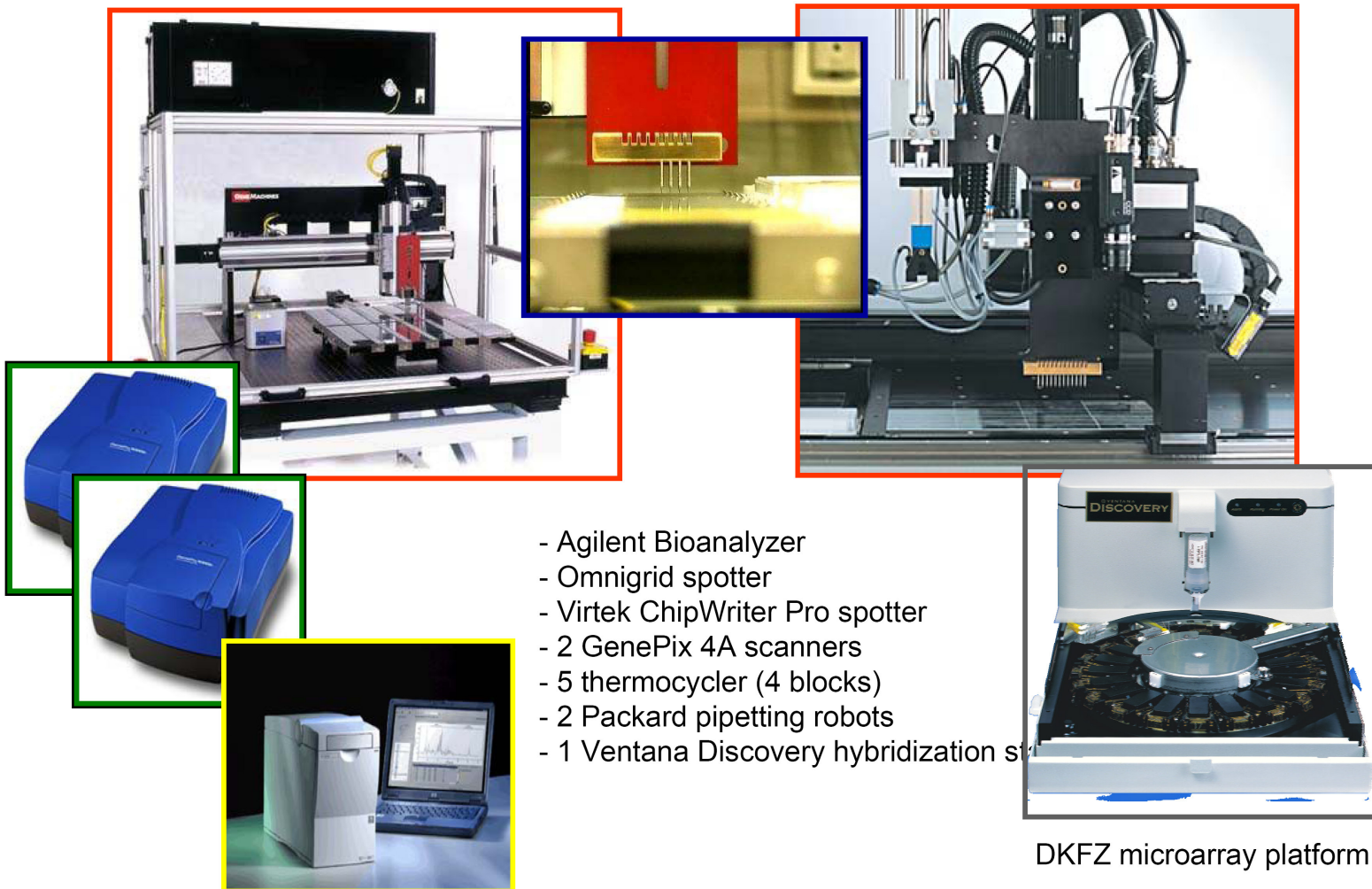

cDNA chips

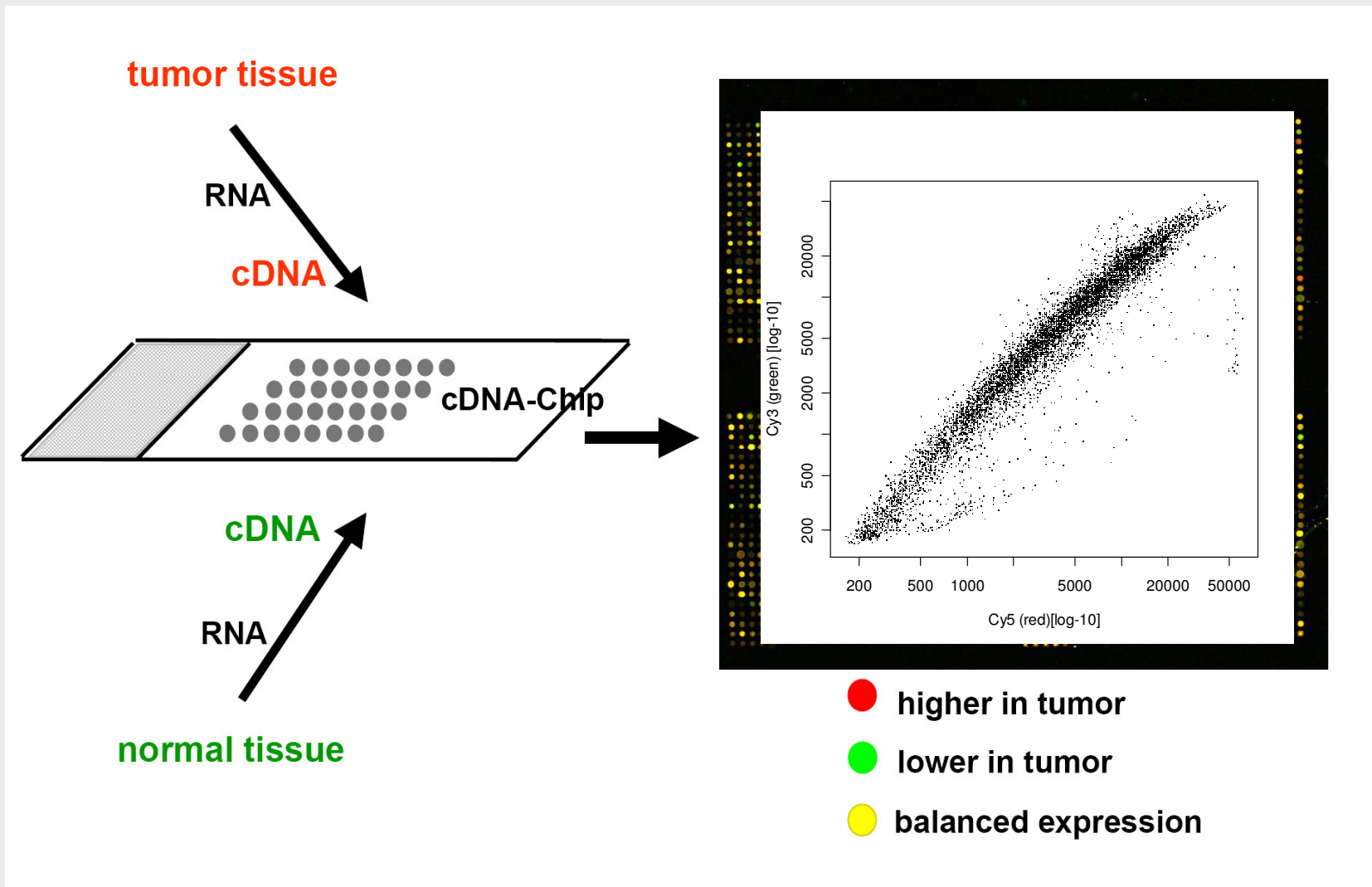
Quality control and pre-processing

Ulrich Mansmann
Department of Biometrics and Bioinformatics
Medical School, University of Munich



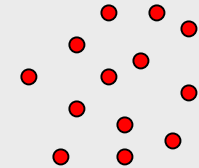
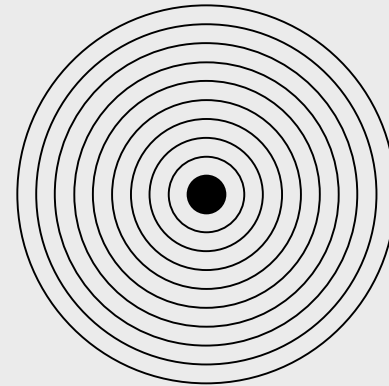
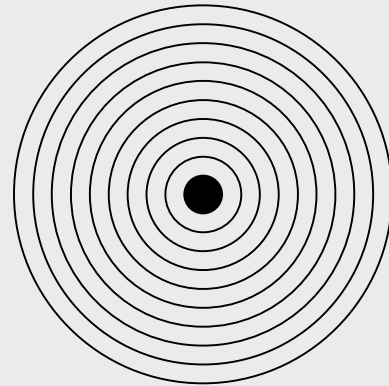
- Agilent Bioanalyzer
- Omnigrid spotter
- Virtek ChipWriter Pro spotter
- 2 GenePix 4A scanners
- 5 thermocycler (4 blocks)
- 2 Packard pipetting robots
- 1 Ventana Discovery hybridization station

DKFZ microarray platform

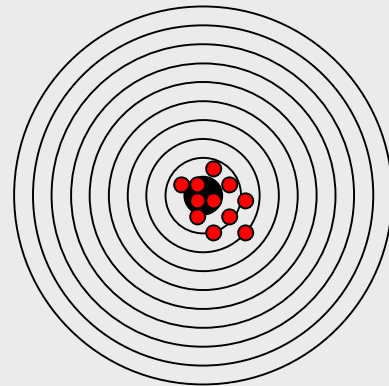


Measurements should be unbiased and precise

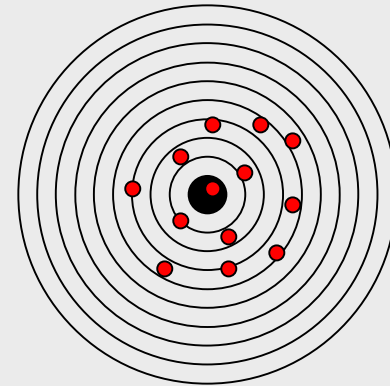
biased



unbiased



low noise



high noise

Swirl Data

This experiment was carried out using zebrafish as a model organism to study early development in vertebrates. Swirl is a point mutant in the BMP2 gene that affects the dorsal/ventral body axis. Ventral fates such as blood are reduced, whereas dorsal structures such as somites and notochord are expanded.

A goal of the Swirl experiment is to identify genes with altered expression in the swirl mutant compared to wild-type zebrafish. Two sets of dye-swap experiments were performed, for a total of four replicate hybridizations. For each of these hybridizations, target cDNA from the swirl mutant was labeled using one of the Cy3 or Cy5 dyes and the target cDNA wild-type mutant was labeled using the other dye.

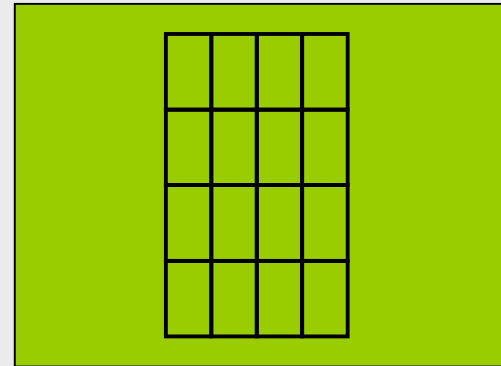
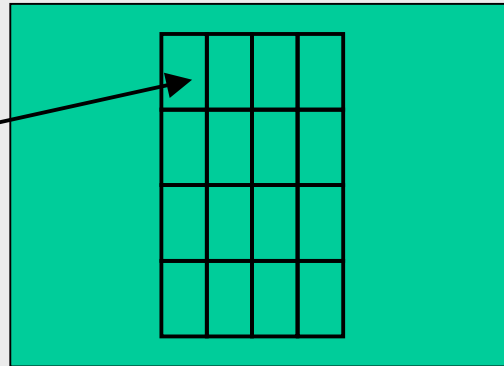
Target cDNA was hybridized to microarrays containing 8,448 cDNA probes, including 768 controls spots (e.g. negative, positive, and normalization controls spots). Microarrays were printed using 4 times 4 print-tips and are thus partitioned into a 4 times 4 grid matrix. Each grid consists of a 22 times 24 spot matrix that was printed with a single print-tip. Here, spot row and plate coordinates should coincide, as each row of spots corresponds to probe sequences from the same 384 well-plate.

Swirl Data

MT - R
WT - G

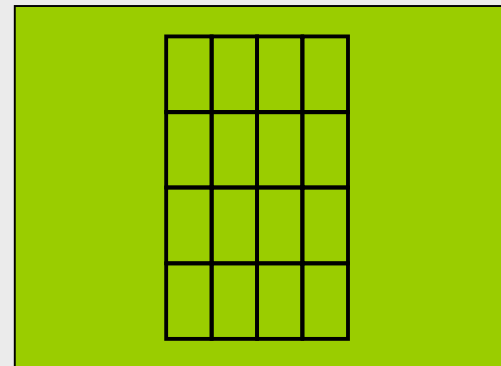
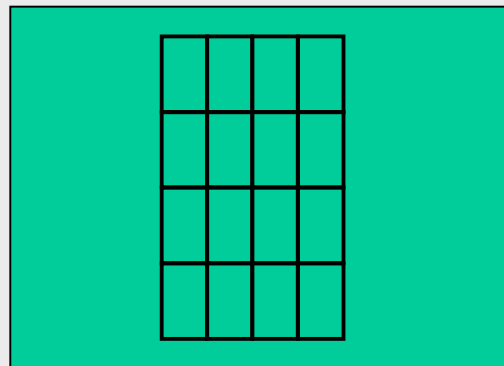
MT - G
WT - R

24 x 22
spots per
print-tip



Hybr. I

technical,
biological
variability



Hybr. II

cDNA - QC - Normalization

Quality control: Noise and reliable signal

- Is the signal dominated by noise? Acceptable amount of noise? Quantifying noise? (biol. / technol. variability), SNR
- Quantifying quality of a signal;
- Guidelines for reasonable thresholds on the quality of a signal;
- Defining strategies for exclusion of probes;
- Probe level: quality of the expression measurement on one particular array
- Gene level: quality of the expression measurement across all arrays
- Array level: quality of the expression measurement on one particular glass slide

Probe-level quality control

- Individual spots printed on the slide
- Sources: faulty printing, uneven distribution, contamination with debris, size of signal relative to noise, poorly measured spots;
- Visual inspection: hairs, dust, scratches, air bubbles, dark regions, regions with haze
 - set points to NA
 - local normalization procedures which account for regional idiosyncrasies.

Visual inspection

4 x 4 sectors

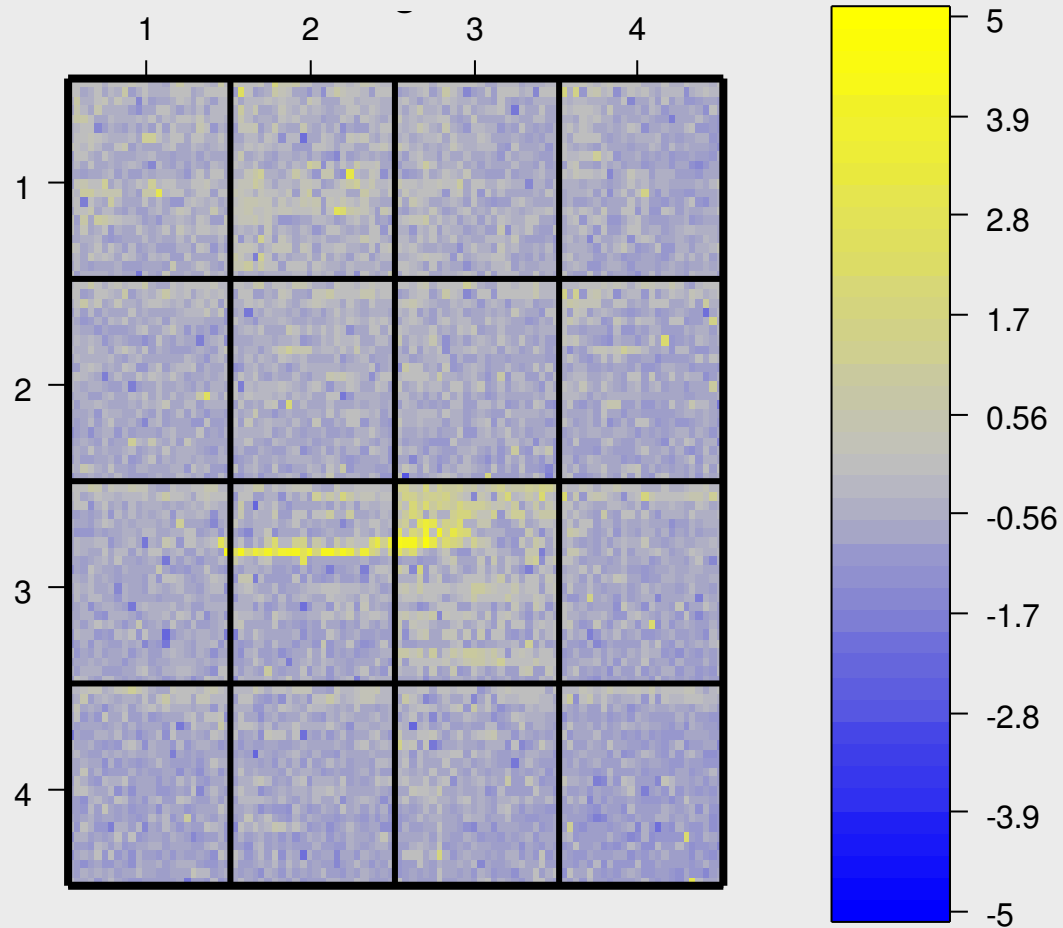
Sector:

24 rows

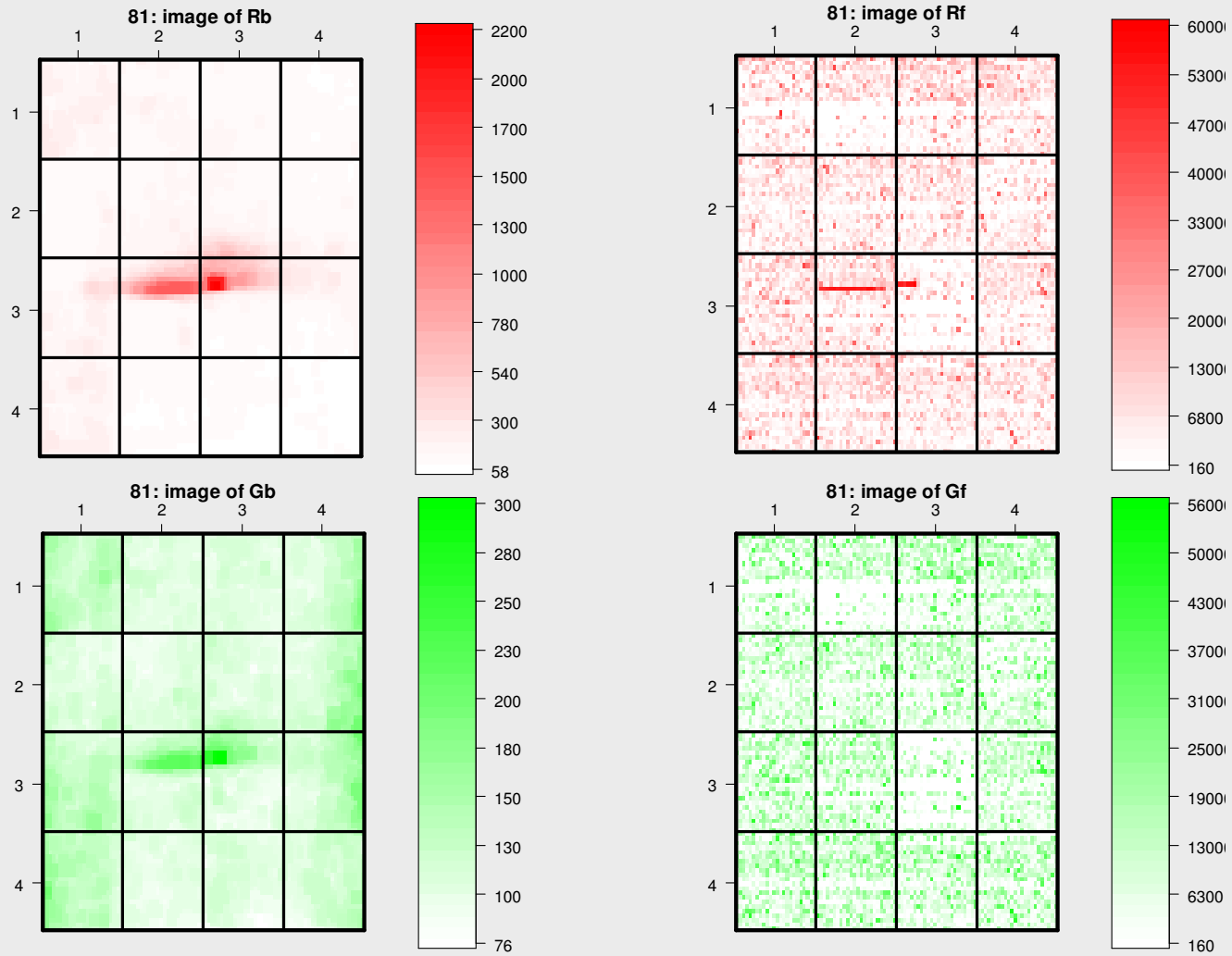
22 columns

8448 spots

Mean signal intensity



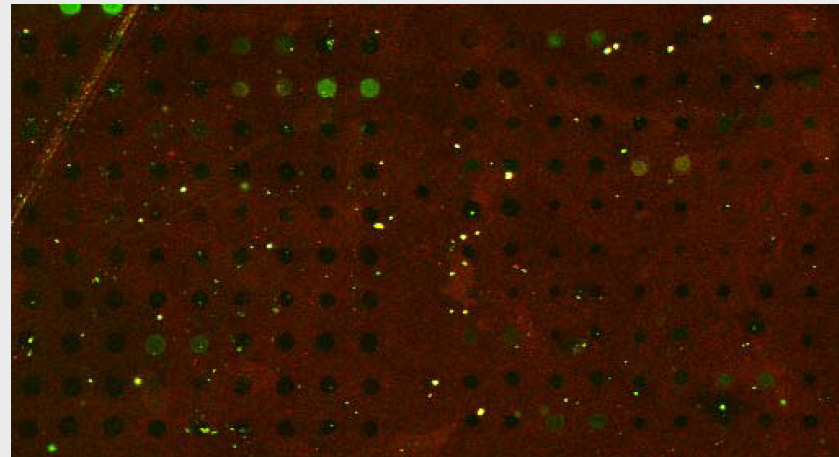
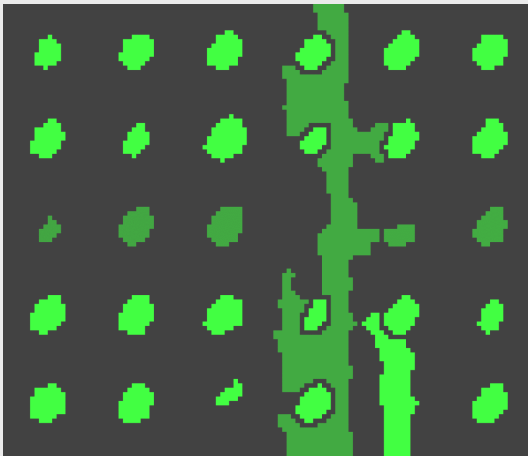
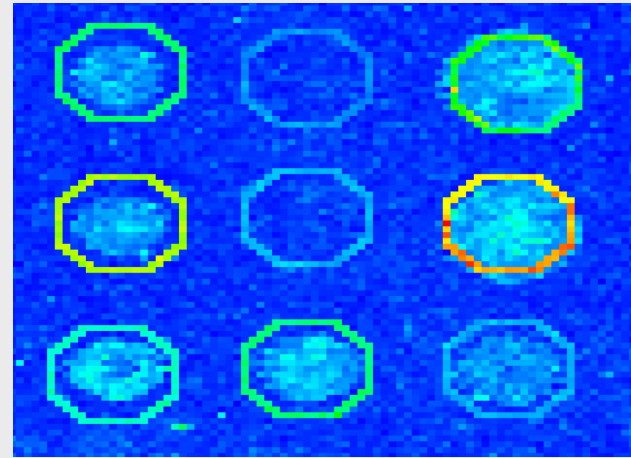
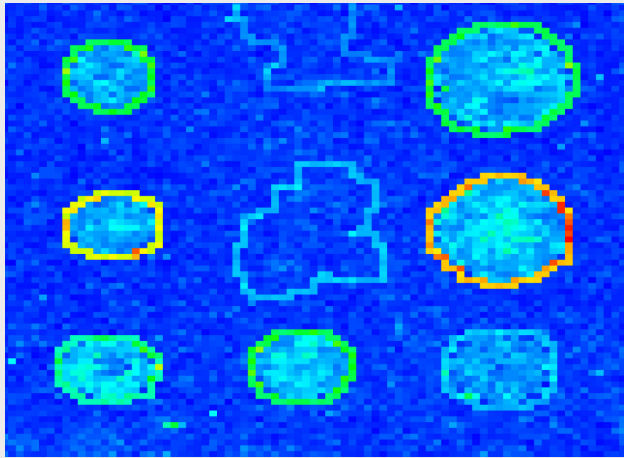
Visual inspection



cDNA - QC - Normalization

Probe-level quality control

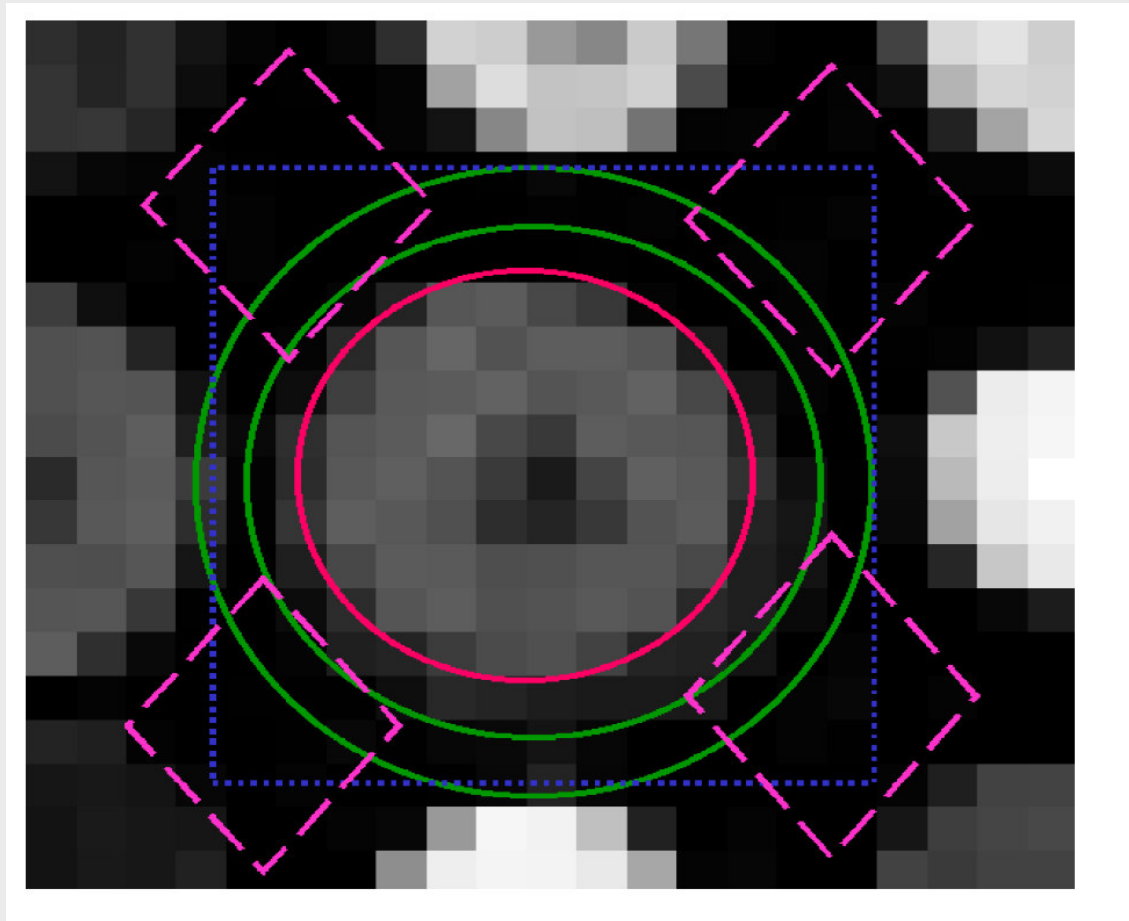
- Individual spots printed on the slide
- Sources: faulty printing, uneven distribution, contamination with debris, size of signal relative to noise, poorly measured spots;
- Visual inspection
- Spot quality
 - Brightness*: foreground/background ratio
 - Uniformity*: variation in pixel intensities and ratios of intensities within a spot
 - Morphology*: area, perimeter, circularity.
 - Spot Size*: number of foreground pixels



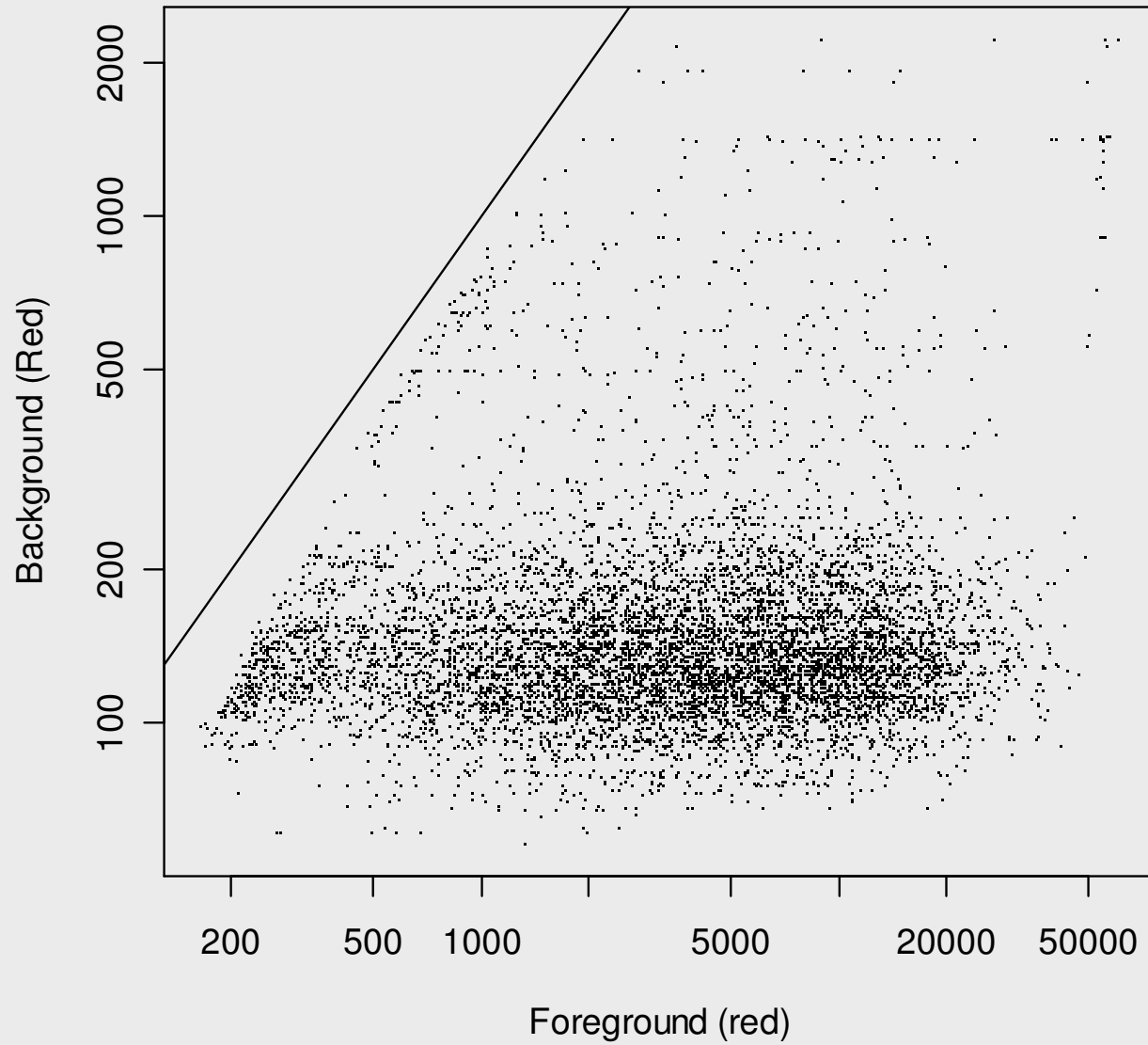
GenePix

QuantArray

ScanAlyse



swirl.1.spot



cDNA - QC - Normalization

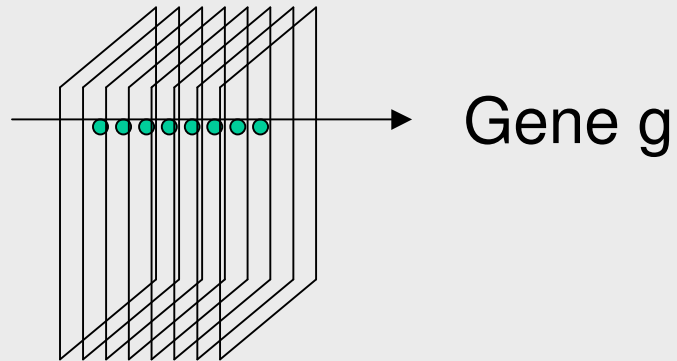
Weak Signal



CY5 channel	CY3 channel	Measure
$\text{FR-BR} < \lambda_{\text{low}}$	$\text{FG-BG} < \lambda_{\text{high}}$	exclude
$\text{FR-BR} < \lambda_{\text{high}}$	$\text{FG-BG} < \lambda_{\text{low}}$	exclude
$\text{FR-BR} < \lambda_{\text{low}}$	$\text{FG-BG} > \lambda_{\text{high}}$	$\text{Log}_2[\lambda_{\text{low}} / (\text{FG-BG})]$
$\text{FR-BR} > \lambda_{\text{high}}$	$\text{FG-BG} < \lambda_{\text{low}}$	$\text{Log}_2[(\text{FR-BR}) / \lambda_{\text{low}}]$
$\text{FR-BR} > \lambda_{\text{high}}$	$\text{FG-BG} > \lambda_{\text{high}}$	$\text{Log}_2[(\text{FR-BR}) / (\text{FG-BG})]$

$\text{Log}_2[510/490] = 0.057, \text{Log}_2[30/10] = 1.58, \quad \text{FG}=2030, \text{BG}=2000$

Gene-level quality control



- Test of poor hybridization in the reference channel may introduce bias on the fold-change
- Test for exclusion of low variance genes.

Gene-level quality control: Poor Hybridization and Printing

- Some probes will not hybridize well to the target RNA
- Printing problems such that all spots of a given inventory well have poor quality.
- A well may be of bad quality – contamination
- Genes with a consistently low signal in the reference channel are suspicious: Median of the background adjusted signal $< 200^*$

*or other appropriate choice

Gene-level quality control: Probe quality control based on duplicated spots

- Printing different probes that target the same gene or printing multiple copies of the same probe.
- Mean squared difference of \log_2 ratios between spot r and s :

$$\text{MSDLR} = \Sigma(x_{jr} - x_{js})^2/J \quad \text{sum over arrays } j = 1, \dots, J$$

recommended threshold to assess disagreement: $\text{MSDLR} > 1$

- Disagreement between copies: printing problems, contamination, mislabelling. Not easy if there are only 2 or 3 slides.
- Jenssen et al (2002) Nucleic Acid Res, 30: 3235-3244.
Theoretical background

Gene-level quality control: Low variance genes

- Good for normalization, but uninformative for the analysis
- Threshold for fold change:
 $\log_2(\max) - \log_2(\min)$, this difference may depend on sample size
taking sample size into account: $\log_2(q_{0.95}) - \log_2(q_{0.05})$
- Variance based criterion (cut point)

Array-level quality control

- Problems: array fabrication defect
 - problem with RNA extraction
 - failed labelling reaction
 - poor hybridization conditions
 - faulty scanner
- Quality measures:
 - Percentage of spots with no signal (~30% excluded spots)
 - Range of intensities
 - $(\text{Av. Foreground})/(\text{Av. Background}) > 3$ in both channels
 - Distribution of spot signal area
 - Amount of adjustment needed: signals have to substantially changed to make slides comparable.

See next lecture

Raw data are not mRNA concentrations

- tissue contamination
- RNA degradation
- amplification efficiency
- reverse transcription efficiency
- hybridization efficiency and specificity
- clone identification and mapping
- PCR yield, contamination
- spotting efficiency
- DNA-support binding
- other array manufacturing-related issues
- image segmentation
- signal quantification
- 'background' correction

W. Huber

Sources of variation

amount of RNA in the biopsy
efficiencies of

- RNA extraction
- reverse transcription
- labeling
- photodetection

PCR yield
DNA quality
spotting efficiency,
spot size
cross-/unspecific hybridization
stray signal

Systematic

- o similar effect on many measurements
- o corrections can be estimated from data

Calibration

Stochastic

- o too random to be explicitly accounted for
- o "noise"

Error model

W. Huber

Normalization

- Identify and remove sources of systematic variation, other than differential expression, in the measured fluorescence intensities.
- Normalization is necessary before analysis is performed, in order to ensure that differences in intensities are indeed due to differential expression and not experimental artefacts.

- **Location** normalisation: corrects for spatial or dye bias

- **Scale** normalisation: homogenises the variability across arrays

$$\text{MAD} = \text{median}\{ |x_1 - m|, \dots, |x_n - m| \}$$

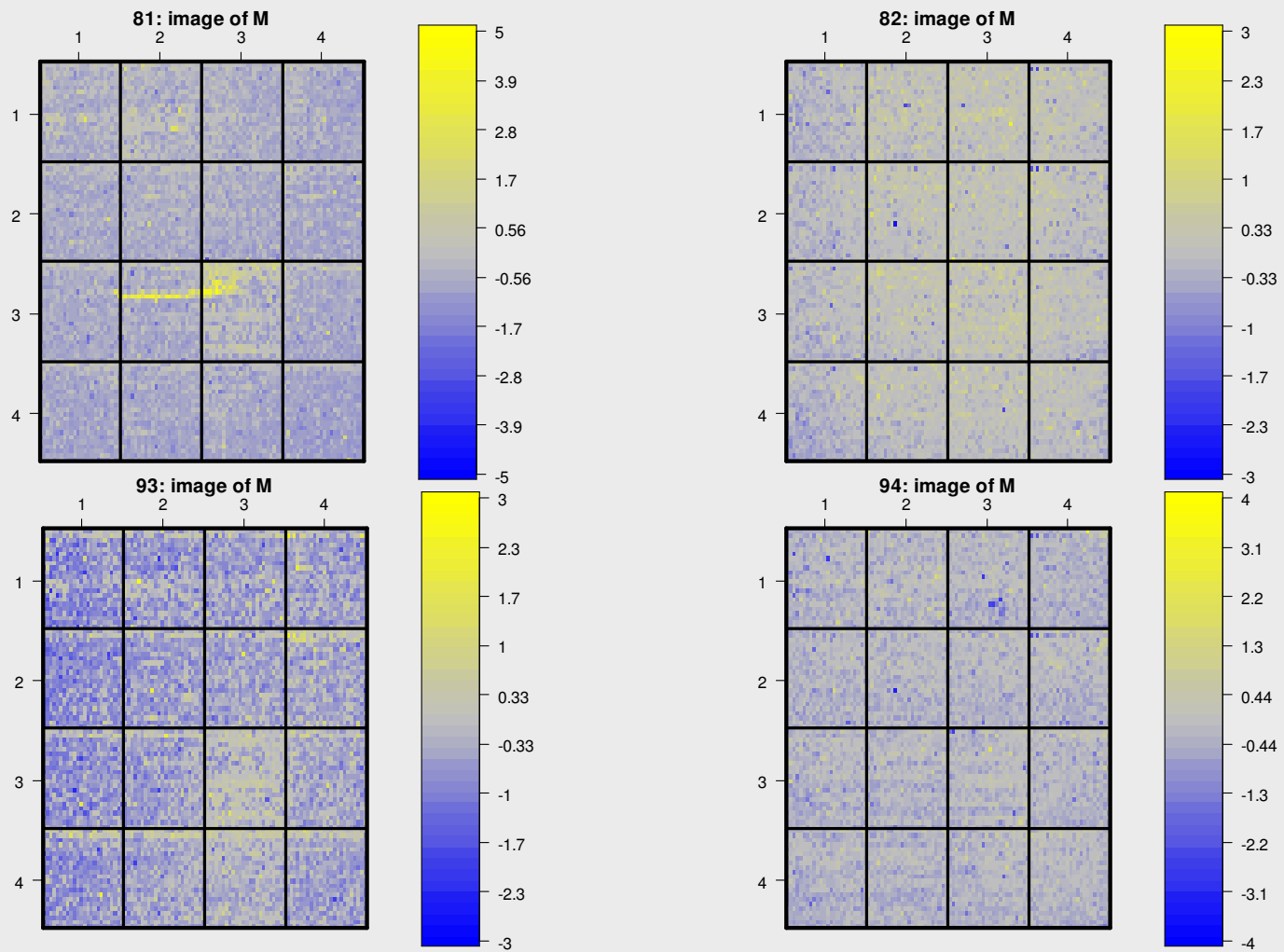
- Location and scale are basic statistical concepts for data description.

- Normalised log-intensity ratios are given by

$$M_{\text{norm}} = (M - \text{loc})/s$$

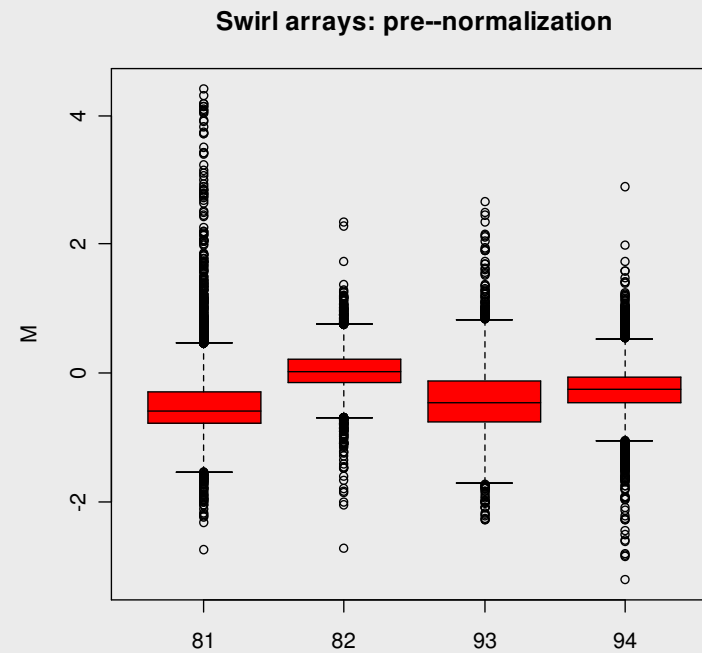
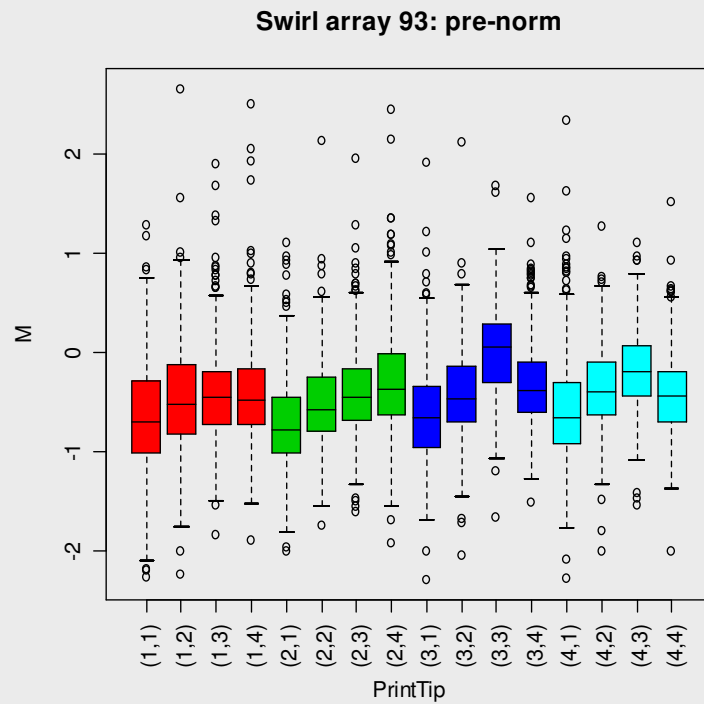
- Normalisation: within arrays (marray) or between arrays (vsnp), single channels between arrays, log expression ratios, etc

Swirl Data



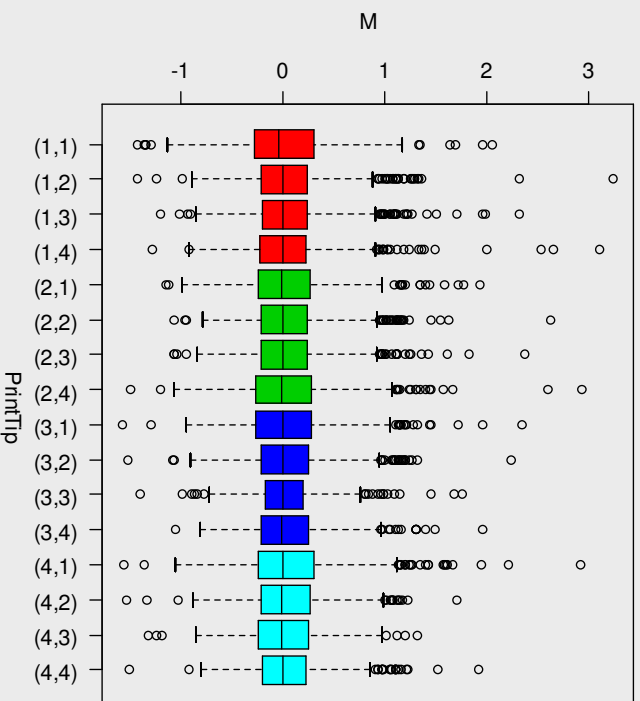
cDNA - QC - Normalization

marray: Pre Normalization

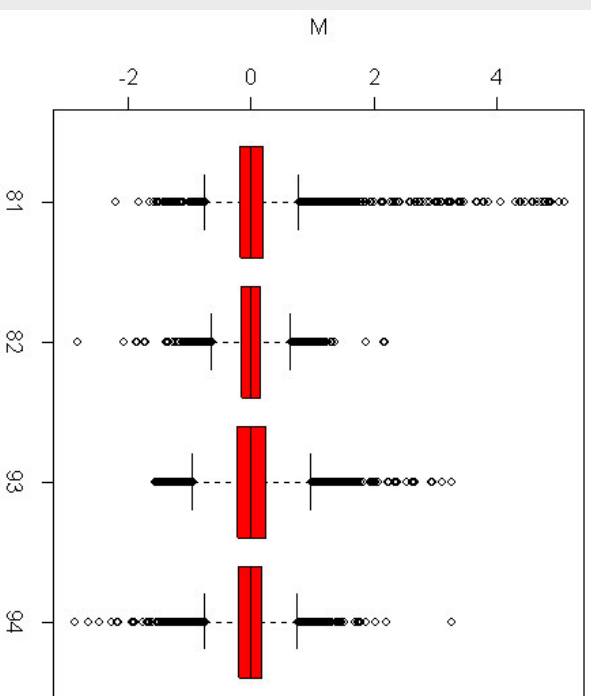


marray: Post Normalization

Swirl array 93: post-norm

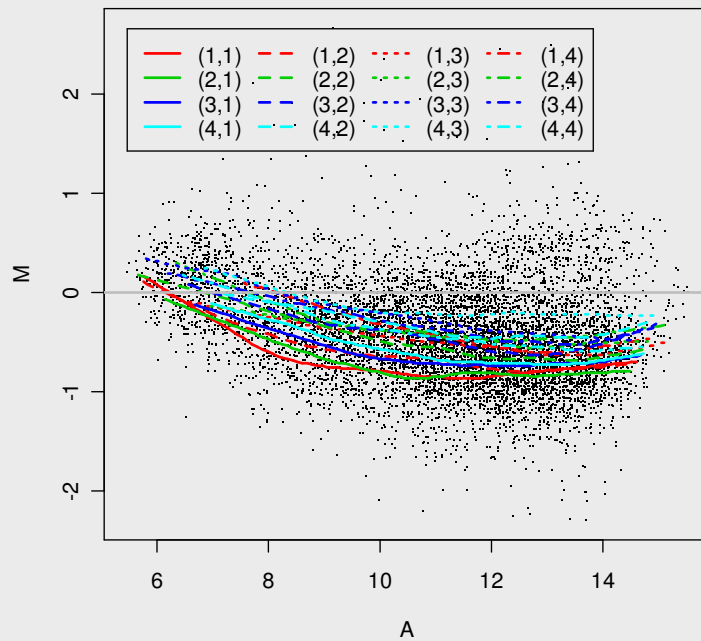


Swirl arrays: post-normalization

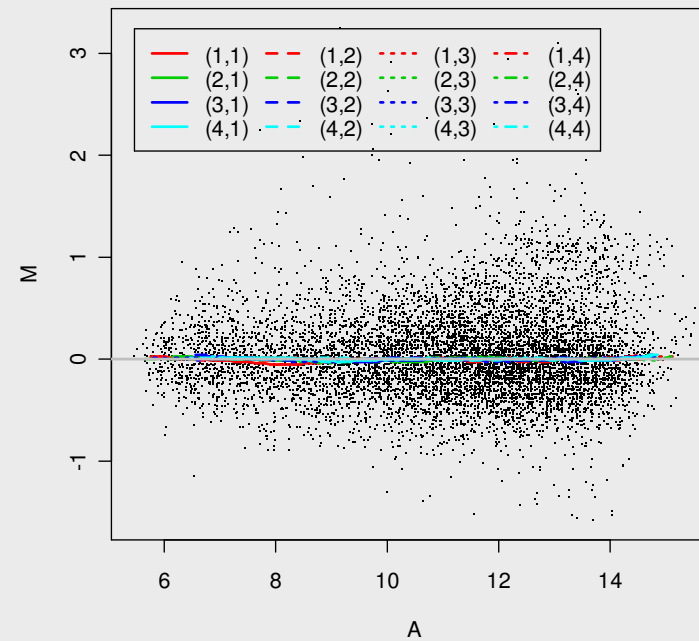


marray: MA Plots

Swirl array 93: pre-norm MA-Plot

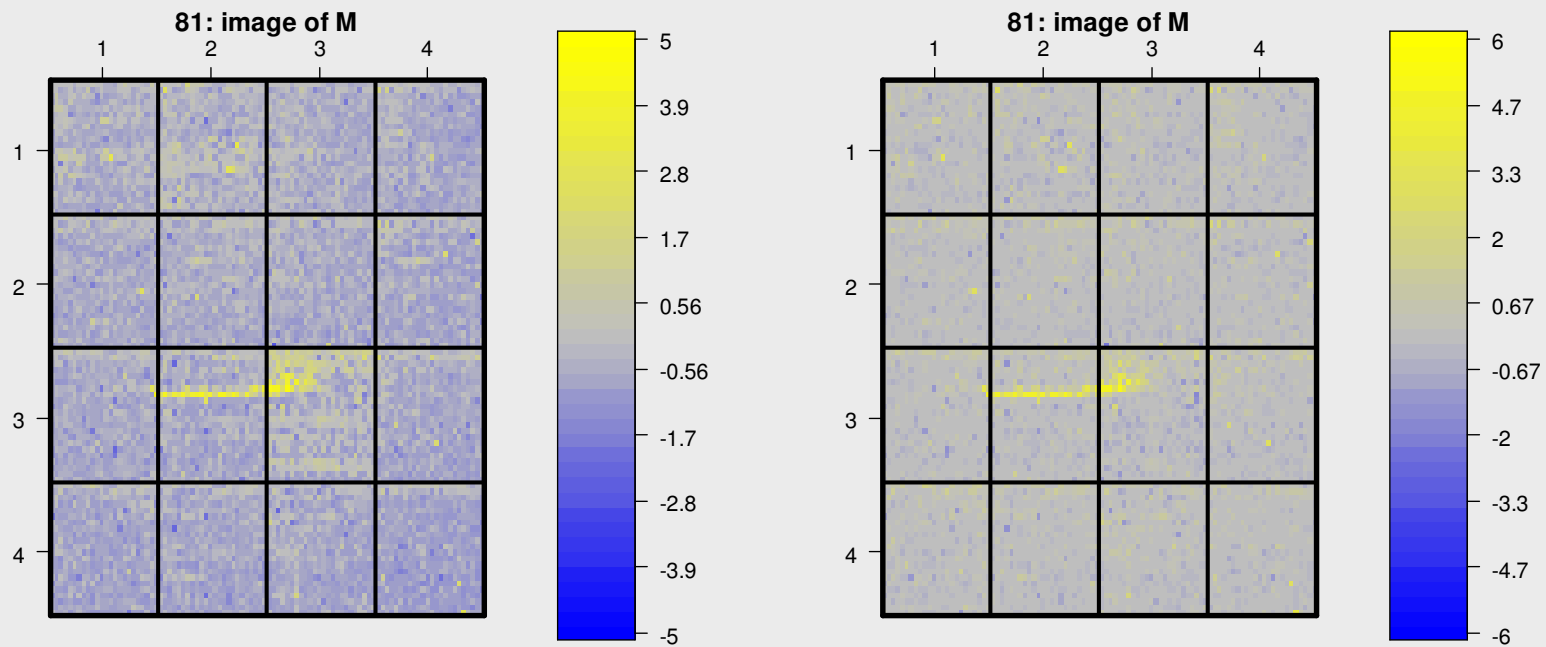


Swirl array 93: post-norm MA-Plot



Non-parametric smoother: loess, lowess, local regression line, generalizes the concept of moving average.

marray:
Signals - raw versus normalised



Normalization procedure was not able to remove scratch

VSN: model and theory

- Huber et al. (2002) *Bioinformatics*, 18:S96–S104
- Model for measured probe intensity
Rocke DM, Durbin B (2001) *Journal of Computational Biology*, 8:557–569
- log-transformation is replaced by a transformation (arcsinh) based on theoretical grounds.
- Estimation of transformation parameters (location, scale) based on ML paradigm and numerically solved by a least trimmed sum of squares regression.
- vsn-normalised data behaves close to the normal distribution
- Following slides are borrowed from Wolfgang Huber

► modeling ansatz

measured intensity = offset + gain × true abundance

$$y_{ik} = a_{ik} + b_{ik} x_k$$

$$a_{ik} = a_i + \varepsilon_{ik}$$

a_i per-sample offset

$$\varepsilon_{ik} \sim N(0, b_i^2 s_1^2)$$

“additive noise”

i - sample / array

k - probe / gene

$$b_{ik} = b_i b_k \exp(\eta_{ik})$$

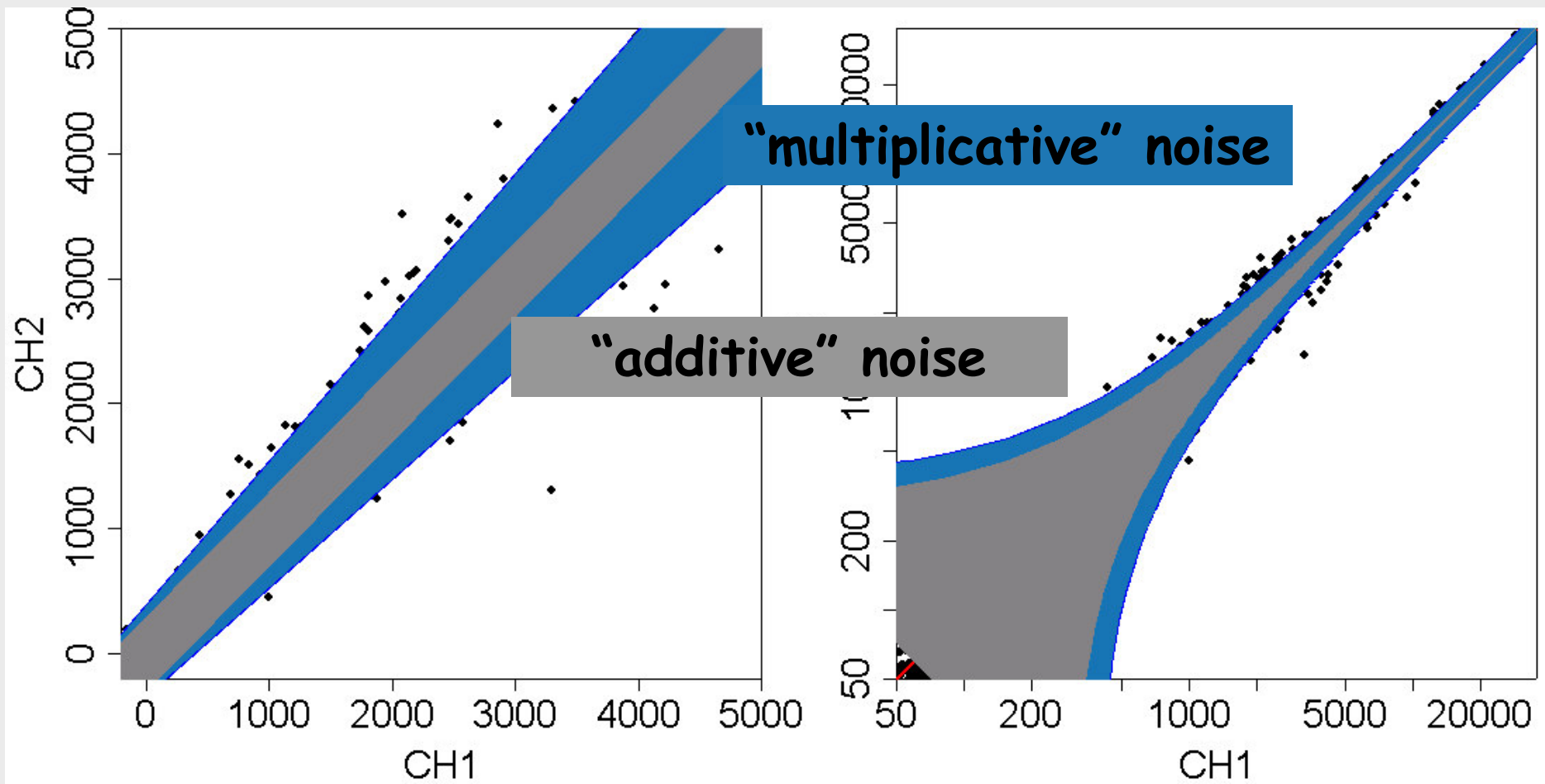
b_i per-sample
normalization factor

b_k sequence-wise
probe efficiency

$$\eta_{ik} \sim N(0, s_2^2)$$

“multiplicative noise”

► The two-component model



raw scale

log scale

▶ variance stabilizing transformations

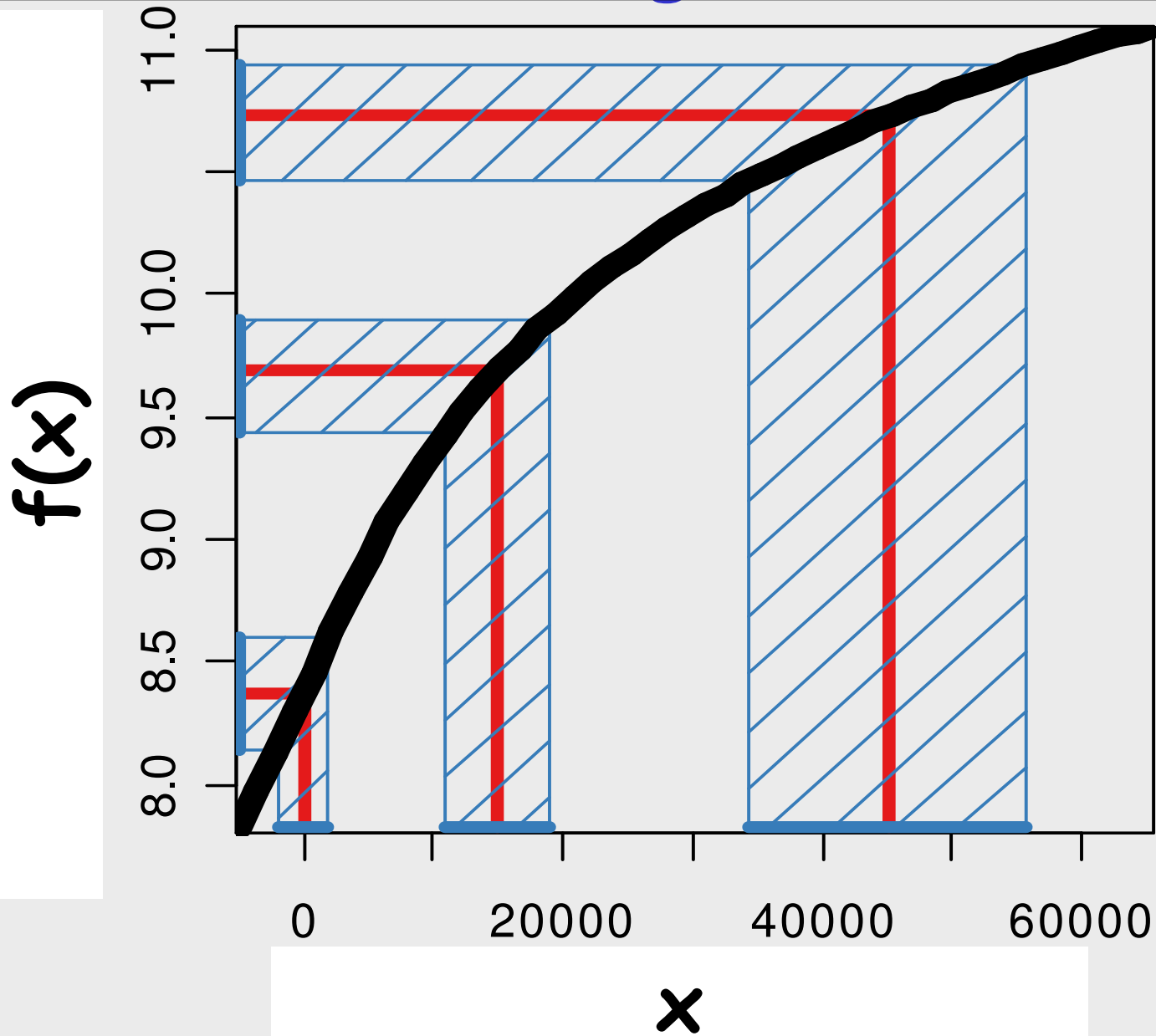
X_u a family of random variables with $EX_u = u$, $\text{Var} X_u = v(u)$. Define

$$f(x) = \int^x \frac{1}{\sqrt{v(u)}} du$$

$\Rightarrow \text{var } f(X_u) \approx \text{independent of } u$

derivation: linear approximation

▶ variance stabilizing transformations



▶ variance stabilizing transformations

$$f(x) = \int^x \frac{1}{\sqrt{v(u)}} du$$

1.) constant variance ('additive') $v(u) = s^2 \Rightarrow f \propto u$

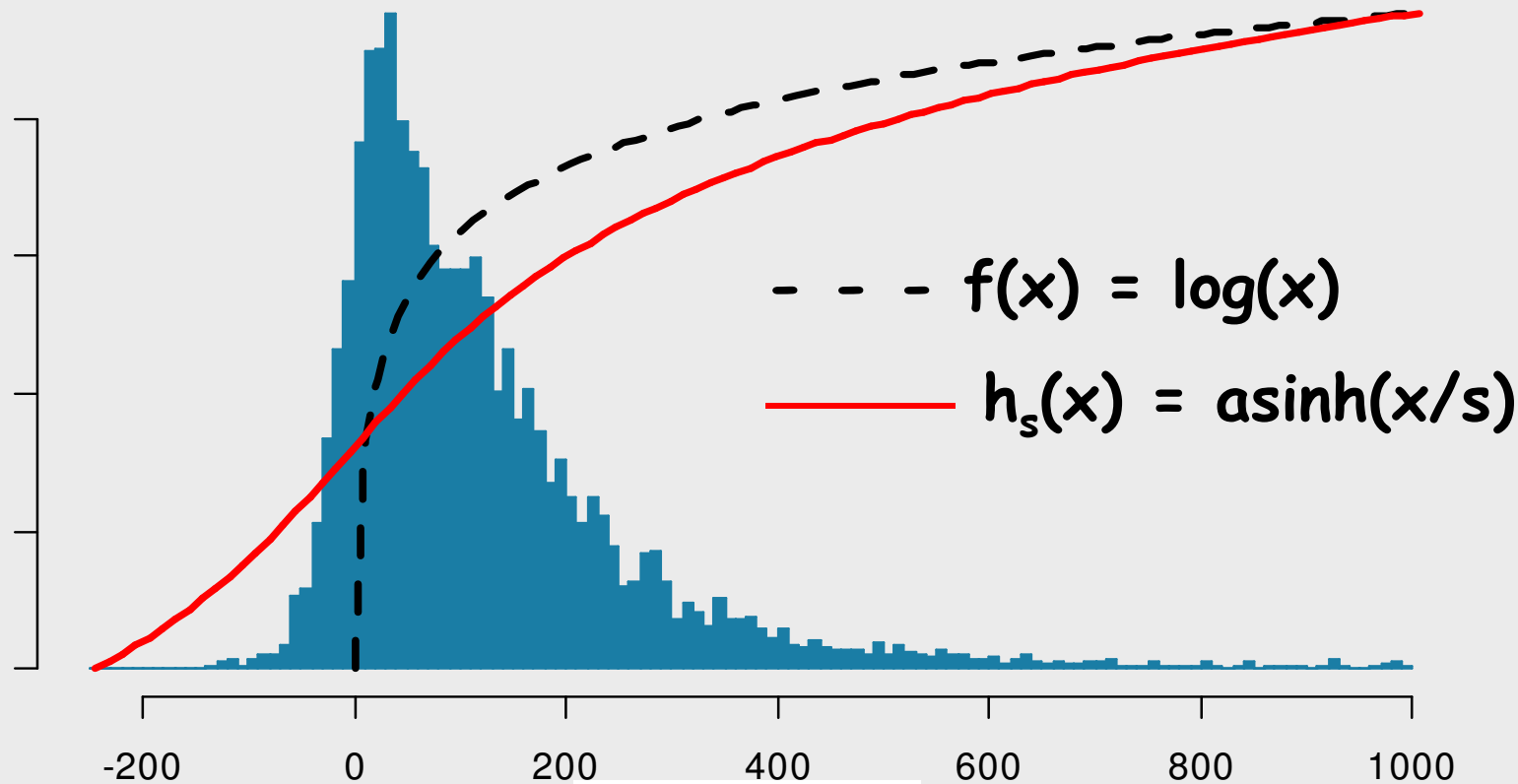
2.) constant CV ('multiplicative') $v(u) \propto u^2 \Rightarrow f \propto \log u$

3.) offset $v(u) \propto (u + u_0)^2 \Rightarrow f \propto \log(u + u_0)$

4.) additive and multiplicative

$$v(u) \propto (u + u_0)^2 + s^2 \Rightarrow f \propto \operatorname{arsinh} \frac{u + u_0}{s}$$

▶ the "glog" transformation



$$\operatorname{arsinh}(x) = \log \left(x + \sqrt{x^2 + 1} \right)$$

$$\lim_{x \rightarrow \infty} (\operatorname{arsinh} x - \log x - \log 2) = 0$$

P. Munson, 2001

D. Rocke & B. Durbin, ISMB 2002

parameter estimation

$$\operatorname{arsinh} \frac{y_{ki} - a_i}{b_i} = \mu_k + \varepsilon_{ki}, \quad \varepsilon_{ki} : N(0, c^2)$$

- o maximum likelihood - but sensitive to outliers
- o model holds differentially
- o robust variance
- Trimmed Sum
- o works as local differentially

measured intensity = offset + gain * true abundance

$$y_{ik} = a_{ik} + b_{ik} x_{ik}$$

$$a_{ik} = a_i + L_{ik} + \varepsilon_{ik}$$

a_i per-sample offset

L_{ik} local background provided by image analysis

$$\varepsilon_{ik} \sim N(0, b_i^2 s_i^2)$$

"additive noise"

$$b_{ik} = b_i b_k \exp(\eta_{ik})$$

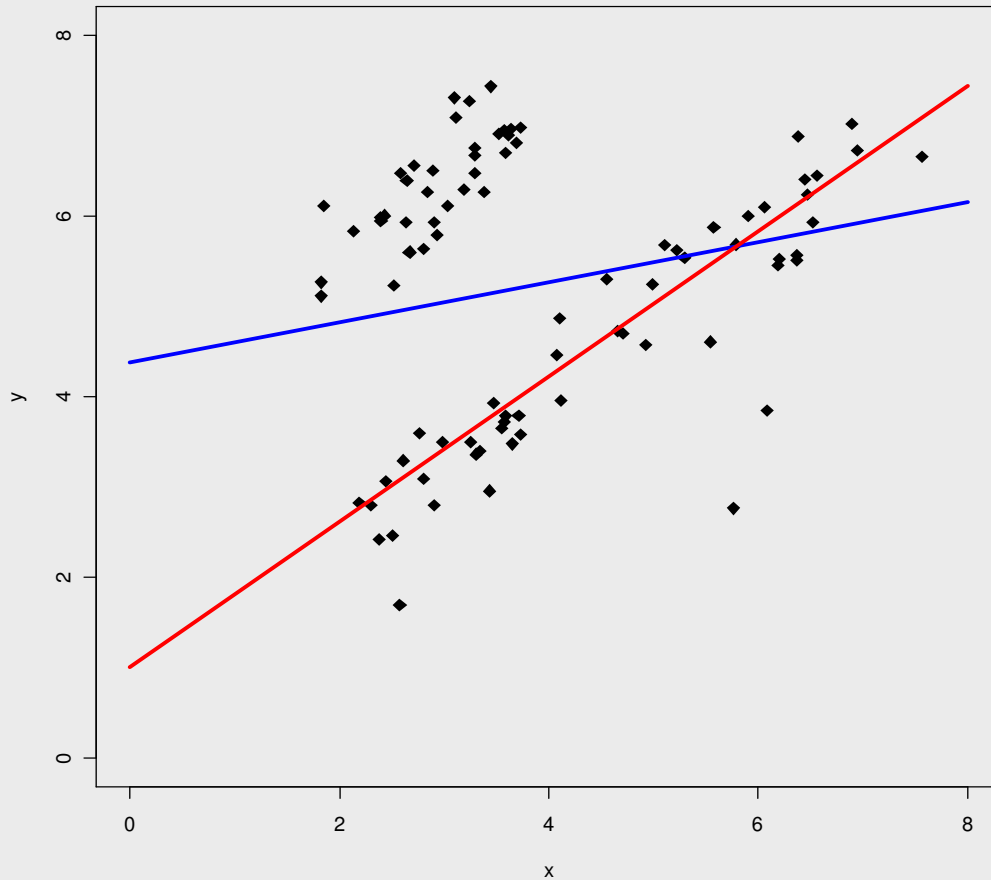
b_i per-sample normalization factor

b_k sequence-wise labeling efficiency

$$\eta_{ik} \sim N(0, s_2^2)$$

"multiplicative noise"

Least trimmed sum of squares regression



minimize

$$\sum_{i=1}^{n/2} (y_{(i)} - f(x_{(i)}))^2$$

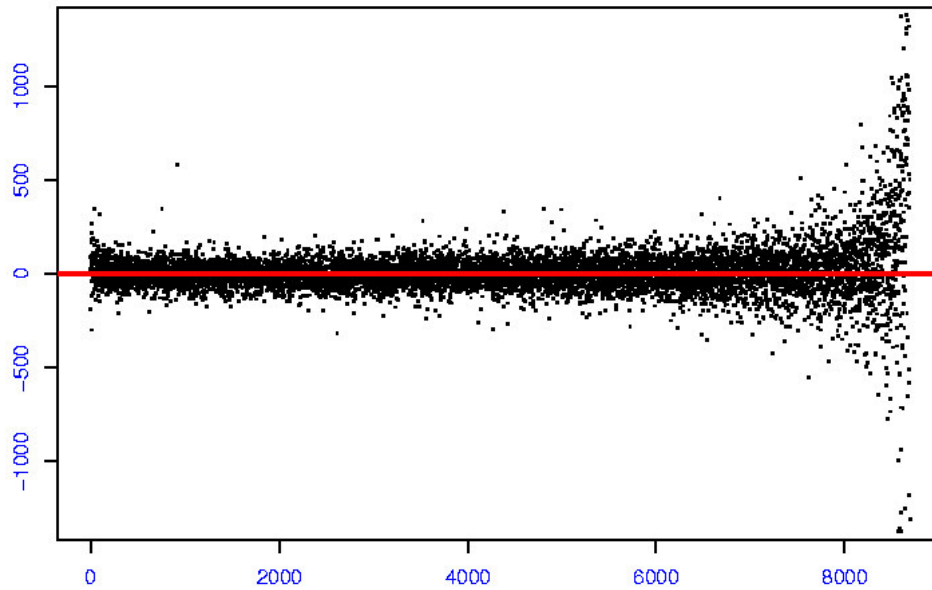
P. Rousseeuw, 1980s

- least sum of squares

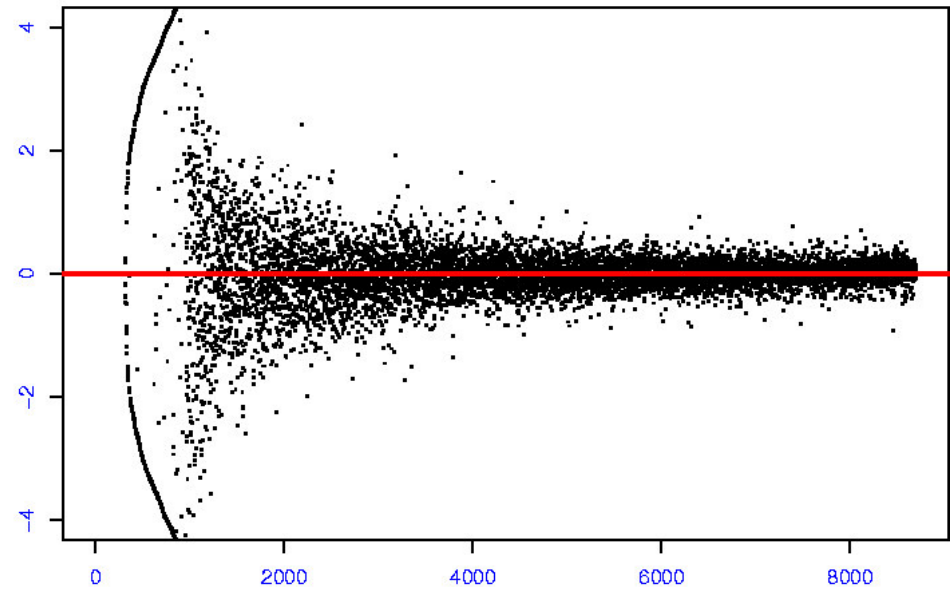
- least trimmed sum of squares

evaluation: effects of different data transformations

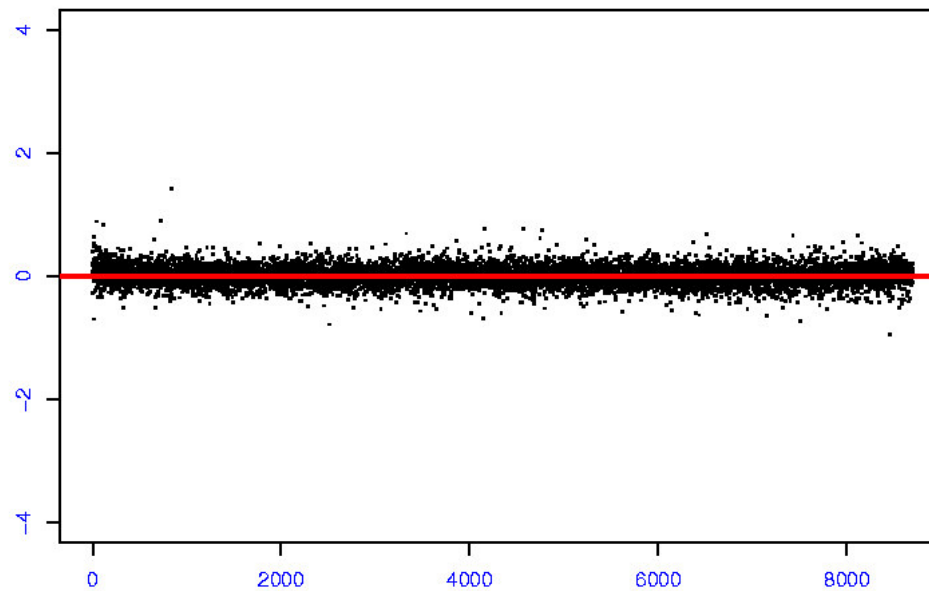
a) Δy



b) $\Delta \log(y)$

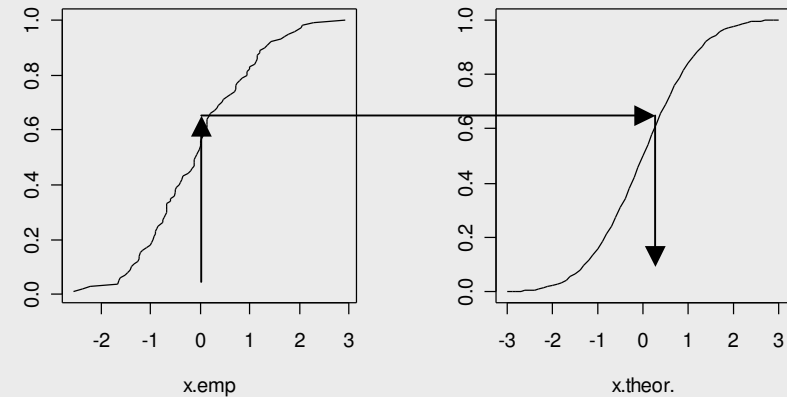
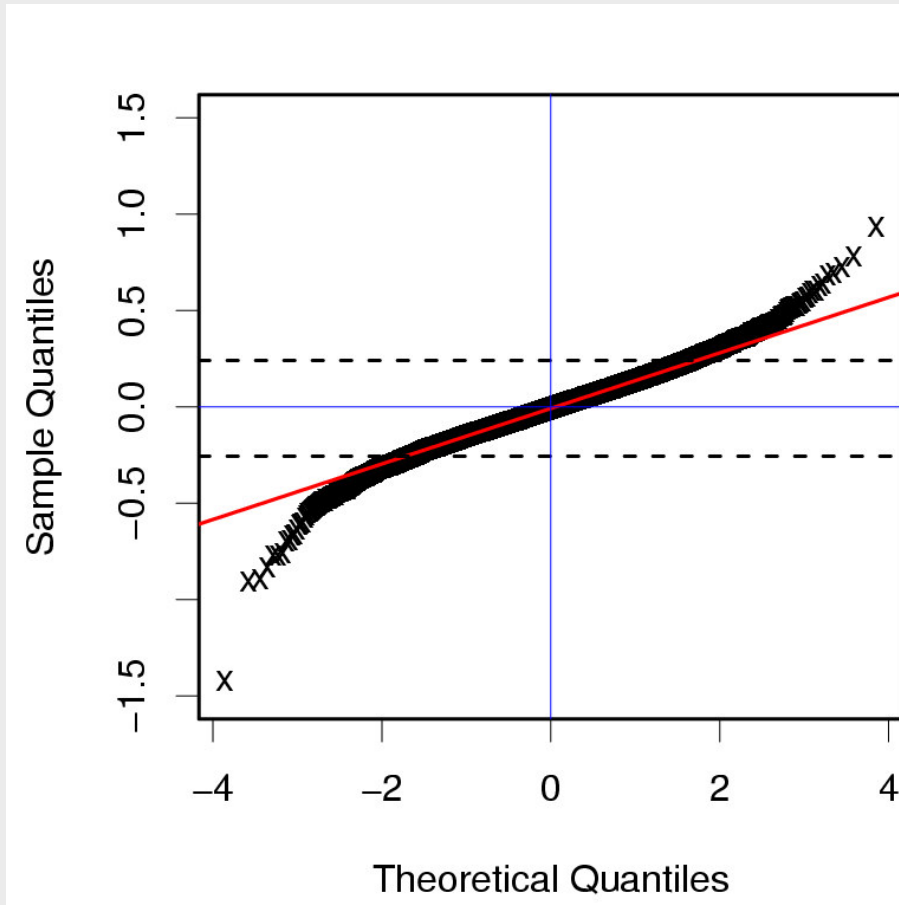


c) $\Delta h(y)$



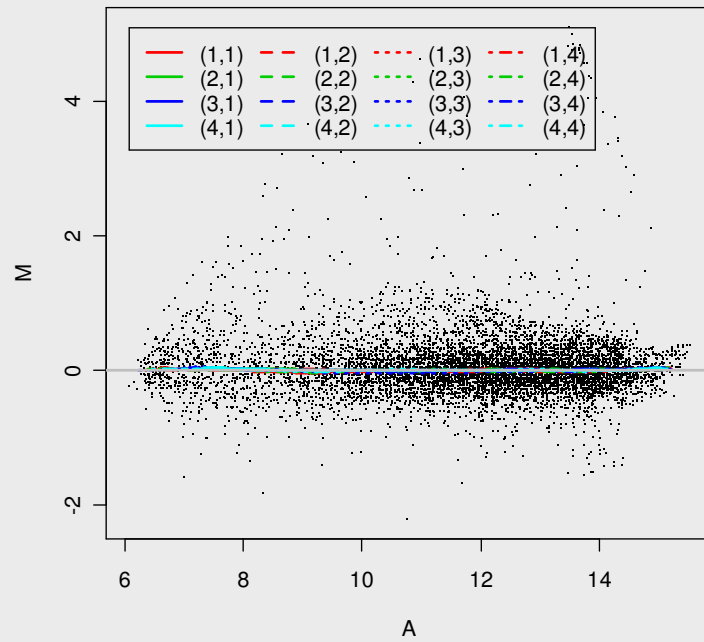
difference red-green
↑
rank(average) →

► Normality: QQ-plot

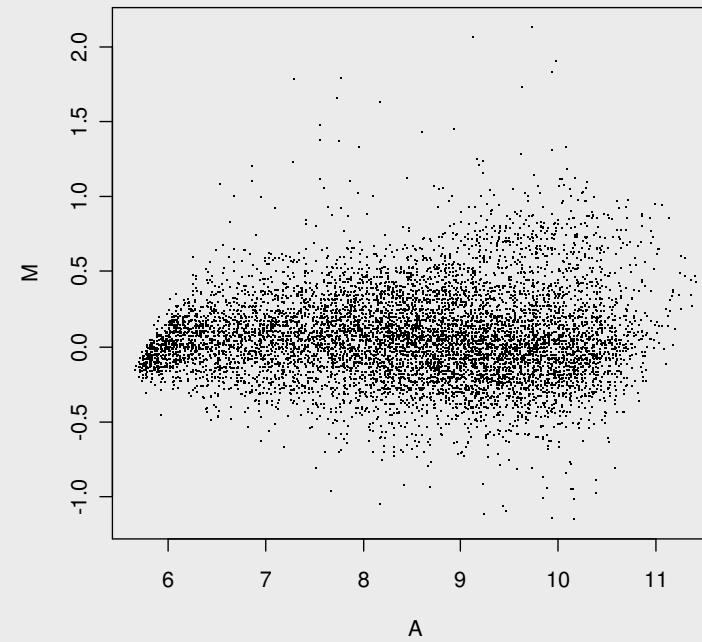


Swirl - VSN

Swirl array 93: post-norm MA-Plot

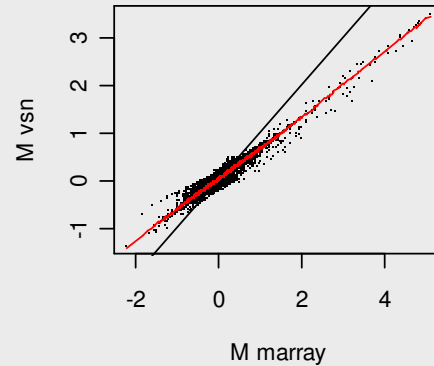


Swirl Array 93: VSN MA-plot

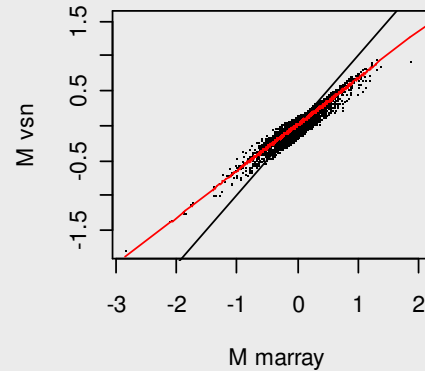


Swirl: marray versus VSN

swirl.1.spot



swirl.2.spot

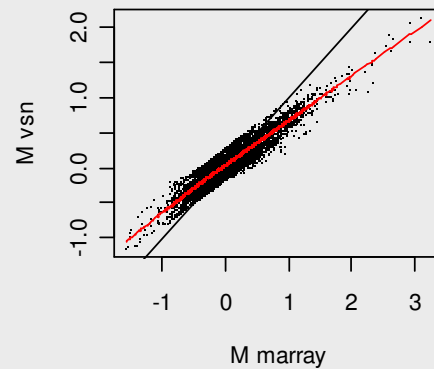


Why is the common slope different from 1?

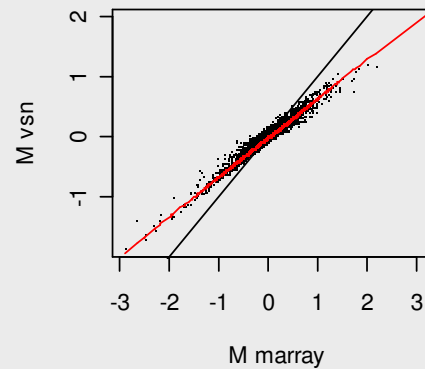
Scaled by a factor
~ 0.64

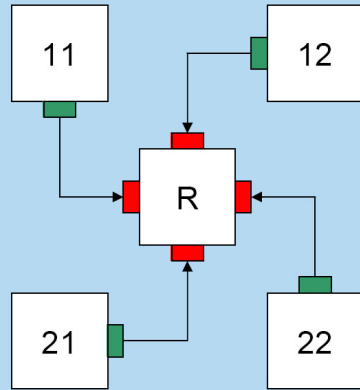
$\text{Var}(Y)/\text{Var}(X) \sim 0.65^2$

swirl.3.spot

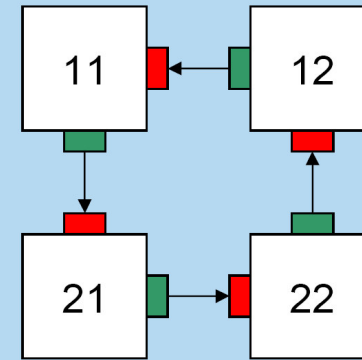


swirl.4.spot

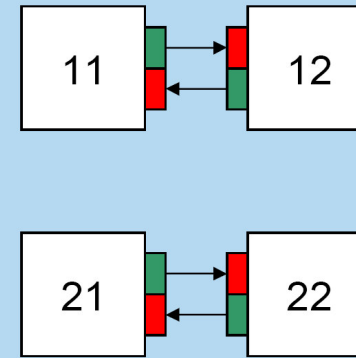
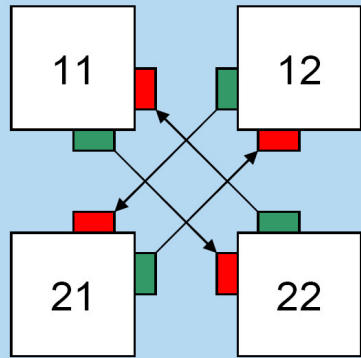




Cross-Swap (CS)



B-Swap (BS)



Summary

- What makes a good measurement: Precision and Unbiasedness
- Make you aware of the need to normalize.
- Normalization is not something trivial, has many practical and theoretical implications which need to be considered.
- What is the best way to normalize?
- How dependent is the result of your analysis from the normalization procedure?
- Parrish and Spencer [J Biopharm Statist 2004; 14: 575-589] considered the problem that normalization may modify the data so that both the technical and the biological variability is reduced.

However, the biological variation should not be changed. If it was reduced, the smaller experimental error variances critically affect significance testing and can elevate false discovery rates.

Parrish and Spencer demonstrate that the quantile normalization can have an impact on the biological variability.

Acknowledgments

Holger Sültmann
Wolfgang Huber

Simon RM et al. (2003) *Design and Analysis of DNA
Microarray Investigations*, Springer New York

Yang YH, Dudoit S. (2004) *Normalization: Bioconductor's
marray package*