

Introduction into Practical Analysis of Microarrays

Benedikt Brors

Division of Theoretical Bioinformatics
German Cancer Research Center
Heidelberg

03/03/08

- 1 Biological motivation of microarray experiments
- 2 Finding differentially expressed genes
- 3 Multi-condition experiments
- 4 Exploratory data analysis
- 5 Scenario 4: time series
- 6 Classification
- 7 What's in this course?

- Will be introduced as needed in subsequent units
- Important for low-level analysis (normalization, quality assessment, ...)
- Just recall: *cDNA* versus *oligonucleotide* microarrays, *spotted* vs. *printed* vs. *in-situ synthesized* chips, *one-channel* vs. *two-channel* readout.
- Terminology: DNA fragment bound to chip surface will be called **probe**, soluble cDNA/cRNA will be called **target**

- You want to compare two conditions (control/treatment, disease/normal etc.) and find differentially expressed genes
- You want to compare more than two conditions (disease subgroups, several treatments, several strains, several knockouts), some of which may interact (control/treatment vs. strain1/strain2)
- You want to find groups that are not defined yet (novel disease subtypes)

- You want to compare two conditions (control/treatment, disease/normal etc.) and find differentially expressed genes
- You want to compare more than two conditions (disease subgroups, several treatments, several strains, several knockouts), some of which may interact (control/treatment vs. strain1/strain2)
- You want to find groups that are not defined yet (novel disease subtypes)

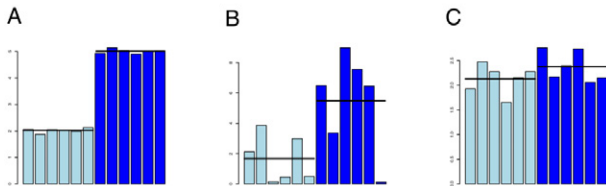
- You want to compare two conditions (control/treatment, disease/normal etc.) and find differentially expressed genes
- You want to compare more than two conditions (disease subgroups, several treatments, several strains, several knockouts), some of which may interact (control/treatment vs. strain1/strain2)
- You want to find groups that are not defined yet (novel disease subtypes)

- You want to investigate time series (developmental stages, transgene induction, cell cycle)
- You want to find predictive patterns for certain conditions (disease subtype markers, disease targets)
- You want to find patterns that are associated with prolonged patients' survival time
- You want to find patterns that tell you when a certain therapy will be of benefit

- You want to investigate time series (developmental stages, transgene induction, cell cycle)
- You want to find predictive patterns for certain conditions (disease subtype markers, disease targets)
- You want to find patterns that are associated with prolonged patients' survival time
- You want to find patterns that tell you when a certain therapy will be of benefit

- You want to investigate time series (developmental stages, transgene induction, cell cycle)
- You want to find predictive patterns for certain conditions (disease subtype markers, disease targets)
- You want to find patterns that are associated with prolonged patients' survival time
- You want to find patterns that tell you when a certain therapy will be of benefit

- You want to investigate time series (developmental stages, transgene induction, cell cycle)
- You want to find predictive patterns for certain conditions (disease subtype markers, disease targets)
- You want to find patterns that are associated with prolonged patients' survival time
- You want to find patterns that tell you when a certain therapy will be of benefit



- You want to find genes that display a large difference in gene expression *between* groups and are homogeneous *within* groups
- Typically, you would use statistical tests (t-test, Wilcoxon test)
- P values from these tests have to be corrected for multiple testing

- If there are more than two conditions, or if conditions are nested, the appropriate statistical method is ANOVA

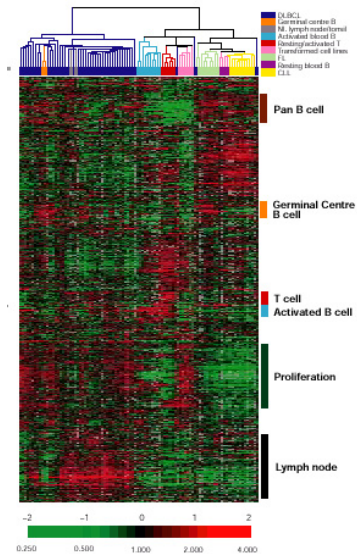


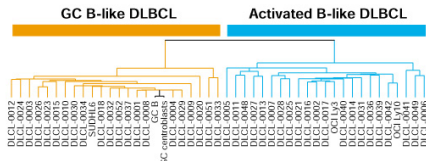
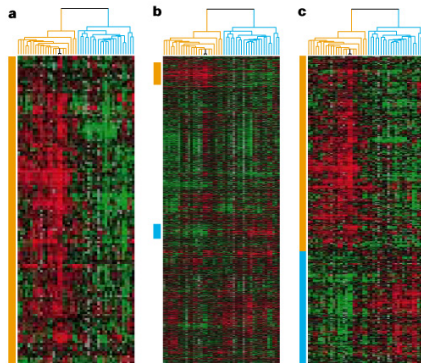
- The problem of multiple testing persists

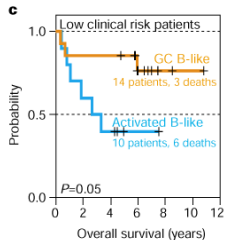
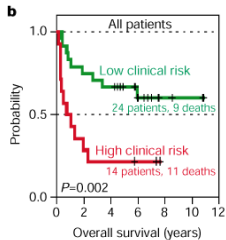
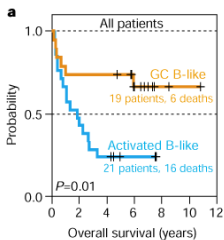
- Methods from this field were the first to be used for microarray data (*Eisenograms*)
- They should be used **only** if no prior knowledge exists that could be incorporated
- They will find patterns in your data, but any patterns, whether they are meaningful or not
- Methods include clustering (*hierarchical, partitioning*) and projection (*principal component analysis, multidimensional scaling*)

- Study was published in *Nature* **403**:503–511 (2000)
- Gene expression profiling of Diffuse Large B-Cell Lymphoma (DLBCL)
- Lymphoma is a blood cancer where *peripheral* blood cells degenerate and divide without control
- DLBCL is an aggressive form of this disease, originating from B-lymphocytes. Overall 5-year survival is about 40%.
- Current clinical risk factors are not sufficient.

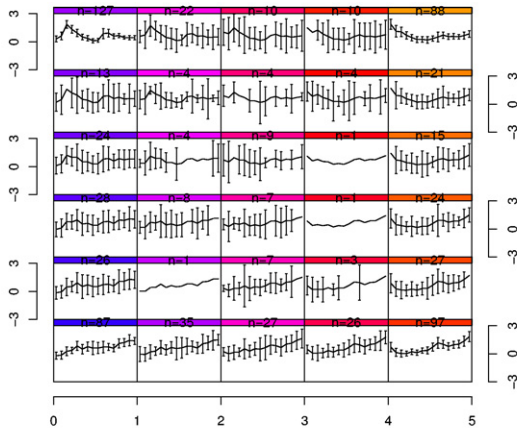
- A special cDNA chip was used, the *Lymphochip*
- spotted cDNA array of approximately 17,000 clones related to Lymphocytes
- 42 samples of DLBCL were analyzed, plus additional samples of normal B cells and of related diseases
- mRNA from these samples was competitively hybridized against control mRNA, stemming from a pool of lymphoma cell line mRNA preparations
- Data were analyzed by clustering





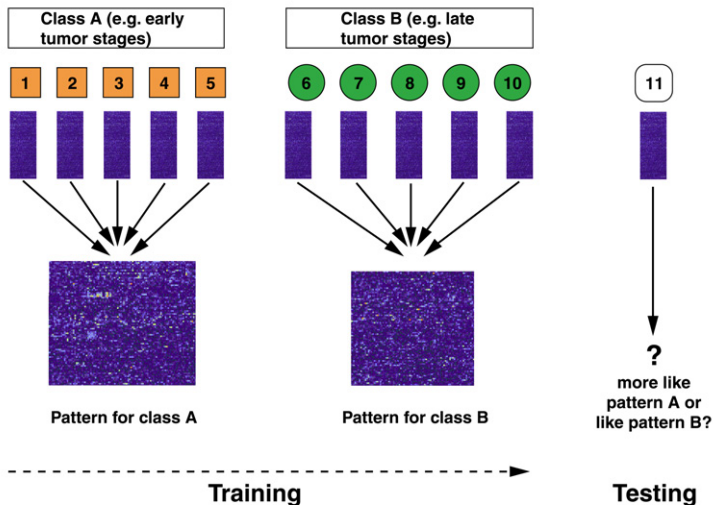


- In time series analysis, you usually want to find patterns of *coexpressed* genes, i.e. with coherent expression patterns
- The meaning of *time series* is different for biologists (2-10 time points) and statisticians (>200 time points)
- As a (non-optimal) solution, you would use clustering methods to find such patterns. Note that they are by no means exhaustive, and that no significance measure can be attached to them
- In contrast to EDA, *partitioning* cluster methods are more popular like **k-means** and **self-organizing maps**.



- If you seek genes whose expression profile is similar to that of a paradigmatic gene, you only need to calculate correlations, and sort by them. There is no need for clustering.
- Special methods exist for periodic changes (\Rightarrow cell cycle), e.g. Fourier analysis

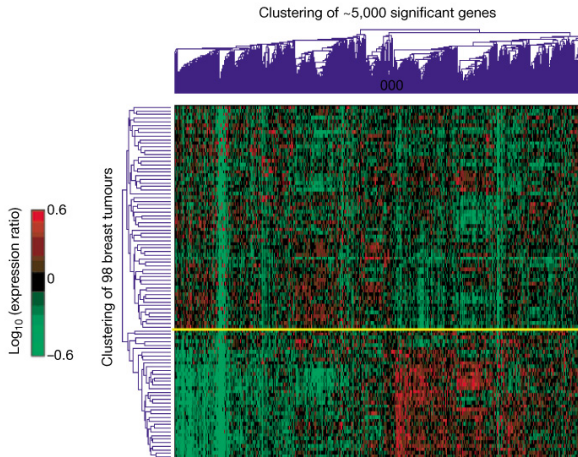
- If you have information about grouping of the samples, it can (and should) be used to get improved results.
- Groupings may be: Treatment/control, disease/normal, disease stage 1/2/3, mutant/wild type, good/poor outcome, therapy success/failure, and many more
- There may be more than two groups
- In Classification, you learn characteristic patterns from a *training set* and evaluate by predicting classes of a *test set*



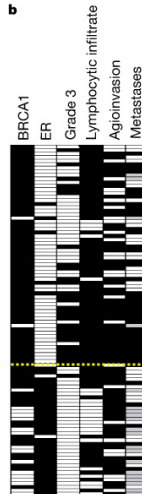
- published in *Nature* **415**:530–536 (2002)
- looks for prognostic markers in breast cancer
- two classes of patients: those with distant metastasis (other than in breast) within 5 years, and those without (also had negative lymph node status)
- In statistical thinking, this is a *classification* problem: given a set of *variables*, can we train a *classifier* such that it predicts for any new sample the *class* as correctly as possible?

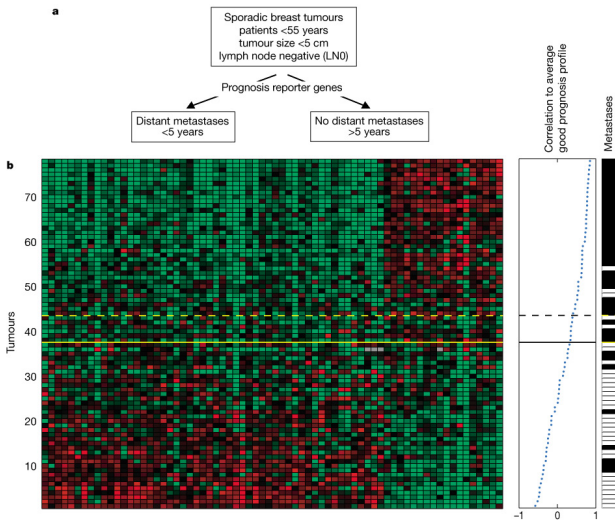
- A custom-made 25,000-clone chip was used; each feature contained a unique 60-mer oligonucleotide. This oligo was transferred to the chip by ink jet-like printing.
- The chips were hybridized competitively; the reference mRNA was obtained from a pool of patient mRNA (98 patients in total).
- Only data from certain genes (231) were used; finding out informative genes is called *feature selection* in machine learning.
- A home-made *ad hoc* classification method was used (no details given here). You can do better with established classification methods (tought later in this course).
- The model was validated by cross validation and by an independent test set.

a



b





Beware: re-analysis yields less optimistic results, cf. Tibshirani & Efron, Stat. Appl. Genet. Mol. Biol. 1:1 (2002).

- Instead of treating outcome as a binary variable (fatal/cured), you can use the *overall survival time* or the *event free survival time* as continuous variables, and try to estimate it by **regression**
- Since the risk to suffer from relapse is decreasing with time, linear regression models are almost always inappropriate
- Specialized models would be, e.g., *Cox regression*
- Regression trees can be used as well

- In pharmacogenomics, you try to find molecular predictors that tell you about probable success (or failure) of a certain therapy
- An example application would be estrogen receptor status for tamoxifen (antihormone) therapy or *HER2/NEU* status for herceptin therapy in breast cancer
- You may regard treatment outcome as a discrete variable and use classification methods, as described above
- Sometimes, it's convenient not to wait for the final endpoint (which may be years away), but to use *surrogate variables*, e.g. the drop of the blood level of a certain protein, or reduction in tumor volume

Finding differentially expressed genes, multiple testing

- Holger, Tue 9.00–10.30 am

Exploratory data analysis

- Tim, Tue 10.30am-12.00pm

Classification, molecular diagnostics

- Holger, Markus, Wed 9.00am–12.30pm

“Advanced topics”

- Achim, Manuela, Thu.



GK Smyth, YH Yang & T Speed (2003)

Statistical Issues in cDNA Microarray Data Analysis.

In: MJ Brownstein, AB Khodurski (eds.), Methods in Molecular Biology, Humana Press 2002.

<http://www.stat.berkeley.edu/~terry/zarray/TechReport/mareview.pdf>



W Huber, A von Heydebreck & M Vingron (2003)

Analysis of microarray gene expression data.

In: Handbook of Statistical Genetics, 2nd ed., Wiley 2003.

<http://www.ebi.ac.uk/huber/docs/hvhv.pdf>



R Gentleman, V Carey, W Huber, R Irizarry & S Dudoit (2005)

Bioinformatics and computational biology solutions using R and Bioconductor. Springer Publ. 2005.