



Generating quality metrics reports for microarray data sets

Audrey Kauffmann

Introduction

- Microarrays are widely/routinely used
- Technology and protocol improvements → trustworthy
- Variance and noise
 - Technical causes:
 - Platform
 - Lab, experimentalist
 - RNA extraction
 - Amplification, labeling, hybridization, scanning...
 - Biological causes:
 - Tissue itself (cell lines, biopsies, blood...)
 - Tissue contamination
 - Clinical covariates (age, sex, race...)
 - Cell cycle...

Who is concerned?

- Experimentalist investigating a set of samples
 - Choose between different technology platforms, expt. Protocols
 - Decide when to repeat (certain parts of) the experiment
- Statistical collaborator analysing the experiment
 - Decide whether to proceed or to ask the experimentalist to go back to the lab
- Microarray core facility
 - Decide whether to consider their product fit for delivery to customer
 - Customer decides whether to be content (pay the bill)
- Integrative biologist analysing data in a public database
 - Has to choose which experiments
 - Which arrays within an experiment to consider
- Public data(base) provider
 - Put a quality score on each of her offerings

At which step of the analysis?

- Importing the data
- Preprocessing: background correction, normalisation, summarization of probesets
- Differential Expression
- Gene set enrichment analysis

At which step of the analysis?

- Importing the data
- Quality Assessment
- Preprocessing: background correction, normalisation, summarization of probesets
- Quality Assessment
- Differential Expression
- Gene set enrichment analysis

At which step of the analysis?

- Importing the data

- Quality Assessment

- Preprocessing: background correction, normalisation, summarization of probesets

- Quality Assessment

Remove outlier(s)

- Differential Expression
- Gene set enrichment analysis

What aspects to be evaluated? Which quality metrics?

Per Slide

- What are we looking at?
 - Intensity-dependent ratio
 - Detection of spatial effects
- How?
 - MAplots
 - Representation of the chip

Between Slides

- What are we looking at?
 - Homogeneity
 - Outlier samples
 - Biological meaning
- How?
 - Boxplots, density plots
 - Heatmap, PCA

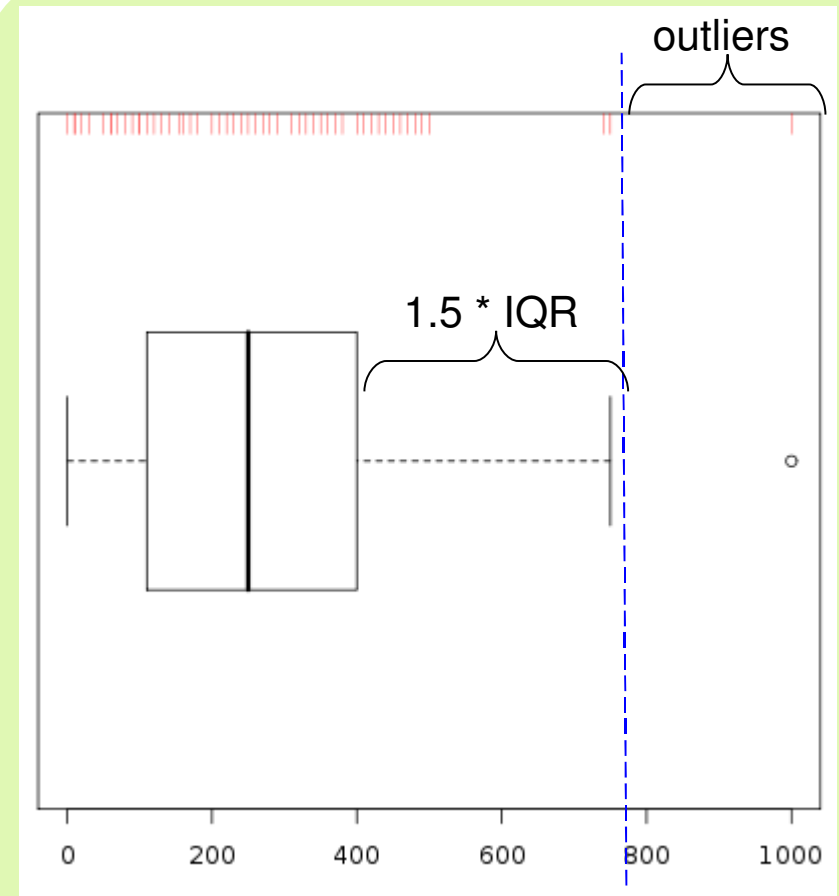
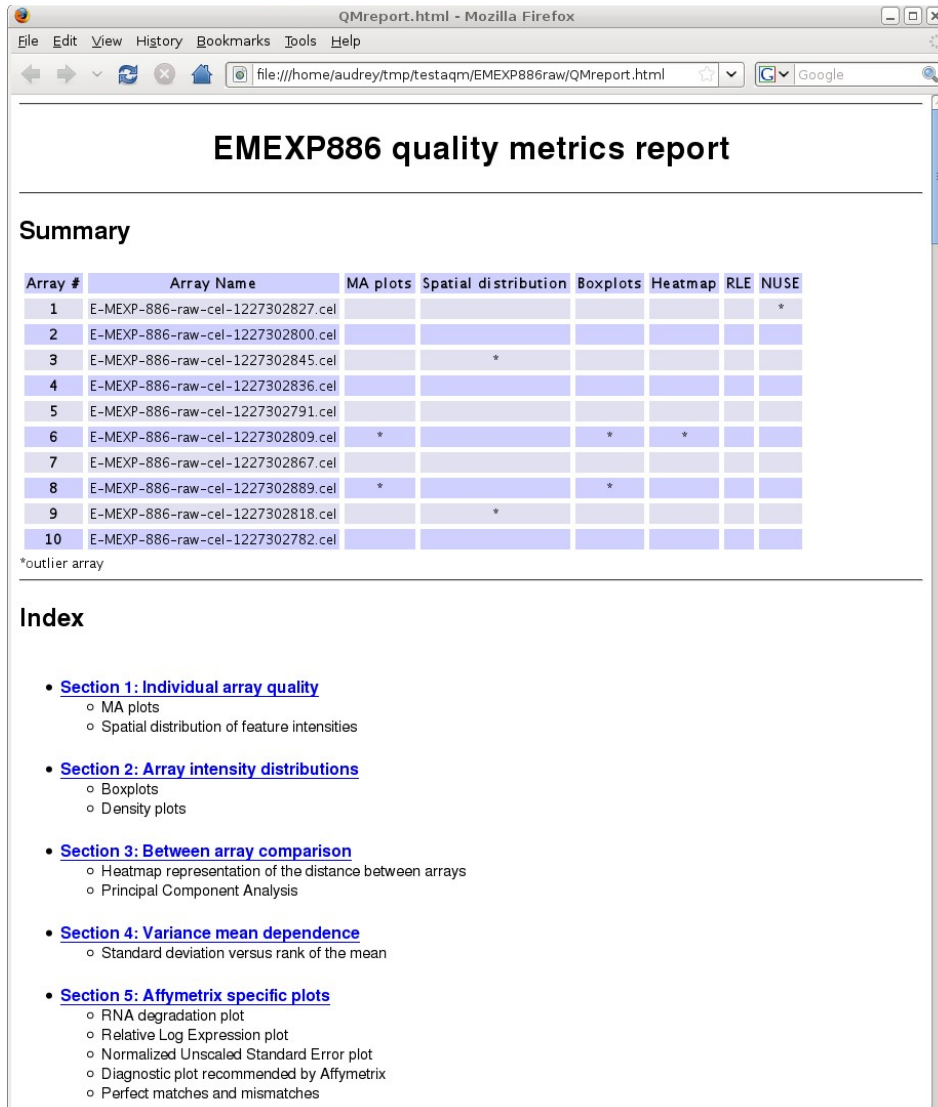
How to easily perform quality assessment?

- **arrayQualityMetrics** - Bioconductor package for Affymetrix, Agilent, Illumina, homemade arrays ...
- From an R object \Rightarrow HTML report
- Plots:
 - MA plot and spatial representations
 - Boxplots and density
 - Heatmap and PCA
 - Variance-mean dependency
 - GC content and probe mapping studies
 - Affymetrix only: NUSE, RLE, RNA degradation, QCstats, PM/MM
- Outlier identification

Hands-on

- Run reports

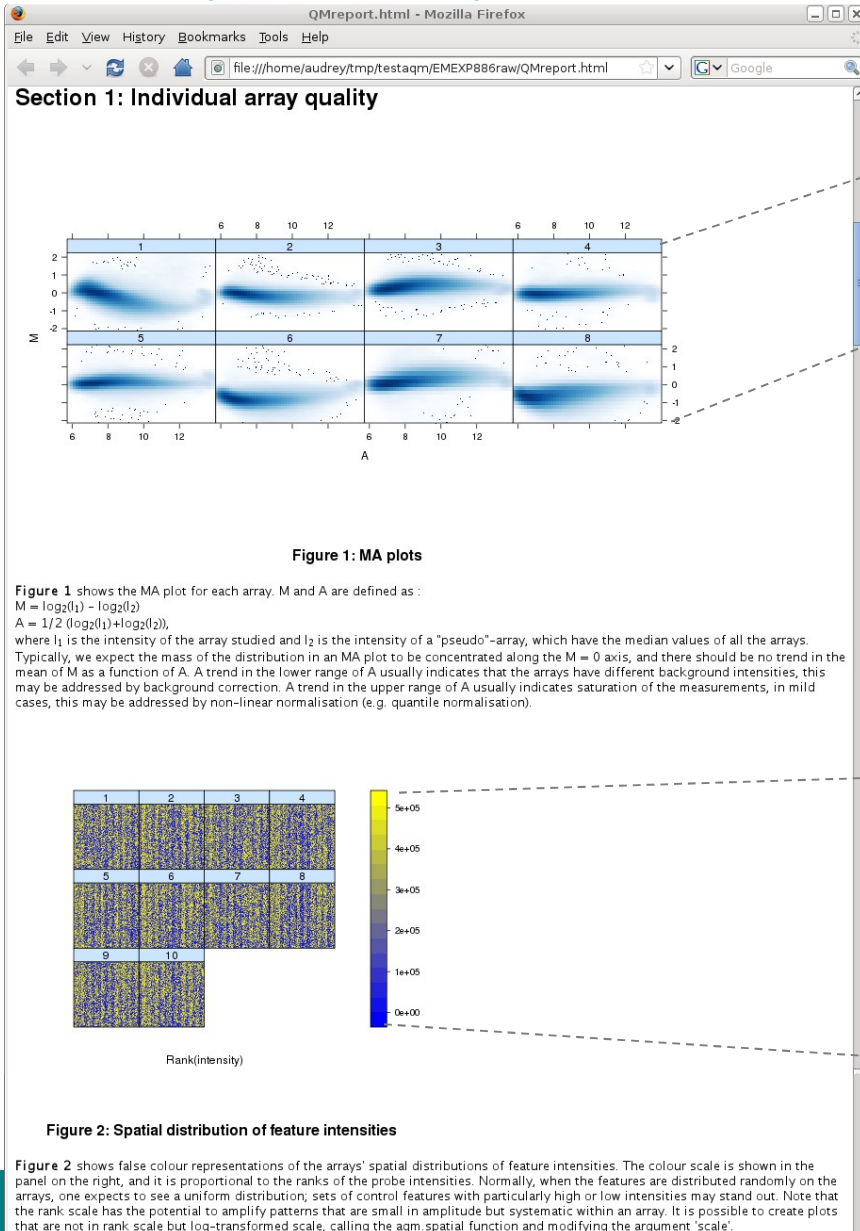
arrayQualityMetrics report – outlier detection



Boxplot of the scores

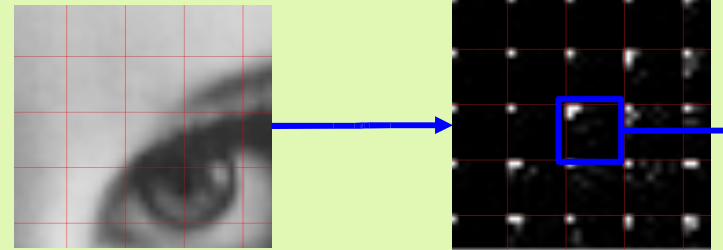


arrayQualityMetrics report – per array



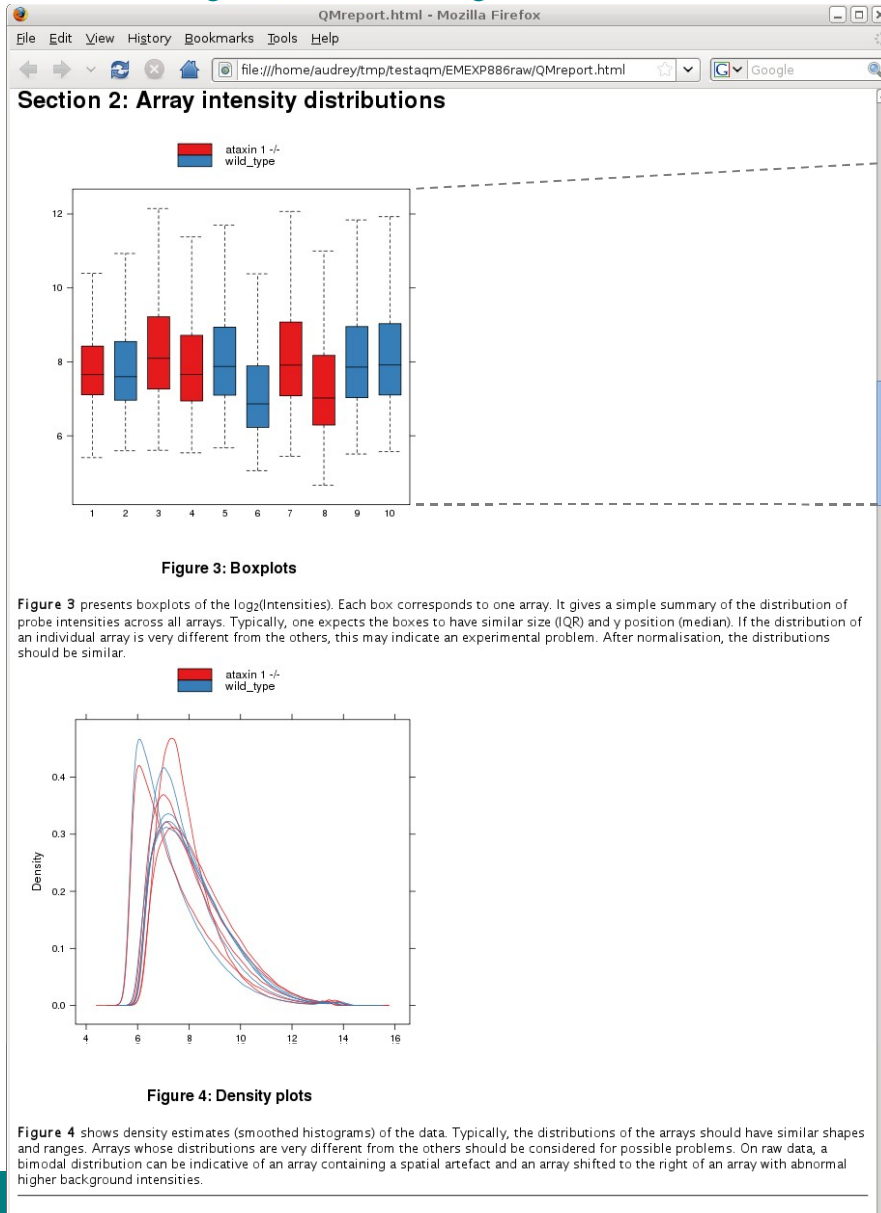
$$S_i = |_{ik}$$

Fourier Transform



$$S_i = \frac{\sum lowFreq_{ik}}{\sum highFreq_{ik}}$$

arrayQualityMetrics report – intensity distributions



$$S_i = \text{median}(I_{ik})$$
$$S_i = \text{IQR}(I_{ik})$$

arrayQualityMetrics report – *Between arrays*

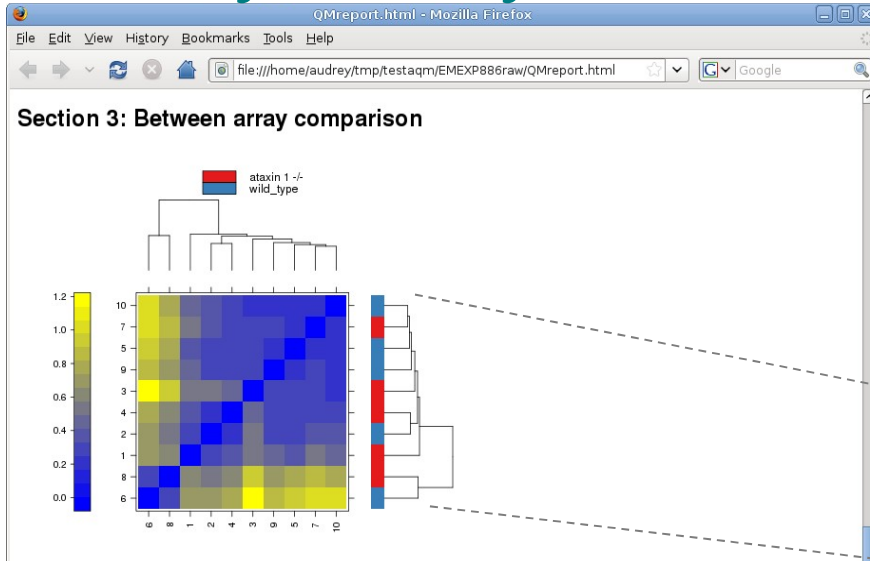


Figure 5: Heatmap representation of the distance between arrays

Figure 5 shows a false colour heatmap of between arrays distances, computed as the mean absolute difference (L₁-distance) of the vector of M-values for each pair of arrays on every probes without any filtering. The colour scale is chosen to cover the range of L₁-distances encountered in the dataset. Arrays for which the sum of the distances to the others is much different from the others, are detected as outlier arrays. The dendrogram on this plot also can serve to check if, the arrays cluster accordingly to a biological meaning.

$d_{xy} = \text{mean} |M_{xi} - M_{yj}|$
 Here, M_{xi} is the M-value of the i -th probe on the x -th array, without preprocessing. Consider the following decomposition of M_{xi} : $M_{xi} = z_i + \beta_{xi} + \epsilon_{xi}$ where z_i is the probe effect for probe i (the same across all arrays), ϵ_{xi} are i.i.d. random variables with mean zero and β_{xi} is such that for any array x , the majority of values β_{xi} are negligibly small (i. e. close to zero). β_{xi} represents differential expression effects. In this model, all values d_{xy} are (in expectation) the same, namely 2 times the standard deviation of ϵ_{xi} .

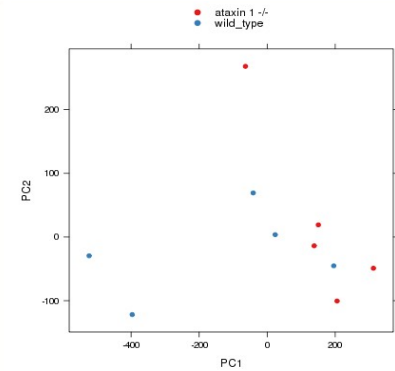


Figure 6: Principal Component Analysis

Figure 6 represents a biplot for the first two principal components from the dataset. The colours correspond to the group of interest given. We expect the arrays to cluster accordingly to a relevant experimental factor. The principal components transformation of a data matrix re-expresses the features using linear combination of the original variables. The first principal component is the linear combination chosen to possess maximal variance, the second is the linear combination orthogonal to the first possessing maximal variance among all orthogonal combination.

For each couple of arrays i and j , k is a probe and the distance between the arrays is:

$$d_{ij} = \text{mean}_k (|I_{ik} - I_{jk}|)$$

$$S_i = \sum d_{ij}$$

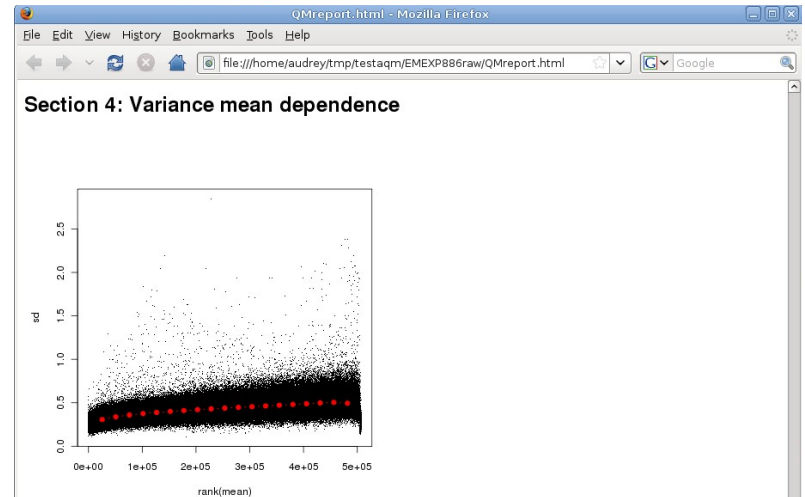
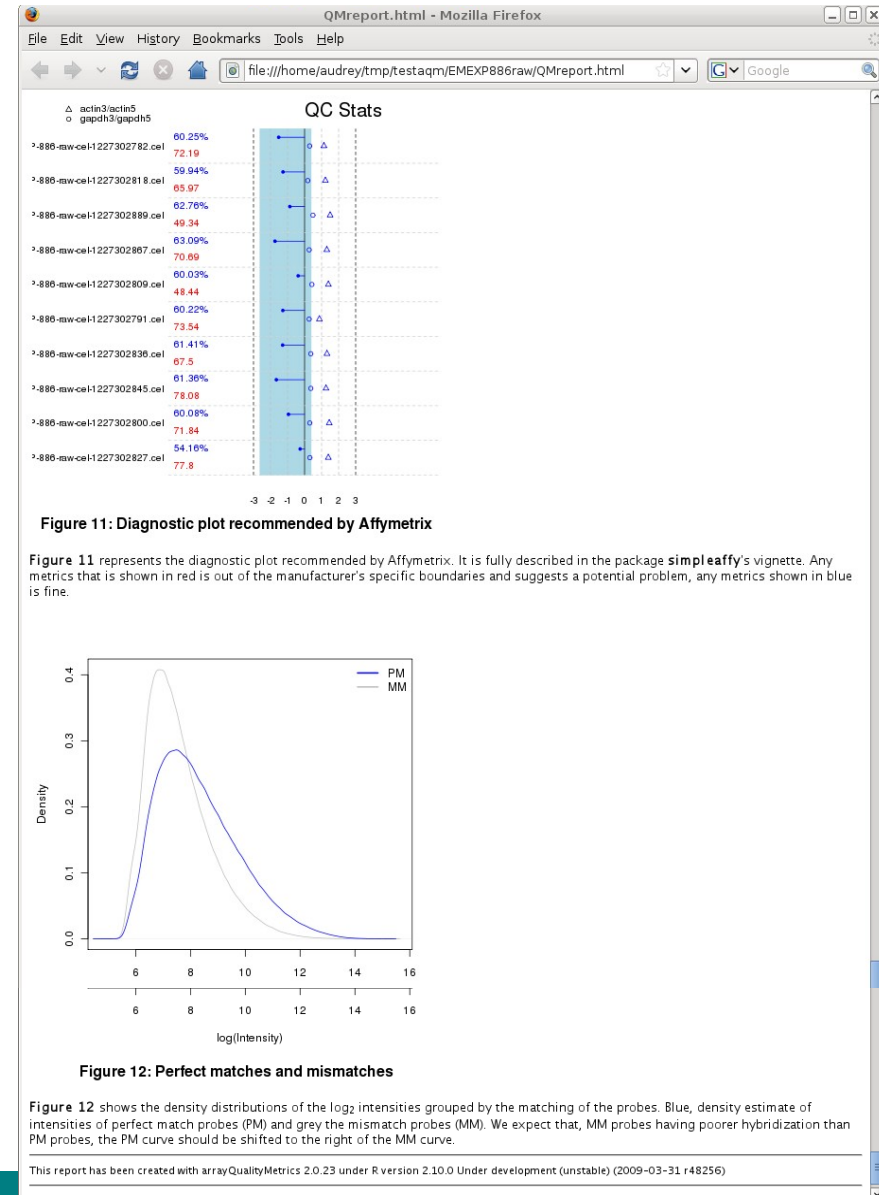
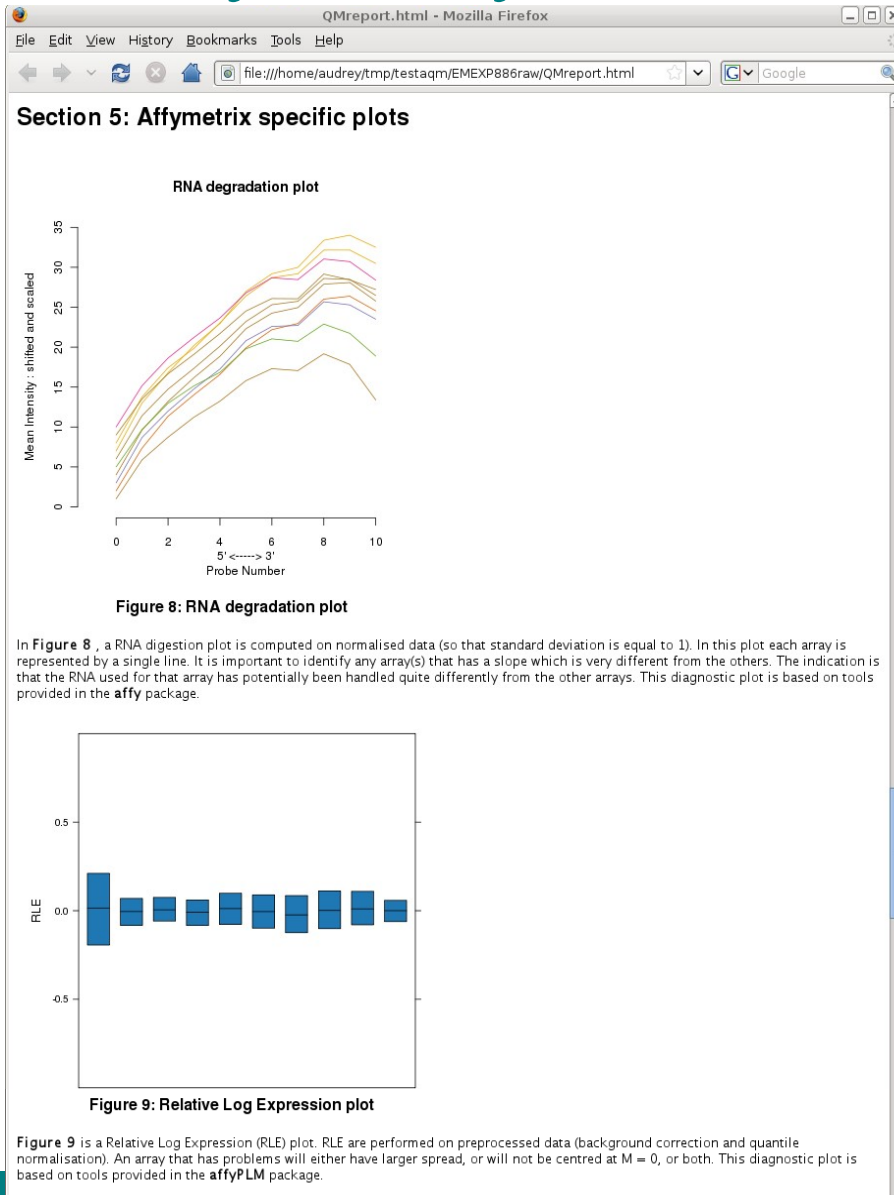


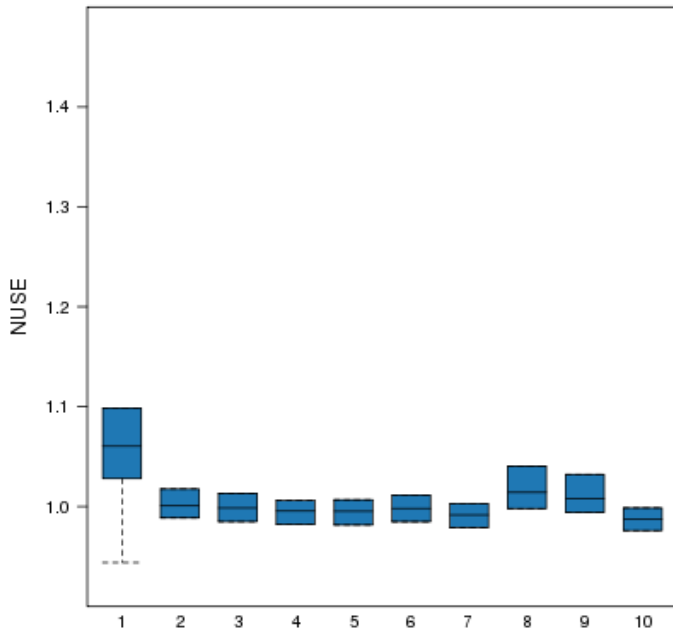
Figure 7: Standard deviation versus rank of the mean

For each feature, Figure 7 shows the standard deviation of the intensities across arrays on the y-axis versus the rank of their mean on the x-axis. The red dots, connected by lines, show the running median of the standard deviation. After normalisation and transformation to a logarithm(like) scale, one typically expects the red line to be approximately horizontal, that is, show no substantial trend. In some cases, a hump on the right hand of the x-axis can be observed and is symptomatic of a saturation of the intensities.

arrayQualityMetrics report – Affymetrix plots



arrayQualityMetrics report – Affymetrix NUSE: Normalised Unscaled Standard Error



$$S_i = \text{median}(I_{ik})$$
$$S_i = \text{IQR}(I_{ik})$$

$$NUSE(\beta_{kj}) = \frac{SE(\beta_{kj})}{\text{med}_j(SE(\beta_{kj}))}$$

- Fitting a probe level model (gene k , array j)
- Differences in variability between genes. An array with elevated SE (standard error) relative to the other arrays is of lower quality

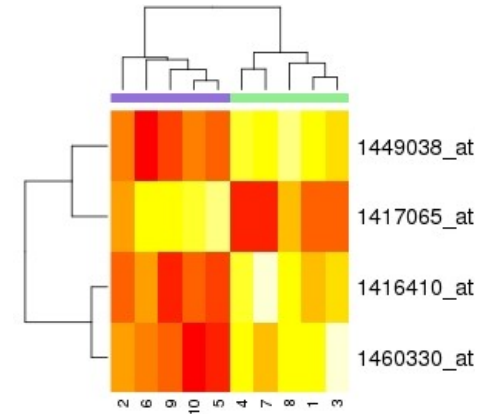
Why is outlier detection important? Example

- ArrayExpress experiment E-MEXP-886, cerebellar gene expression:
 - 5 WT mice (15 weeks of age)
 - 5 Atnx1 KO mice (15 weeks of age)
- Affymetrix MOE-430A (mouse) Genechip
- Ataxin 1 (Atxn1): protein of unknown function associated with cerebellar neurodegeneration in Spinocerebellar Ataxia type 1 (SCA1), which impairs the of eye movement

Results from E-MEXP-886 analysis

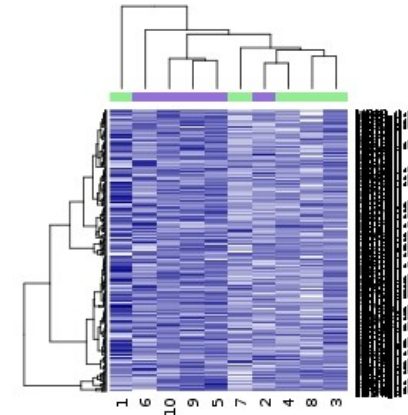
- Moderated t-test

	# Genes	
	P < 0.01	P < 0.001
10 samples	34	4



- Most enriched KEGG Pathway

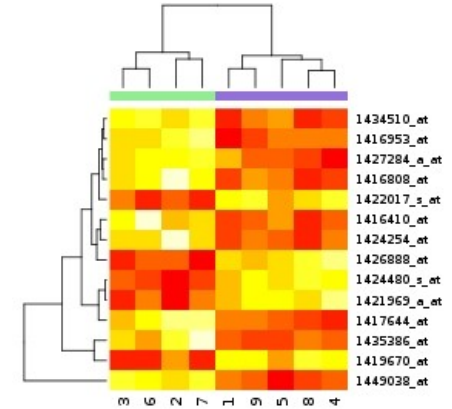
	Neuroactive ligand-receptor interaction	
	# Significant genes	Corrected t-value
10 samples	4	-5.65



Outlier array's impact on results

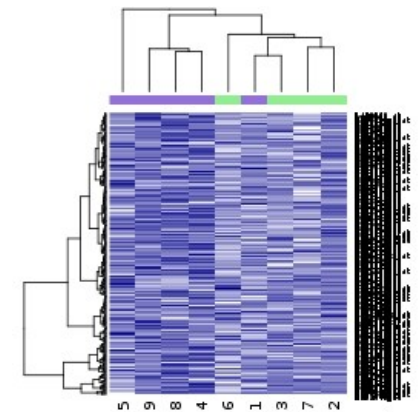
- Moderated t-test (limma)

	# Genes	
	P < 0.01	P < 0.001
10 samples	34	4
Without array 1	190	14



- Most enriched KEGG Pathway

	Neuroactive ligand-receptor interaction	
	# Significant genes	Corrected t-value
10 samples	4	-5.65
Without array 1	23	-11.53



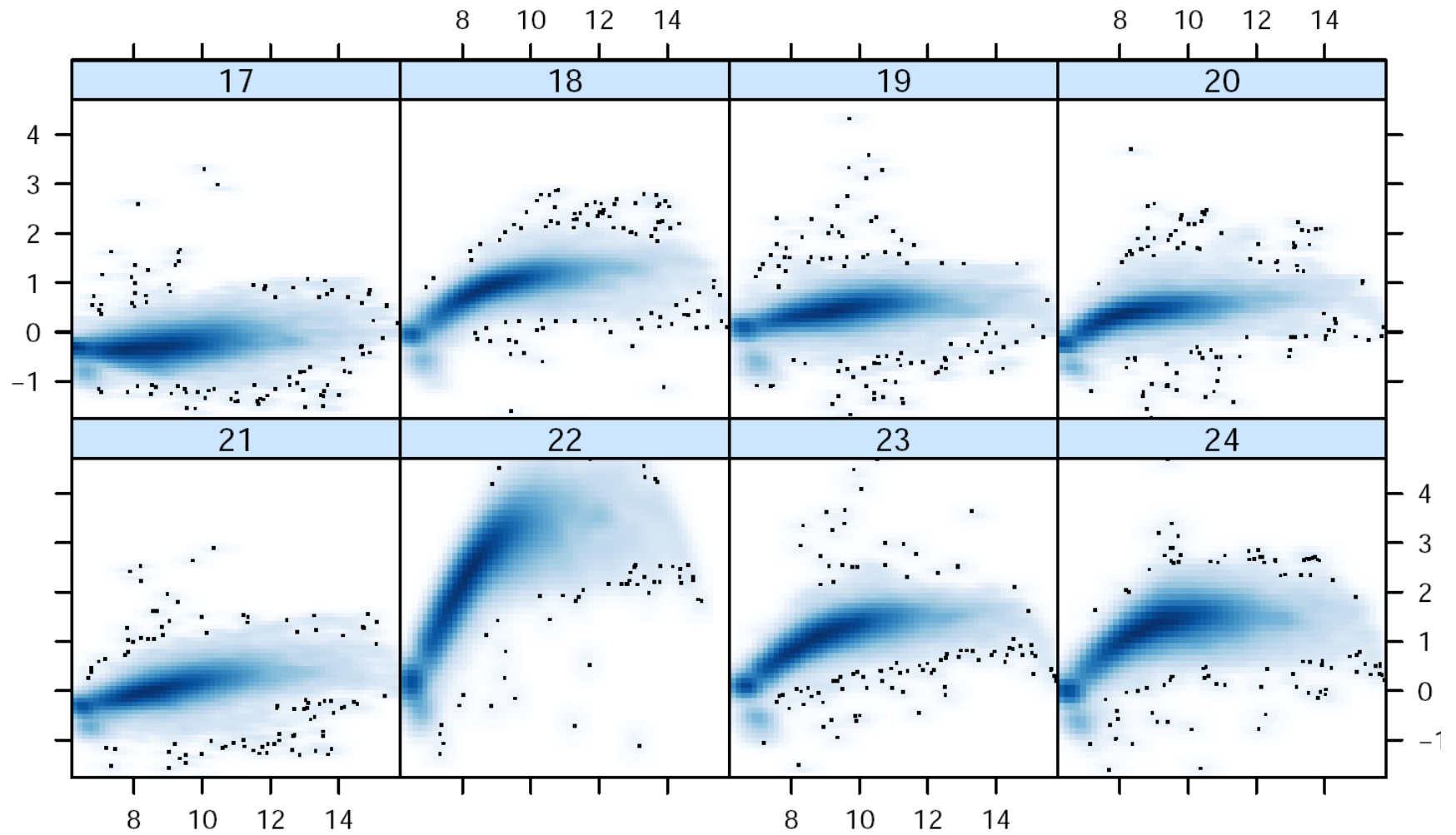
Validation of outlier's array specificity

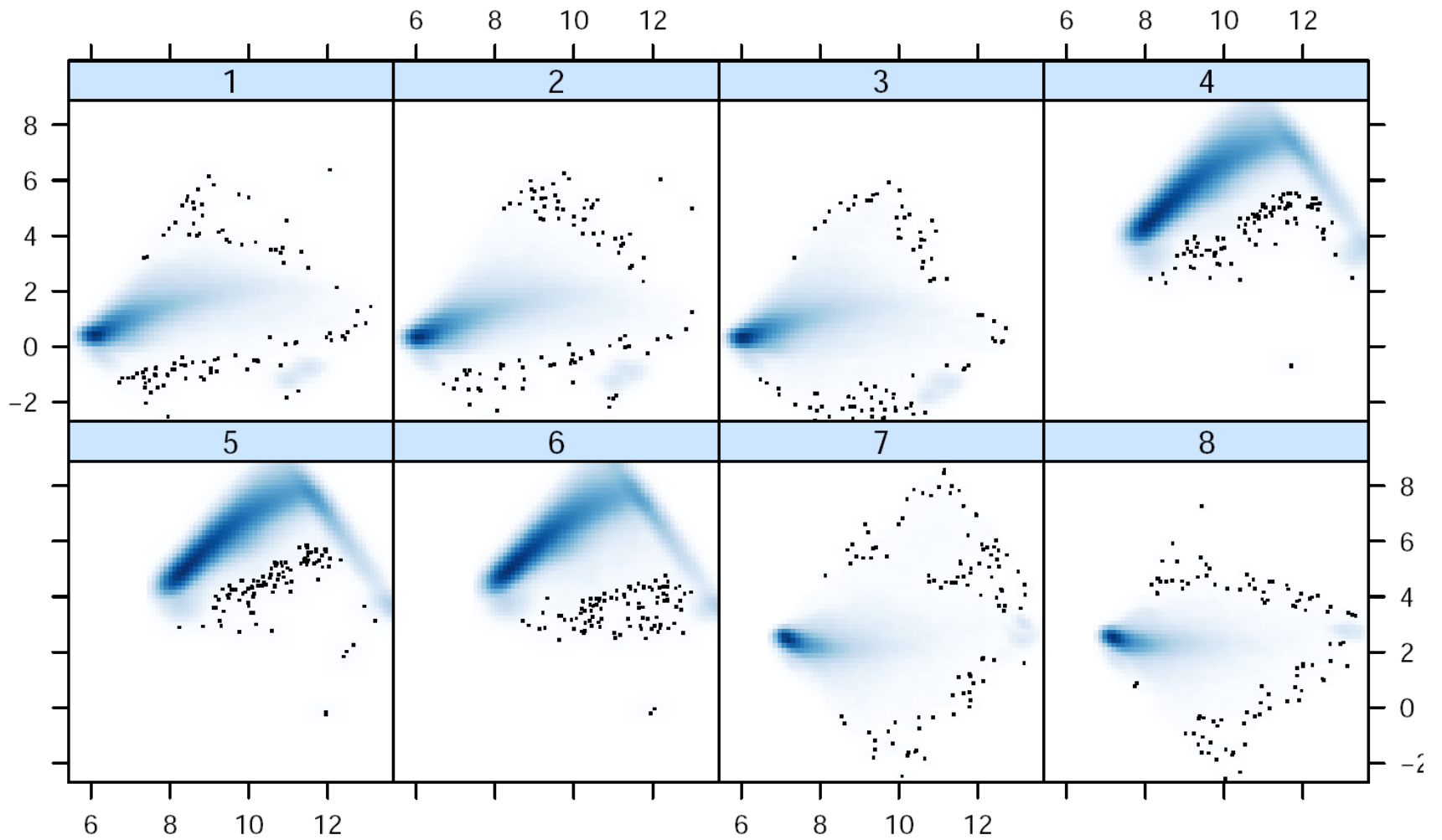
- Array #1 specific effect?

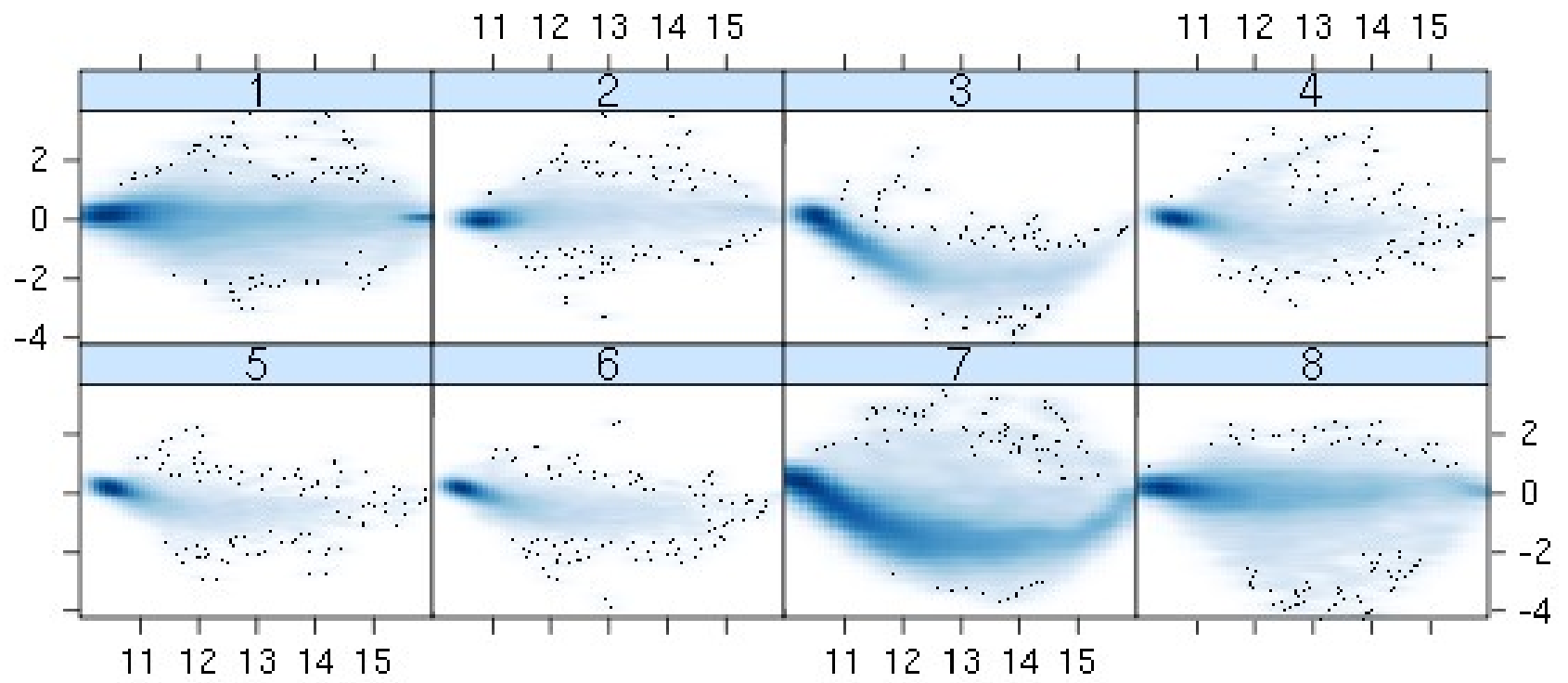
⇒ Remove one by one each array and perform the moderated t-test

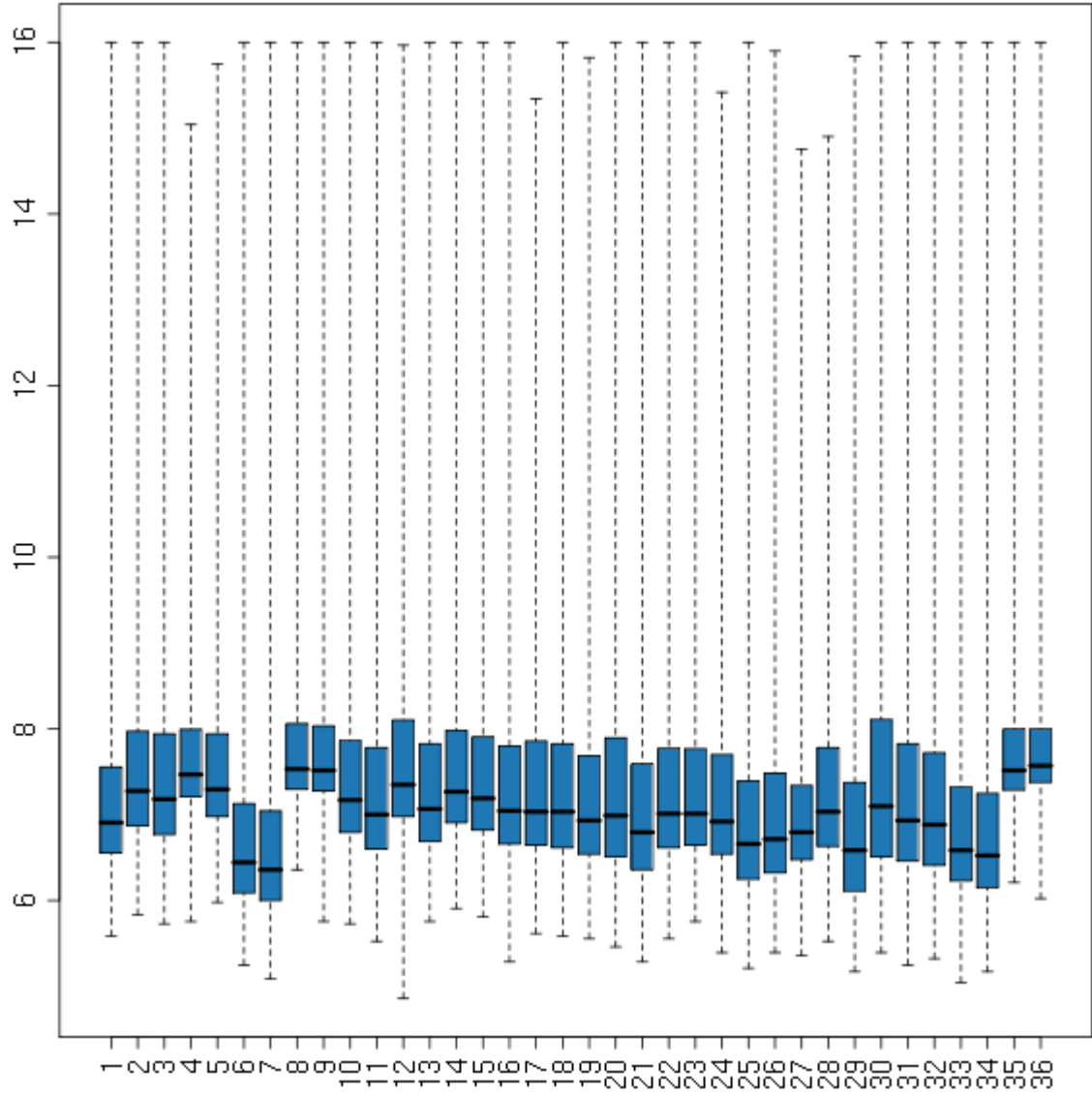
	# Genes	
	P < 0.01	P < 0.001
10 samples	34	4
Without sample 1	190	14
Without sample 2	39	3
Without sample 3	29	2
Without sample 4	21	1
Without sample 5	12	1
Without sample 6	87	5
Without sample 7	23	4
Without sample 8	34	4
Without sample 9	17	2
Without sample 10	23	2

Horror Picture Show

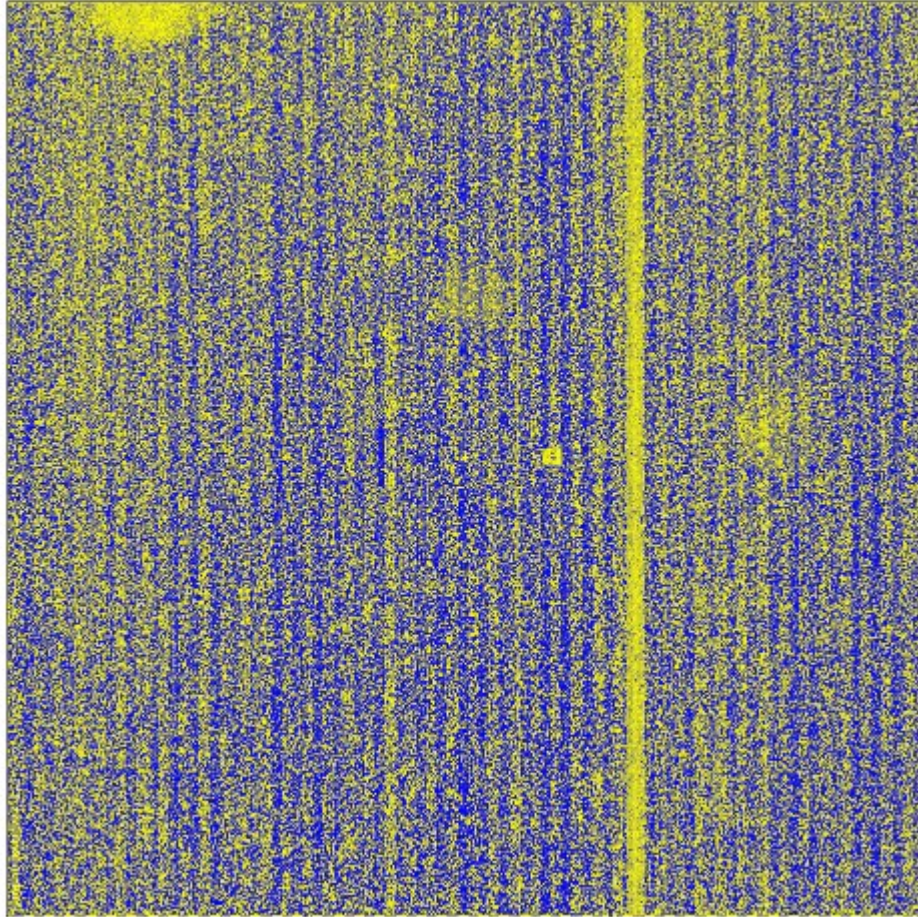




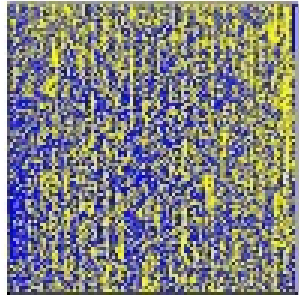




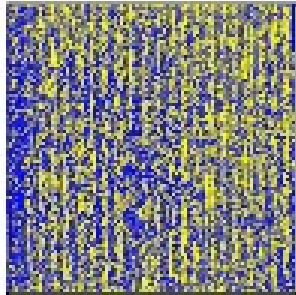
2



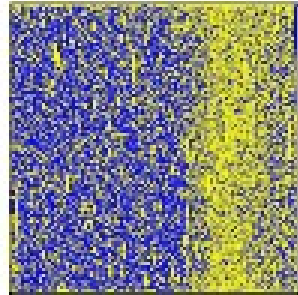
1



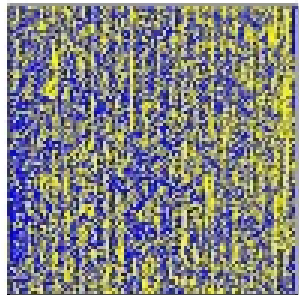
2



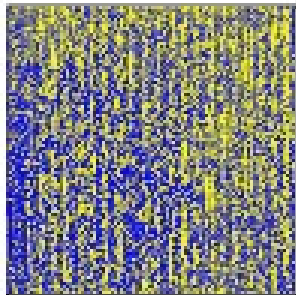
3



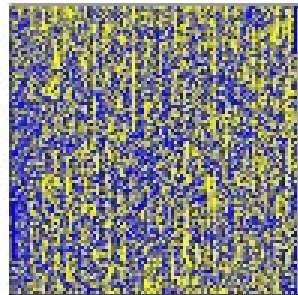
4

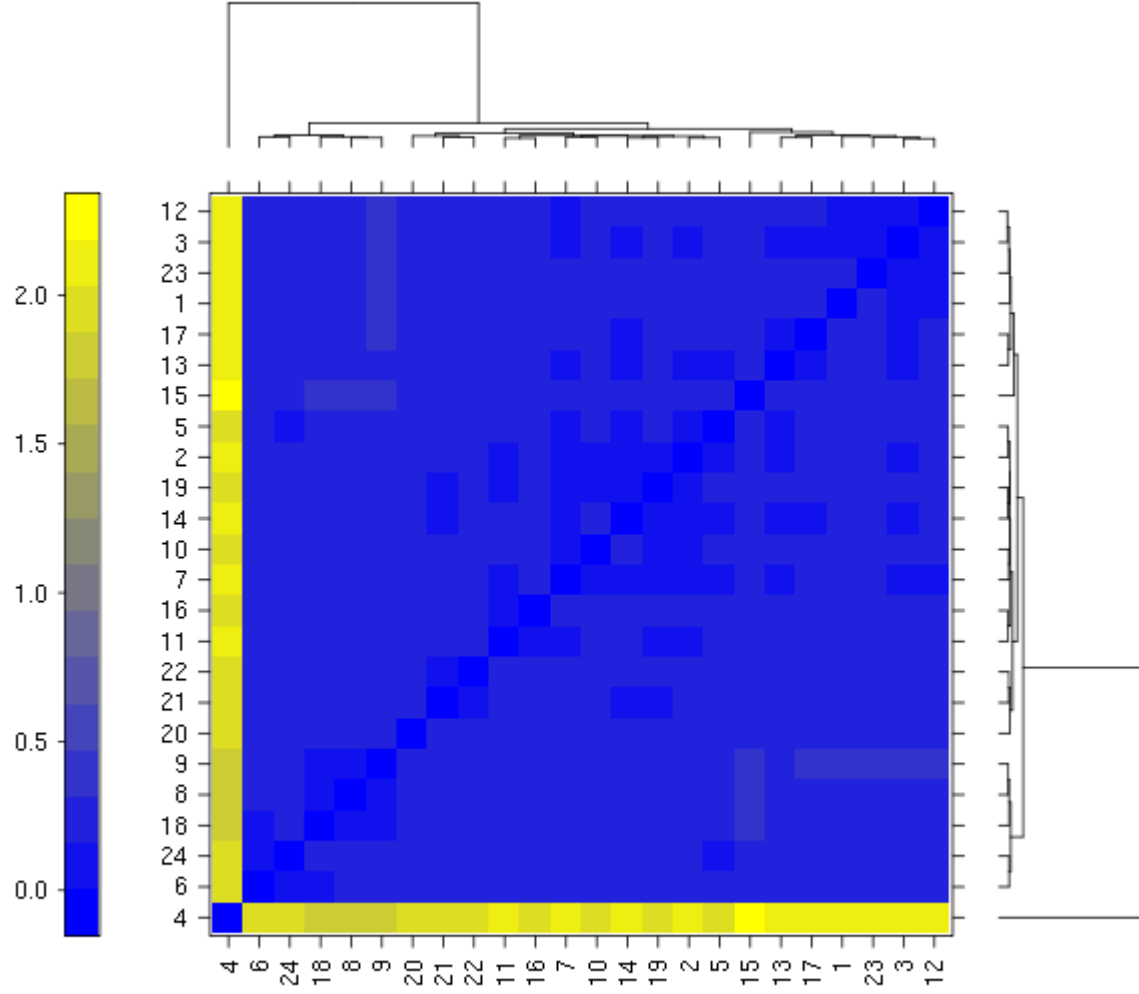


5

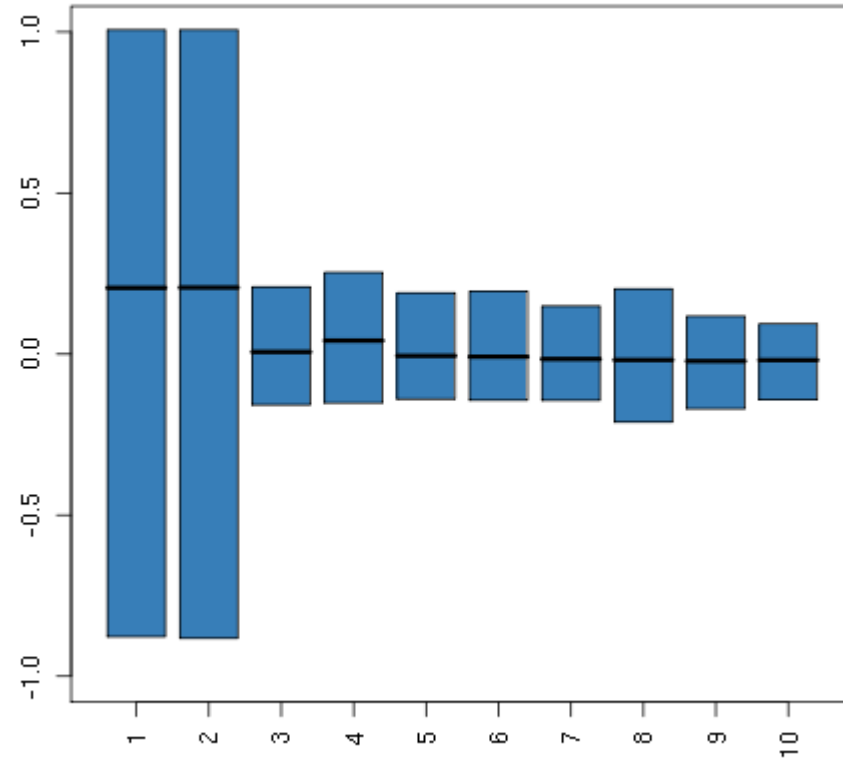


6

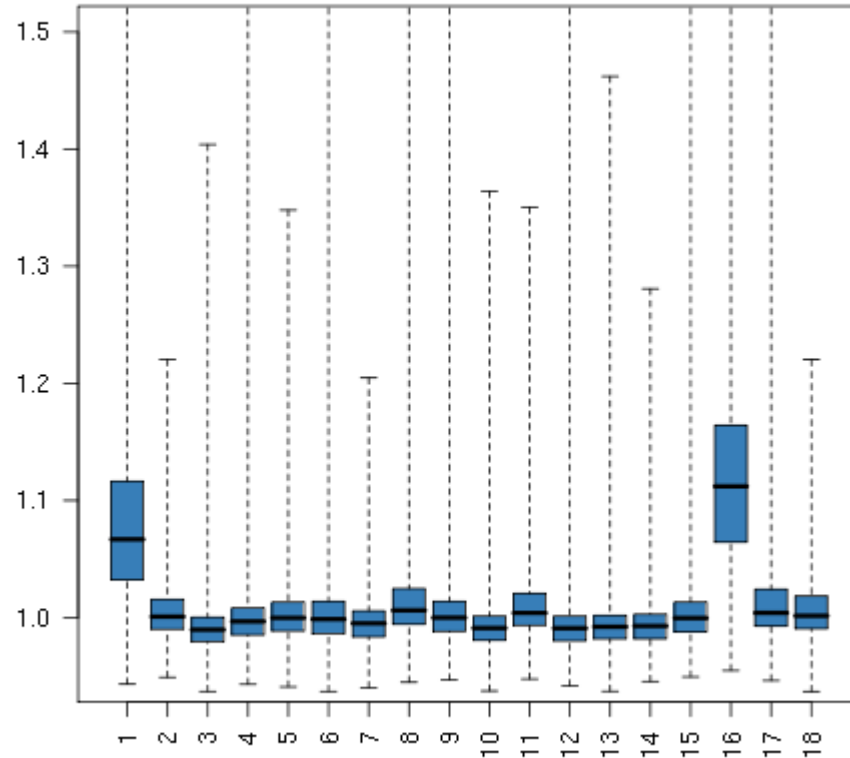




RLE



NUSE



Conclusions

- Quality assessment is important
 - Still needed
 - Give a first “taste” of the data
 - By removing outlier, the statistical power is increased
- arrayQualityMetrics package
 - One command line
 - Before preprocessing: to decide which normalisation
 - After normalisation: to check normalisation efficiency
 - Comprehensive report
 - Outlier detection