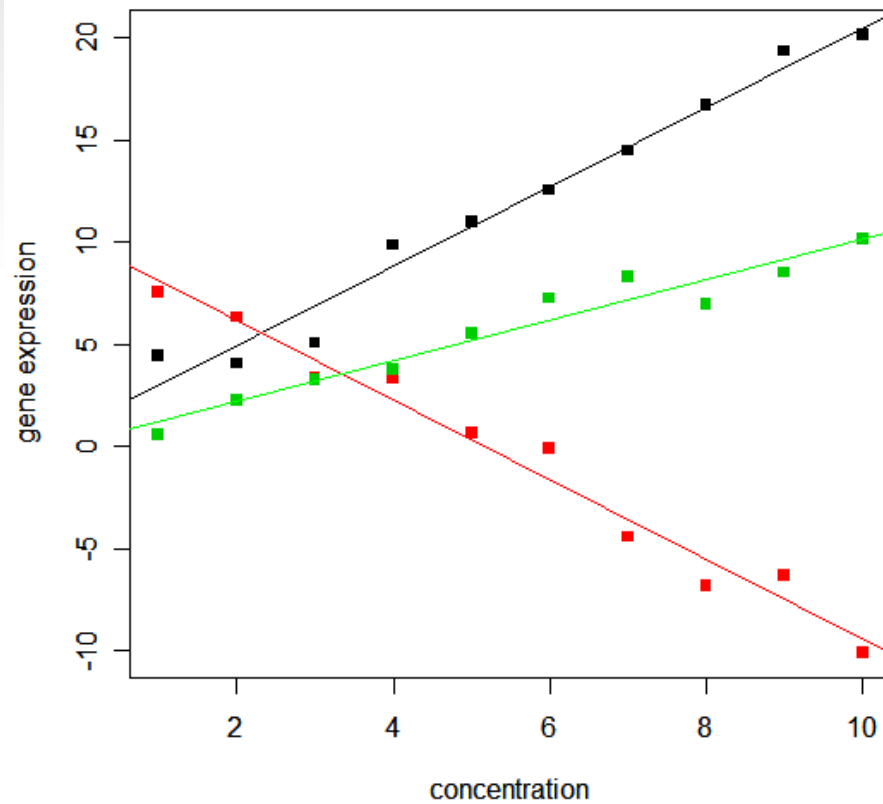Genes ranked by correlation to the continuous variable

# Regression

# Regression



gene1 slope

# Regression



gene1 slope

# Regression



gene1 slope
gene2 slope
gene3 slope
…

# Survival *form*

### Select your data

browse server   no data selected.

Or go to Upload Data form:   Upload [datamatrix]

←   1. Expression

### Select time and series variables

Time variable name:   No classes available ⌄

Censored variable name:   No classes available ⌄

←   2. Survival times & censoring variable

### Select test

● Cox

←   3. Cox proportional hazards regression model

### Select multiple-test correction

● Benjamini and Hochberg (BH), FDR

○ Benjamini and Yekutieli (BY)

○ Hochberg

○ Holm

○ Bonferroni

←   4. Multiple-test correction

### Select adjusted p-value

Adj. p-value (0.0-1.0)   0.05

←   5. adjusted p-value

# Survival *results*



Arrays ranked according to the survival time

Genes ranked by their relationship with survival time

- Cox model coefficients
- Estimate for the statistics
- p-values

# Time course analysis / Dose analysis, *form*

**Select your data**

browse server    no data selected.
Or go to Upload Data form:  Upload [datamatrix]

← Expression data

**Series variables**

Name of Continuous Variable:  No classes available ∨

Name of Variable defining Series:  No classes available ∨

← Time variable and series classification

**Options**

Polynomial degree  1 ∨

← Complexity of model

Significance level for gene selection (0.0-1.0)  0.05

← Significance level

Multiple testing adjustment
- ⦿ Benjamini and Hochberg (BH), FDR
- ◯ Benjamini and Yekutieli (BY)
- ◯ Hochberg
- ◯ Holm
- ◯ Bonferroni
- ◯ Hommel

← Multiple testing

Significance Level for model variable(0.0-0.9)  0.05

Cluster method
- ⦿ Hierarchical clustering
- ◯ K-Means

Number of clusters (k-value)  9

← Clustering

# Time course analysis / Dose analysis, *example*

- ✓ Arabidopsis.
- ✓ 4 series: **control and 3 treatments**
  (cold, salt, heat).
- ✓ 3 time points.
- ✓ 3 replicates.

## What do we want?

# Time course analysis / Dose analysis, *results*

## Control

```
#genes   cluster
STMEQ29  1
STMID05  2
STMGB57  1
STMEY09  3
STMHY68  4
STMGI03  5
STMCU02  1
STMGB35  6
STMDI90  1
STMJI76  7
STMJO83  8
STMCS44  1
STMIA31  3
STMJF53  1
STMIQ37  4
STMJC14  9
STMCM86  2
STMGQ83  3
STMCK87  1
STMCU87  1
STMHN19  9
STMED61  3
STMIC27  5
STMCH79  6
STMDU84  8
STMIO93  9
STMEG09  1
STMIX47  9
STMIP63  6
STMEV77  4
```

## Cold vs Control

```
#genes   cluster
STMJH42  1
STMDE66  2
STMHZ45  1
STMGL58  3
STMIF71  1
STMEG62  4
STMFB37  5
STMEQ29  1
STMDW06  6
STMEL85  5
STMEG74  1
STMCO26  7
STMHX33  3
STMDV94  3
STMID12  1
STMCV66  2
STMGH56  2
STMEJ16  5
STMCD46  1
STMIT95  1
```

## Salt vs Control

```
#genes   cluster
STMJH42  1
STMDE66  2
STMHZ45  1
STMGL58  3
STMIF71  1
STMEG62  4
STMEQ29  1
STMDW06  5
STMEG74  1
STMCO26  6
STMID05  5
STMID12  1
STMDH27  7
STMEJ16  1
STMCD46  1
STMIT95  3
STMHJ39  1
STMGB57  3
STMIT31  1
STMEZ42  1
STMIM44  5
STMHN16  5
STMEY09  2
STMCE01  1
STMIY82  2
STMEU24  8
STMHH10  1
STMGQ20  5
STMGI03  9
STMCY10  5
STMHV34  5
STMHY91  2
STMJN05  1
STMEF65  1
```
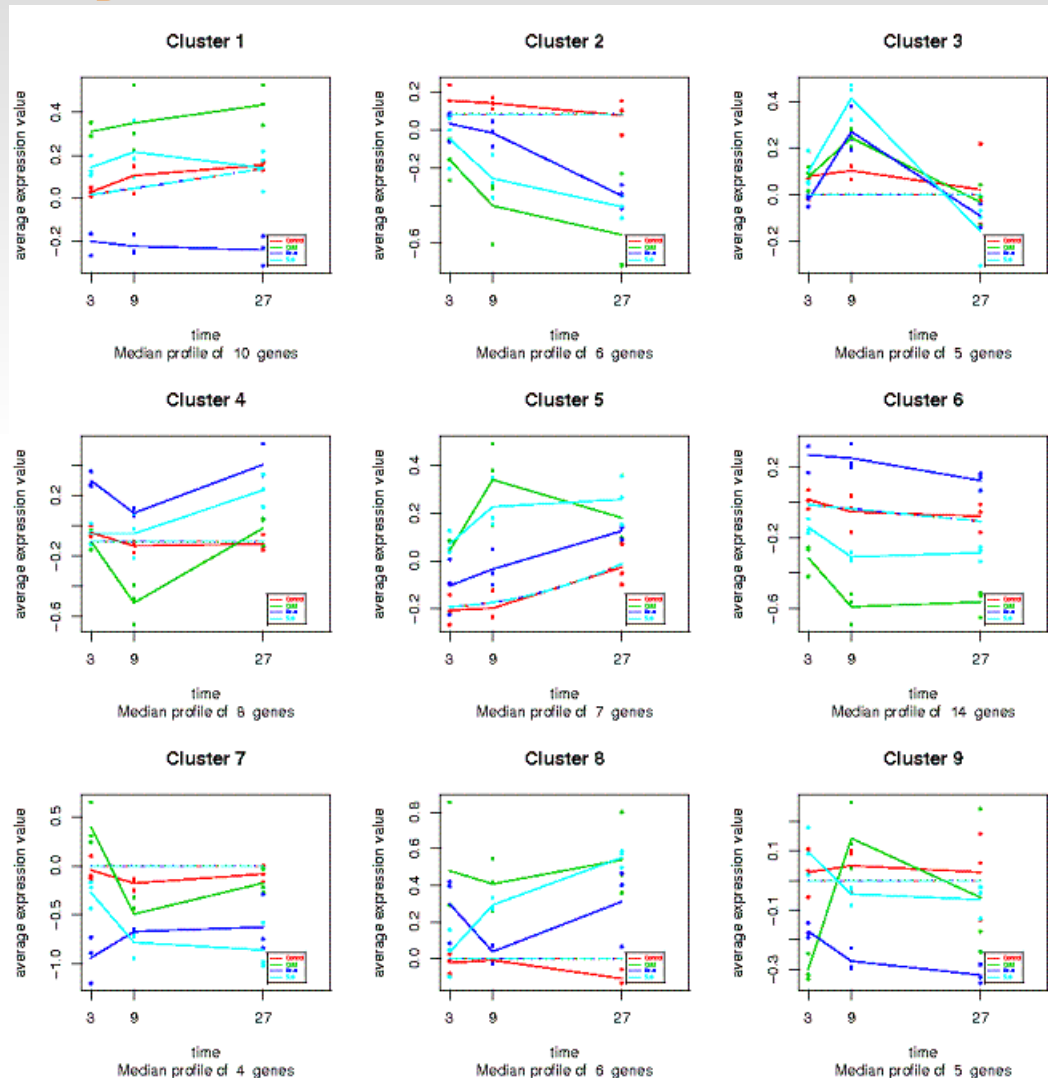
## Heat vs Control

```
#genes   cluster
STMJH42  1
STMDE66  2
STMHZ45  1
STMGL58  3
STMIF71  1
STMEG62  4
STMFB37  5
STMEQ29  1
STMDW06  6
STMEL85  5
STMEG74  1
STMCO26  7
STMHX33  3
STMDV94  3
STMID12  1
STMCV66  2
STMGH56  2
STMEJ16  5
STMCD46  1
STMIT95  1
STMJE19  2
STMHJ39  1
STMGU26  5
```

# Time course analysis / Dose analysis, *results*

## Cold vs Control

| #genes | cluster |
|--------|---------|
| STMFB37 | 1 |
| STMDW06 | 2 |
| STMEL85 | 1 |
| STMCO26 | 3 |
| STMHX33 | 1 |
| STMHQ28 | 1 |
| STMCV66 | 4 |
| STMDH27 | 4 |
| STMIT95 | 5 |
| STMJE19 | 6 |
| STMHJ39 | 5 |
| STMGU26 | 1 |
| STMEZ42 | 5 |
| STMCV36 | 4 |
| STMIM44 | 6 |
| STMEM39 | 5 |
| STMHY68 | 4 |
| STMGH85 | 1 |
| STMGQ20 | 2 |
| STMGI03 | 7 |
| STMJN05 | 8 |
| STMDB75 | 7 |
| STMFB38 | 4 |
| STMDJ03 | 9 |
| STMGB35 | 5 |
| STMDU19 | 9 |
| STMDE59 | 4 |
| STMCF08 | 6 |
| STMHK44 | 7 |
| STMJI76 | 3 |
| STMEM80 | 6 |
| STMIA39 | 6 |
| STMIO60 | 8 |

210 significant genes

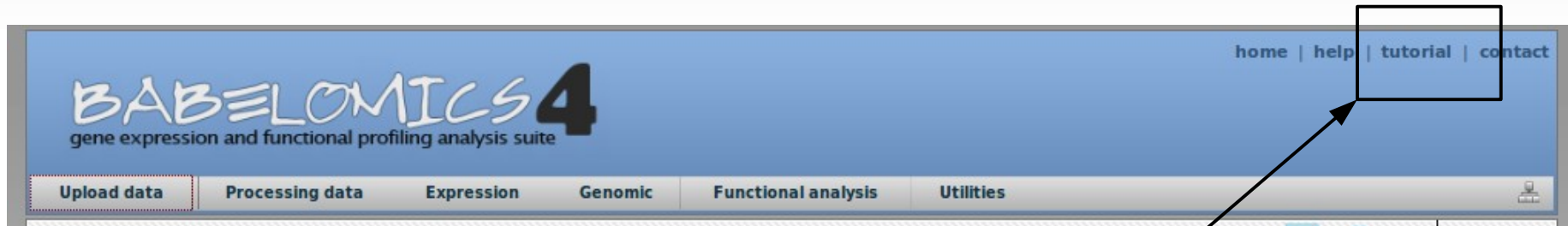# Time course analysis / Dose analysis. Redirecting to functional tools

# Remember...

➢ Babelomics allows us to analyze Differential Expression in **different experimental scenarios.**

➢ Differential Expression needs **normalized data** from Normalization Babelomics or other tool.

➢ These results can be **functional interpretated** using several tools in Babelomics: FatiGo, Logistic Models, Snow,...

# Let's practise!

## http://babelomics.bioinfo.cipf.es/



# Go to the tutorial:
**Expression Data Analysis / Differential Expression**