Genexpressionsanalyse

A cell and its population of genes:



DNA forms double strands by a process called hybridization:



Labeling









placenta complex

2 color fluoresent dye labeled mRNA



Affymetrix oligonucleotide arrays





Illumina bead arrays



3 micron diameter beads are coated with Oligonucleotides

IL6

L12

IFR1







E E

Beads embedded in a slide with multiple arrays located on each slide

Pooled Beads are randomly located and assayed to identify the location of each bead within the array using a 29 base tag sequence







Beads located in the end of a microfibre and collected into a bundle each bundle in a 96 sample/array matrix format



• <u>Illumina flyer</u>

SNP detection

(slide from NCBI,

http://www.ncbi.nlm.nih.gov/projects/genome/probe/doc/TechBeadArray.shtml



T Allele-specific oligo (ASO) 1
C Allele-specific oligo (ASO) 2
Primer 3
Locus-specific oligo (LSO)
Paramagnetic particles
Three assay oligonucleotides are designed for
each SNP locus. Two oligos are specific to

each SNP locus. Two oligos are specific to each allele of the SNP site (Allele-Specific Oligos or ASOs). A third oligo that hybridizes between 1 and 20 bases downstream from the ASO site is the Locus-Specific Oligo (LSO). All three oligos contain regions of genomic complementarity and universal PCR primer sites; the LSO also contains a unique address sequence complementary to a particular bead type.

After hybridization, a polymerase fills the gap between the ASO and LSO. A DNA ligase seals the nick between the extended sequence of the ASO and the LSO to form PCR templates that can be amplified with universal PCR primers.

Universal PCR primers P1 and P2 are labeled with Cy3 and Cy5 dyes, respectively. After thermal cycling and downstream processing the single-stranded, dye-labeled PCR products are hybridized to their complement bead type through their unique address sequences. After the hybridization, the BeadArray Reader is used to analyze fluorescence signal.



RNA-seq





Estimation of gene expression from RNA-seq: **RPKM values**

Number of reads which map per kilobase of exon per million mapped reads for each gene



Two samples from the same kidney carcinoma



Cancer – normal comparison



MA-plot: Minus vs. Average

$$M = \log_2(R/G) = \log_2(R) - \log_2(G)$$
$$A = \frac{1}{2}\log_2(RG) = \frac{1}{2}(\log_2(R) + \log_2(G))$$

Pre-Norm Dilutions Dataset (array 20B v 10A)





Normalization of microarray data

Anja von Heydebreck

Dept. Computational Molecular Biology, MPI for Molecular Genetics, Berlin

Systematic differences between arrays

The boxplots show distributions of logratios from 4 redgreen 8448-clone cDNA arrays hybridised with \ge zebrafish samples.

Some are not centered at 0, and they are different from each other.



Experimental variation

amount of RNA in the sample efficiencies of -RNA extraction -reverse transcription

- -labeling
- -photodetection

Systematic

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

Normalization

Normalization:

Correction of systematic effects arising from variations in the experimental process

Ad-hoc normalization procedures

- 2-color cDNA-arrays: multiply all intensities of one channel with a constant such that the median of log-ratios is 0 (equivalent: shift log-ratios). Underlying assumption: equally many up- and downregulated genes.
- One-color arrays (Affy, radioactive): multiply intensities from each array *k* with a constant *c*_{*k*}, such that some measure of location of the intensity distributions is the same for all arrays (e.g. the trimmed mean (Affy *global scaling*)).

log-log plot of intensities from the two channels of a microarray

comparison of kidney cancer with normal kidney tissue, cDNA microarray with 8704 spots

- red line: median normalization
- blue lines: two-fold change



Assumptions for normalization

- When we normalize based on the observed data, we assume that the majority of genes are unchanged, or that there is symmetry between up- and downregulation.
- In some cases, this may not be true. Alternative: use housekeeping genes, which supposedely don't change, or (spiked) controls, and base normalization on them.

1. Loess normalization

M-A plot (minus vs. add): $\log(R) - \log(G) = \log(R/G)$ VS. $\log(R) + \log(G) = \log(RG)$

With 2-color-cDNA arrays. often "banana-shaped" scatterplots on the logscale are observed.



Loess normalization

≥ 2 10 12 14 8

zebrafish data

• Intensity-dependent trends are modeled by_N a regression curve, $M = f(A) + \varepsilon$.

• The normalized $\begin{bmatrix} 2 \\ log-ratios are \\ computed \\ as the residuals <math>\varepsilon$ of the loess regression.



Loess regression

- Locally weighted regression.
- For each value x_i of X, a linear or polynomial regression function f_i for Y is fitted based on the data points close to x_i . They are weighted according to their distance to x_i .
- Local model: $Y = f_i(X) + \varepsilon$.
- Fit: Minimize the weighted sum of squares $\Sigma w_j(x_j)(y_j f_i(x_j))^2$
- Then, compute the overall regression as:

$$Y = f(X) + \varepsilon$$
, where $f(x_i) = f_i(x_i)$.

Loess regression





Print-tip normalization

- With spotted arrays, distributions of intensities or log-ratios may be different for spots spotted with different pins, or from different PCR plates.
- Normalize the data from each (e.g. printtip) group separately.



Print-tips correspond to localization of spots

Slide: 25x75 mm

Spot-to-spot: ca. 150-350 µm



4x4 or 8x4 sectors

17...38 rows and columns per sector

ca. 4600...46000 probes/array

sector: corresponds to one print-tip

Print-tip loess normalization



 \geq

Q-Q Plots

- [Wikipedia] A Q–Q plot is a plot, which is a graphical method for comparing two probability distributions by plotting their quantiles against each other.
- A point (x,y) on the plot corresponds to one of the quantiles of the second distribution (y-coordinate) plotted against the same quantile of the first distribution (xcoordinate)

Q-Q Plots, continued

- If the two distributions being compared are similar, the points in the Q–Q plot will approximately lie on the line *y* = *x*.
- Q–Q plots can be used to compare collections of data, or theoretical distributions. The use of Q–Q plots to compare two samples of data can be viewed as a non-parametric approach to comparing their underlying distributions.



Quantile Normalization

A comparison of normalization methods for high density oligonucleotide array data based on variance and bias

- Bolstad, Irizarry, Astrand, Speed
- Bioinformatics 2003, vol 19(2), 185-193

Quantile Normalization

	Hyb 1	Hyb 2	Hyb 3	
Gen 1				
Gen 2				
Gen 3				
Gen 4				

Quantile normalization: Algorithm

- 1. Given n arrays of length p, form X of dimension $p \times n$ where each array is a column;
- 2. sort each column of X to give Xsort;
- 3. take the means across rows of Xsort and assign this mean to each element in the row to get Xsort;
- 4. get Xnormalized by rearranging each column of Xsort to have the same ordering as original X

Questions

- Differential genes between two conditions (e.g., healthy/diseased)
- Coexpressed genes: pathways, coregulation
- Genes characteristic for particular conditions (e.g., tumor staging)

Chromatin IP (ChIP chip)



Tuning the model

ChIP-chip data for yeast: bound/unbound binding ratios (R/G)

- ~ 200 transcription factors
- 32 transcription factor descriptions (*TRANSFAC*)



Drosophila: Eve-2 Promoter

