
Group Testing: Global Tests, Holistic Approaches

Ulrich Mansmann, Manuela Hummel
IBE, Medical School, University of Munich

Content of the Lecture

Biological relevant information may rather be encoded in groups and not predominantly in the expression of single genes.

- **How to define gene groups:**
 - Exploratory research produces functional groups and genomic signatures: confirm the relevance of the specific group
 - Bioinformatic algorithms can be used to define pathways and functional groups
- **How to assess the relevance of groups of genes:**
 - Outstanding gene expression in a specific group compared to other genes
 - Differential gene expression not of single genes but over a specific group of genes
 - Relevance of specific gene group for biological phenomena

Gene Groups

- **Pathways**
networks of interacting genes (KEGG, cMAP, BioCarta)
- **Gene Ontology**
biological process, molecular function, cellular component
- **Regions in the genome**
- **Signatures for classification**
- **Groups defined by literature search**
- ...

Group Testing Strategies

- **Gene set enrichment approaches**

Gene-wise analysis for differential expression followed by over-representation analysis for gene groups

- **Holistic approaches**

Find gene groups that include differentially expressed genes

- **Special approaches for the Gene Ontology**

Find interesting categories while taking into account the special structure of the GO

Some Methods

- **Gene set enrichment approaches**
 - Fisher-test, hypergeometric test (*e.g. Draghici et al., 2003*)
 - Gene set enrichment analysis (GSEA) (*Subramanian et al., 2005*)
- **Holistic approaches**
 - globaltest (*Goeman et al., 2004*)
 - GlobalAncova (*Mansmann and Meister, 2005*)
 - Category approach (*Gentleman, 2006*)
- **Special approaches for the Gene Ontology**
 - Decorrelating the GO (*Alexa et al., 2005*)
 - Parent-child approach (*Grossmann et al., 2005*)
 - Focus-level approach (*Goeman et al., 2006*)

Fisher-test approach

In comparison with the whole population of genes, is the category of interest enriched with genes found to be differentially expressed?

	in category	not in category
differential		
not differential		

- Fisher-test
- Hypergeometric test

Gene Set Enrichment Analysis

Problem:

Is the gene expression in gene group A different from that in gene group B?
or: Is group A, compared to B, enriched with differentially expressed genes?

Important: Genes in both groups are different!

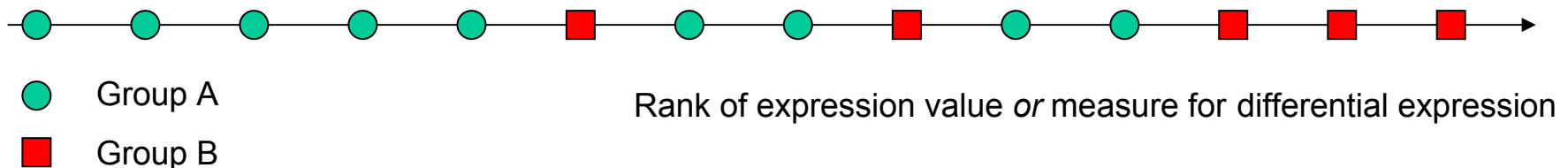
Basic idea:

Order genes with respect to expression value
or a measure for differential expression

Difference between groups → separation in the ordered list

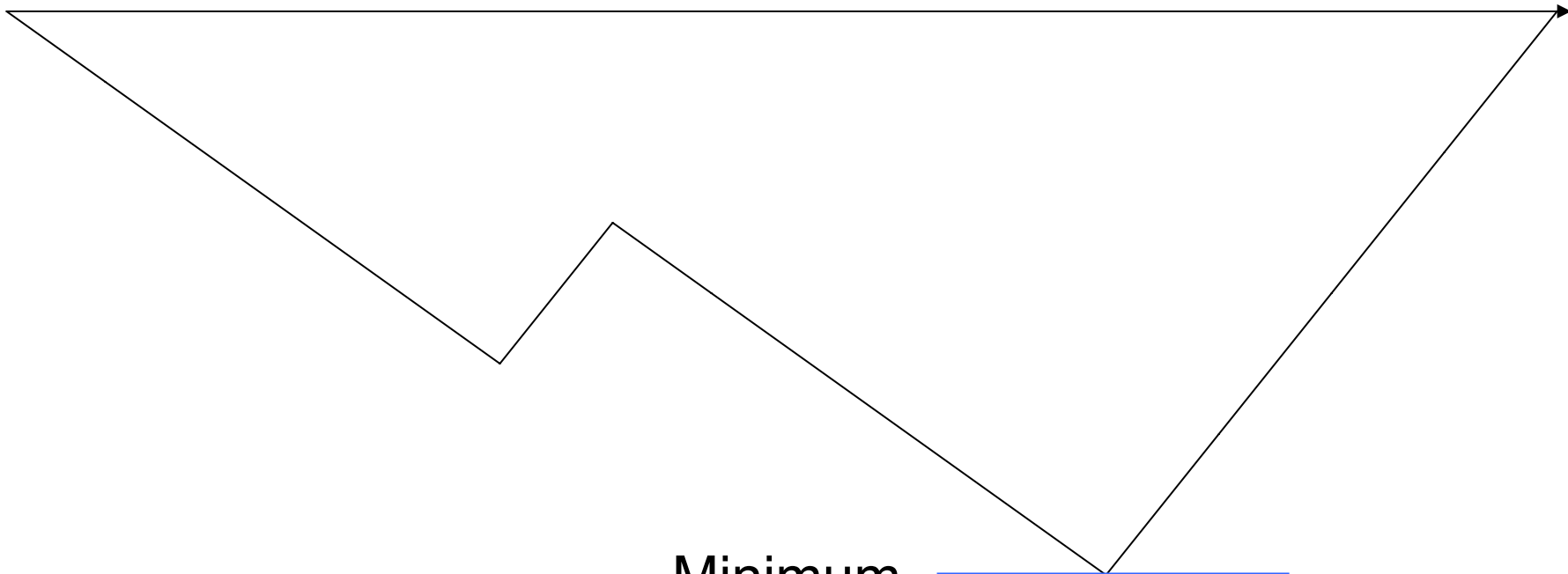
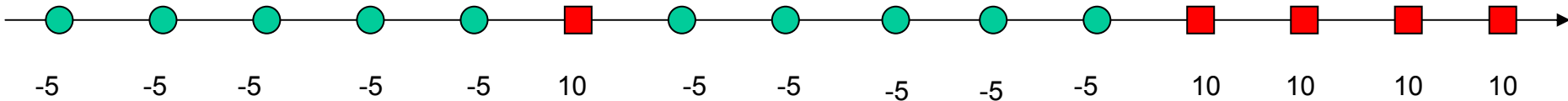
The Position of a gene in group A will have the tendency to be high or low

No difference → positions will be nicely mixed



Gene Set Enrichment Analysis

Genes ordered by rank of expression *or* by a measure for differential expression



- Group A ($n_A=10$)
- Group B ($n_B=5$)

Minimum

Is the minimum extreme with respect to random group mixing?

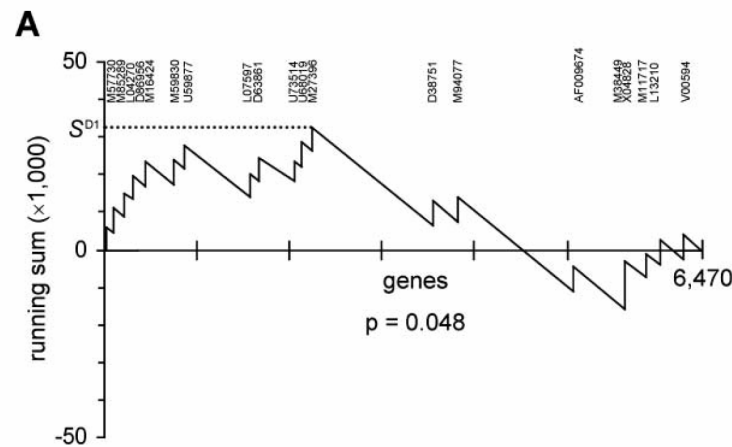
Group testing

The Algorithm Formalized

- n_A genes in group A, n_B genes in group B
- **Order genes** with respect to expression values *or* a measure for differential expression
- Create a **score vector** vv with value $-n_B$ at each position of a member from group A and with value n_A at each position of a member from group B
- Calculate the **cumulative sum** $yy = \text{cumsum}(vv)$
- Draw a line for yy starting and ending at 0
- Look at the **most extreme value** of the cumulative sum $M_{vv} = \max\{|\min(yy)|, \max(yy)\}$ which will be large in case of a good separation between both groups
- **Permutation test**
 - Permute the vector vv (= permute genes) to get vv^* , calculate yy^* and M_{vv^*}
Determine the permutation based p-value: $p_{\text{perm}} = \# \{M_{vv^*} \geq M_{vv}\} / \# \text{ permutations}$
 - **or**: Permute phenotype labels and compute vv^* , yy^* and M_{yy^*} for each permutation
- **Kolmogorov-Smirnov test**
Asymptotic test to determine whether the distribution of gene ranks is the same for group A and B

Example I: Cyclin D1 Action

- Lamb J et al. (2003) *A mechanism of Cyclin D1 Action Encoded in the Patterns of Gene Expression in Human Cancer*, Cell, 114: 323-334
- Cyclin D1 activity in Human Tumors: does the cyclin D1 target gene set play a prominent role in different tumor entities?
Being present as highly expressed genes
- Group A: Cyclin D1 expression signature: cyclin D1 target gene set.
Group B: all other genes



Example II: Colon Cancer

Study:

18 patients with UICC II colon cancer, 18 with UICC III colon cancer, snap-frozen material, laser microdissection, HG-U133A-arrays, 22.283 probesets representing ~18.000 genes

Question 1:

Are there specific cancer related pathways with a more distinct differential gene expression between UICC II/III?

Analysis:

Use Gene Set Enrichment approach described before and rank genes according to a measure for differential expression (absolute values of t-statistics)

Gene Set Enrichment – Colon Cancer

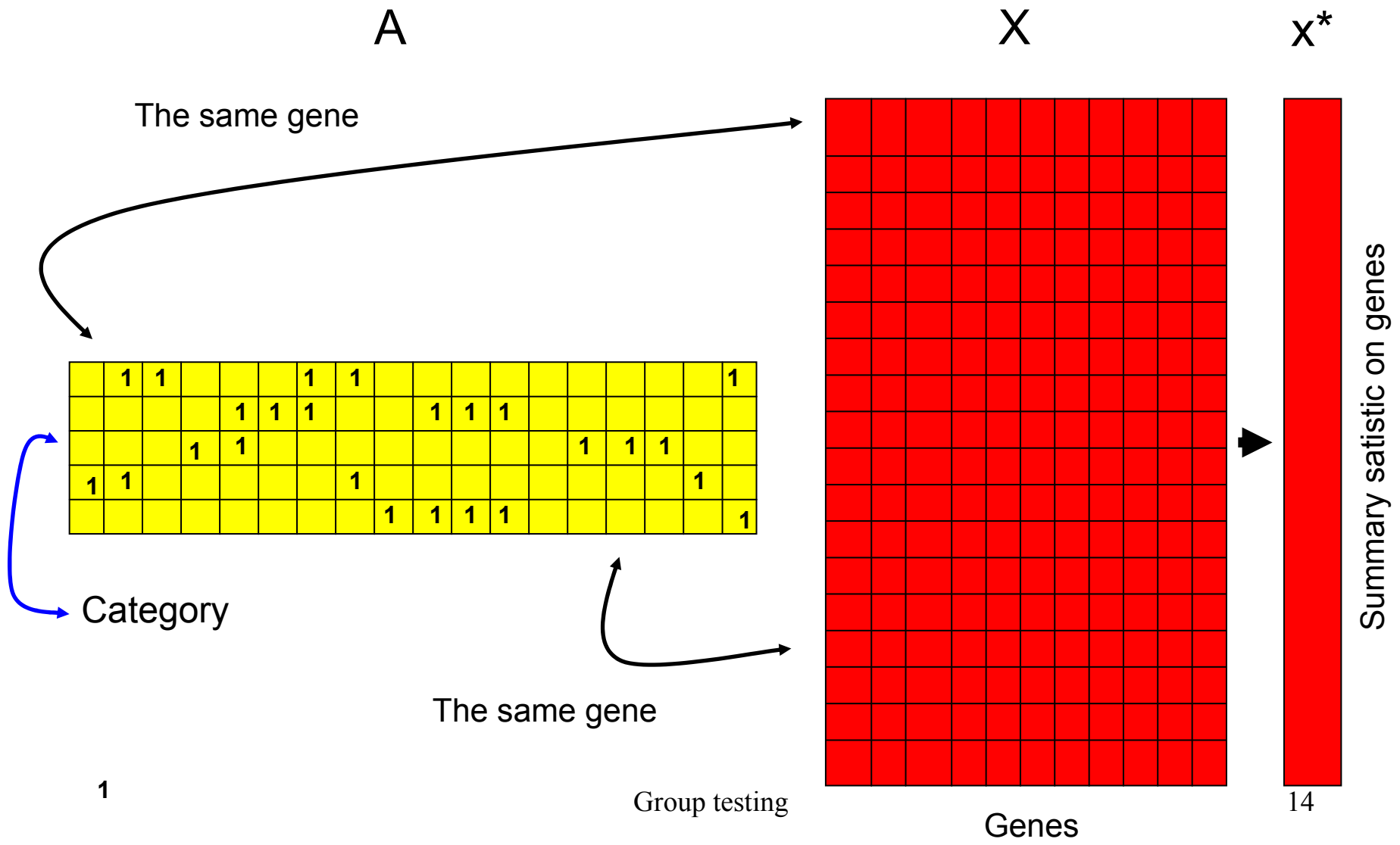
1407 probe sets are studied which belong to 9 cancer specific pathways

	group.A	group.B	M_{yy}	p.value
androgen_receptor_signaling	118	1289	6983	0.0568
Apoptosis	238	1169	17801	0.7438
cell_cycle_control	51	1356	10413	0.3616
notch_delta_signalling	50	1357	9010	0.6492
p53_signalling	45	1362	12390	0.0924
ras_signalling	311	1096	15486	0.6252
tgf_beta_signaling	100	1307	22615	0.0128
tight_junction_signaling	406	1001	15456	0.4414
wnt_signaling	214	1193	16318	0.8432

Gentleman's Categories

- A set of categories is merely a grouping of genes
- The groups do not need to be exhaustive or disjoint
- The mapping from a set of genes to a set of categories can be presented by an **incidence matrix A** (CxG)
 - C: Number of categories
 - G: Number of genes
- The elements of A:
 - $A[i,j] = 1$ if gene j is in category i , else 0**
- Row sums: Numbers of genes in each category
- Column sums: Number of categories a gene is in

Gentleman's Categories



Gentleman's Categories

- $z = A \cdot X$ or $z = A \cdot x^*$
- x^* could be the vector of **gene-wise t-statistics** between two phenotypes
- z is a vector of length C , represents **per category sum**, we are **interested in large or small z 's**
- H_0 : no difference between the expression means of the clinical groups
- Components of x^* are approximately $N(0,1)$
- The elements of $z = A \cdot x^*$ are sums of $N(0,1)$ [unfortunately not independent summands]
- **Permutation test**
 - Permute the columns of A . This is the same as permuting the gene labels (the labels or rows of X and x^*)
 - **or**: Permute the sample labels (the columns of X) and recompute x^*

Gentleman's Categories

- **Comparisons:**

within category comparison: For a given category, is the observed test statistic unusual?

overall comparison: Are any of the observed category statistics unusually large or small with respect to the entire reference distribution?

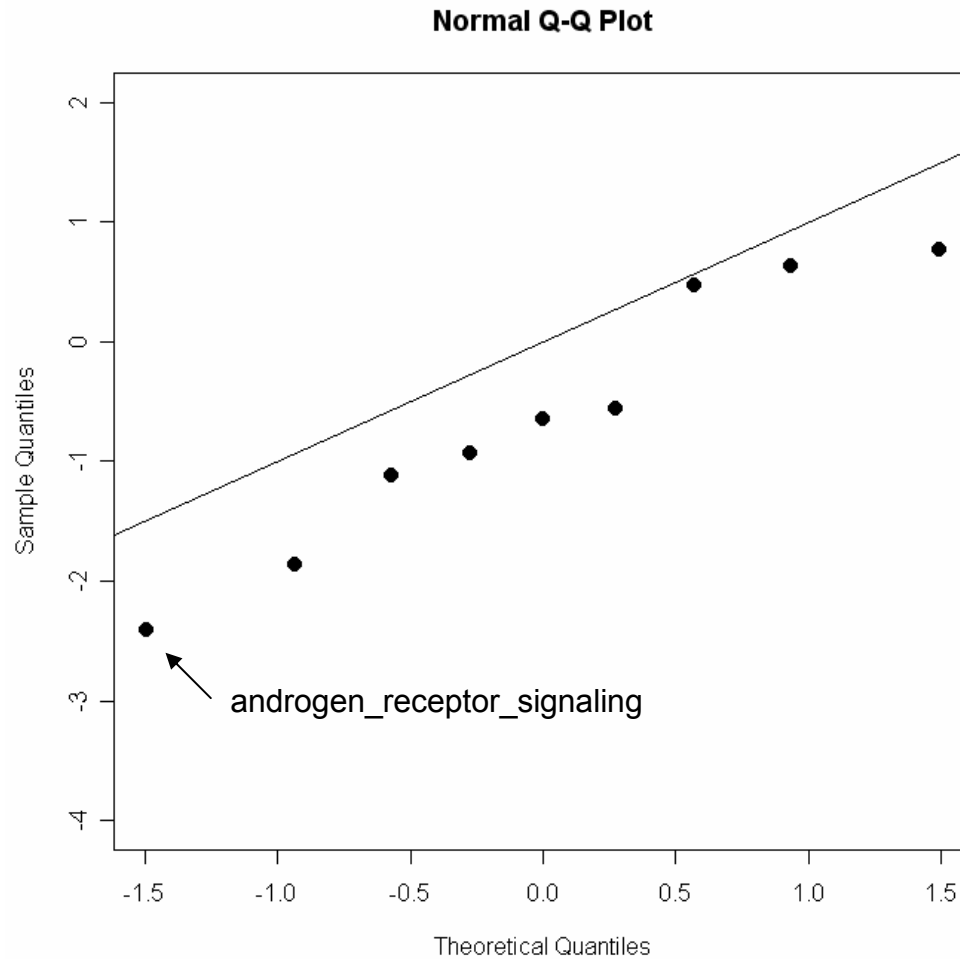
- **Note:** The approach is inherently multivariate, one data set gives G test statistics and these are transformed to yield C z_i 's
- The approach is well suited to fit the reasoning in a proper statistical framework
- In principle the approach can be used with arbitrary summary statistics x^* and arbitrary functions to yield z

Gentleman's Categories – Colon Cancer

Results for the colon data:

	z.statistic	p.Lower	p.Upper
androgen_receptor_signaling	-2.4124695	0.082	0.918
apoptosis	0.6270588	0.598	0.402
cell_cycle_control	0.7682091	0.690	0.310
notch_delta_signalling	-0.6442985	0.393	0.607
p53_signalling	-0.9325874	0.261	0.739
ras_signalling	0.4736045	0.586	0.414
tgf_beta_signaling	-1.1235767	0.376	0.624
tight_junction_signaling	-0.5652049	0.462	0.538
wnt_signaling	-1.8580599	0.288	0.712

Gentleman's categories – Colon Cancer



Goeman's Global Test

- Tests if global expression pattern of a group of genes is significantly related to some outcome of interest (groups, continuous phenotype)
- If this relationship exists, then the **knowledge of gene expression X** helps to improve the **prediction of the phenotype Y** of interest.
If the prediction can not be improved by knowing the gene expression then there will not be differential gene expression.
- Test statistic:
 - $Q \sim (Y-\mu)'R (Y-\mu)$
 - $\sim \sum [X_i'(Y-\mu)]^2$ sum over genes in the group
 - $\sim \sum \sum R_{ij}(Y_i-\mu) (Y_j-\mu)$ sum over subjects

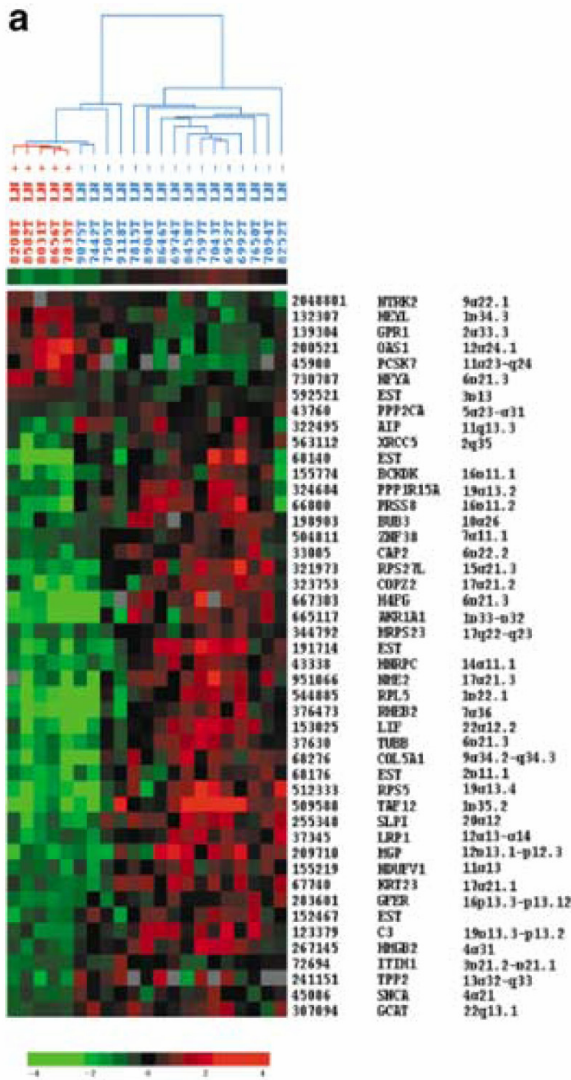
μ : Mean of phenotype,

X_{mi} Expression for gene m in subject i

R : X'X matrix of correlations between gene expression of subjects

Goeman JJ. Et al. (2003) *A global test for groups of genes: Testing association with a clinical outcome*, Bioinformatics, 20:93-99; Bioconductor package: *globaltest*

Example III: Lymph Node Metastases



Bertucci F et al. (2004) *Gene expression profiling of colon cancer by DNA microarrays and correlation with histoclinical parameters*, *Oncogene* 23, 1377–1391

Bertucci et al. present a gene signature consisting of 46 genes which is claimed to be able to discriminate between LN- and LN+ colorectal cancer.

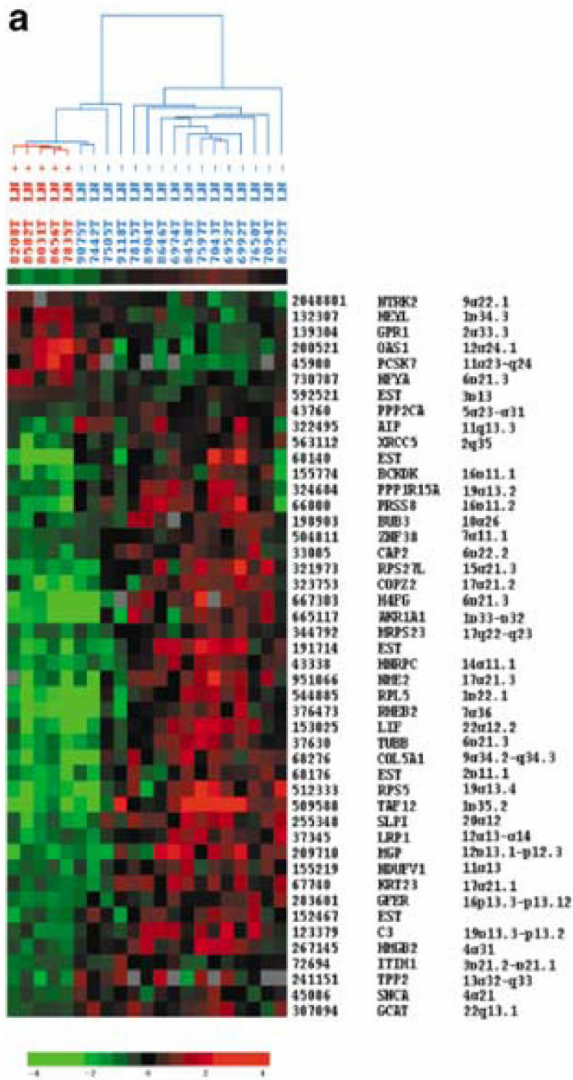
Is it possible to prove with new data that the signature has discriminative value? Can we reject the Null hypothesis

$$P[Y|X] = P[Y]$$

Y : LN+/LN-

X: expression pattern of 46 genes

Goeman's Global Test - Lymph Node Metastases



Global test is not significant ($p=0.43$)

No clear answer on the predictive power of the signature

No evidence for a difference is not evidence for no difference!

Question of power

Reasons for a non-significance: bad experiment or ...?

Example IV: Colon Cancer

Study:

18 patients with UICC II colon cancer, 18 with UICC III colon cancer, snap-frozen material, laser microdissection, HG-U133A-arrays, 22.283 probesets representing ~18.000 genes

Question 2:

Is there differential gene expression in the p53 signalling pathway between UICC II and UICC III colon cancer?

Analysis:

Goeman's global test

Goeman's Global Test – Colon Cancer

- Test for differential gene expression in *p53 signalling* pathway, 45 probesets

```
> gt <- globaltest(expressions,UICC.stage,genesets=pathways["p53_signalling"])  
> gt
```

Global Test result:

Data: 36 samples with 1407 genes; 1 pathway tested

Model: logistic

Method: Asymptotic distribution

	Genes	Tested	Statistic	Q	Expected	Q	sd of Q	P-value
p53_signalling	45	45	20.446		9.1621	4.3335	0.021357	

- **Informative plots**

Sample plot: how good fits a sample to its phenotype

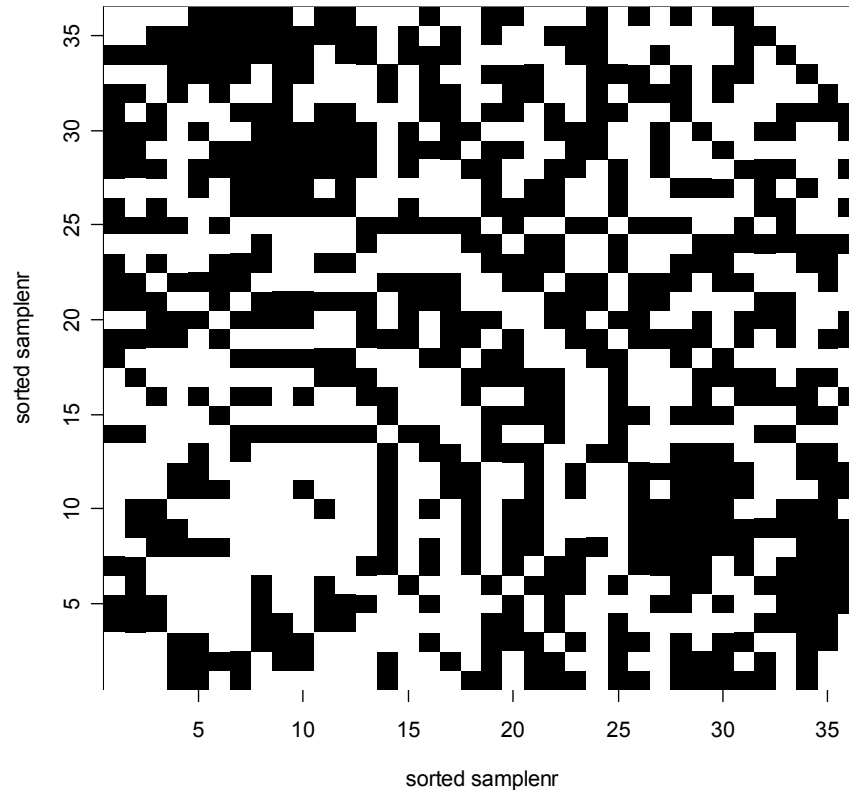
Checkerboard: correlation between samples

Gene plot: influence of single genes to test statistics

Goeman's Global Test – Colon Cancer

```
> checkerboard(gt)
```

Simultaneous correlation of phenotype and expression



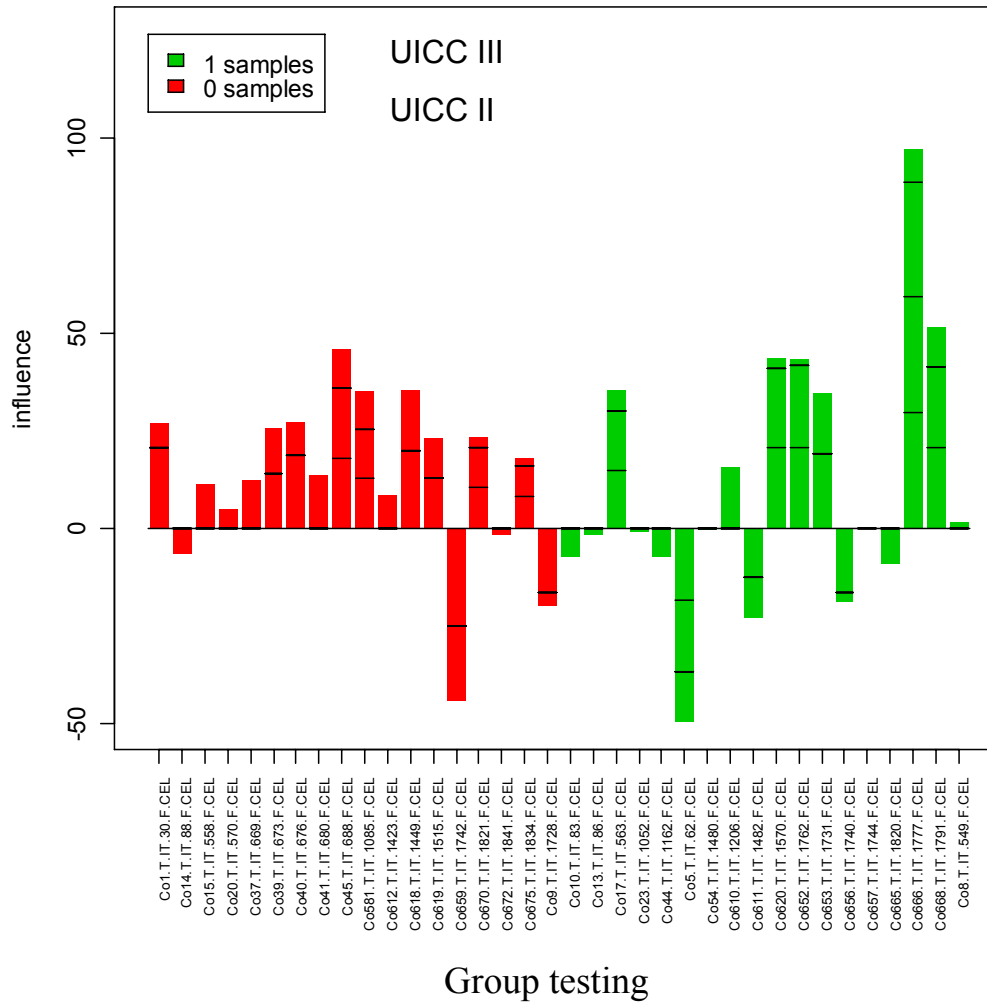
$$\sum \sum R_{ij} (Y_i - \mu) (Y_j - \mu)$$

Values dichotomized around median

Group testing

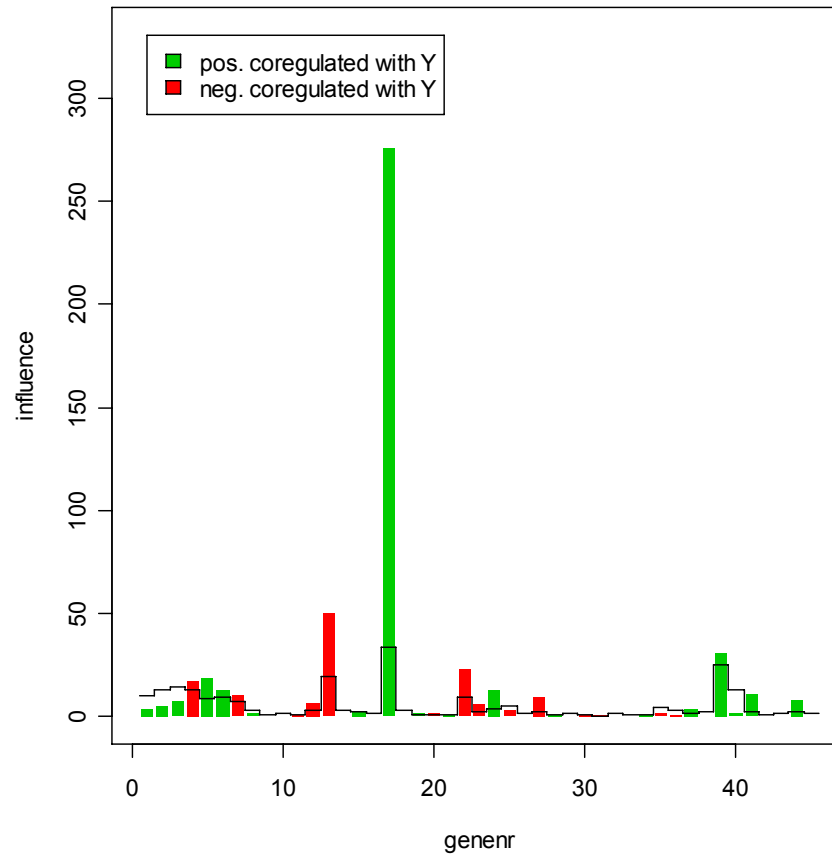
Goeman's Global Test – Colon Cancer

> sampleplot(gt)



Goeman's Global Test – Colon Cancer

```
> geneplot(gt)
```



$$\sum [X_i'(Y-\mu)]^2$$

GlobalANCOVA

- $P[X|Y]$ how is **gene expression X** influenced by **structure of phenotype Y**?
- General framework: p genes, n probes, $p \gg n$

$$X = \begin{pmatrix} x_{11} & \cdots & x_{1n} \\ \vdots & & \vdots \\ x_{p1} & \cdots & x_{pn} \end{pmatrix} = M + E \quad E[E] = 0 \quad M = \begin{pmatrix} \mu_{11} & \cdots & \mu_{1n} \\ \vdots & & \vdots \\ \mu_{p1} & \cdots & \mu_{pn} \end{pmatrix} = \begin{pmatrix} \mu^{(1)} \\ \vdots \\ \mu^{(p)} \end{pmatrix}$$

- The expectation for gene i and n probes follows a linear model $\mu^{(i)'} = D\beta^{(i)}$
- Example for D^t in the two group case with two covariates

Probe	P1	P2	P3	P4	P5	P6	P7	P8
Gene (i) specific mean	1	1	1	1	1	1	1	1
Group	0	0	0	0	1	1	1	1
Sex	M	M	F	F	F	F	M	M
Localization	C	R	C	R	C	R	C	R

GlobalANCOVA

Design	Full model	Reduced model
Groups	~ group + cov	~ cov
Dose – Response	~ dose + cov	~ cov
Group * Dose	~ group*dose + cov	~ group + dose + cov
time trends in groups	~ group*time + cov	~ group + time + cov
gene gene interaction	~ gene + cov	~ cov
differential co-expression	~ group*gene + cov	~ group + gene + cov
etc.		


GlobalANCOVA studies the following question:

Is the observed gene expression sufficiently explained by the reduced model or does the full model improve the model fit substantially?

GlobalANCOVA fits gene-wise linear models but summarizes the fit of all single models to a global measure

Example V

Prion infection of mice (Xiang et al.)



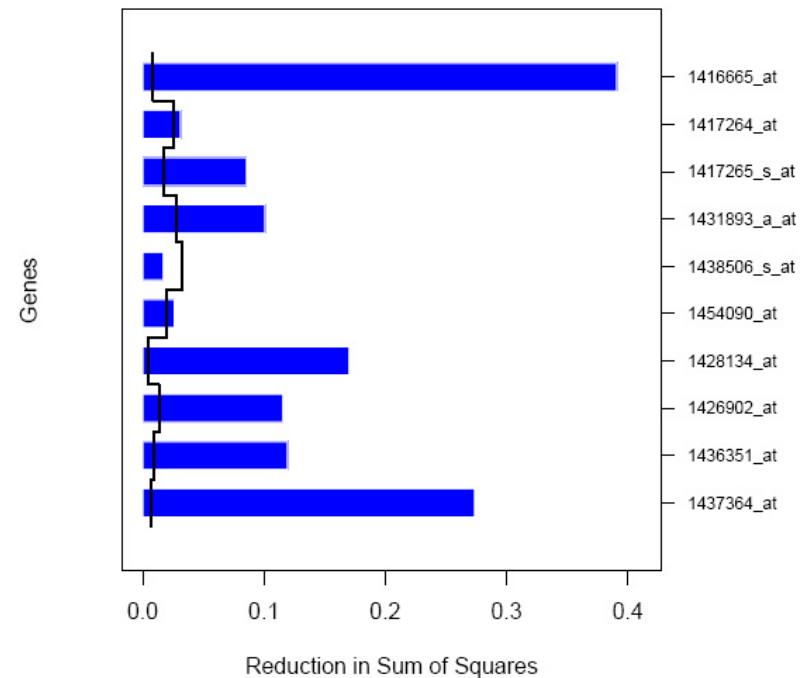
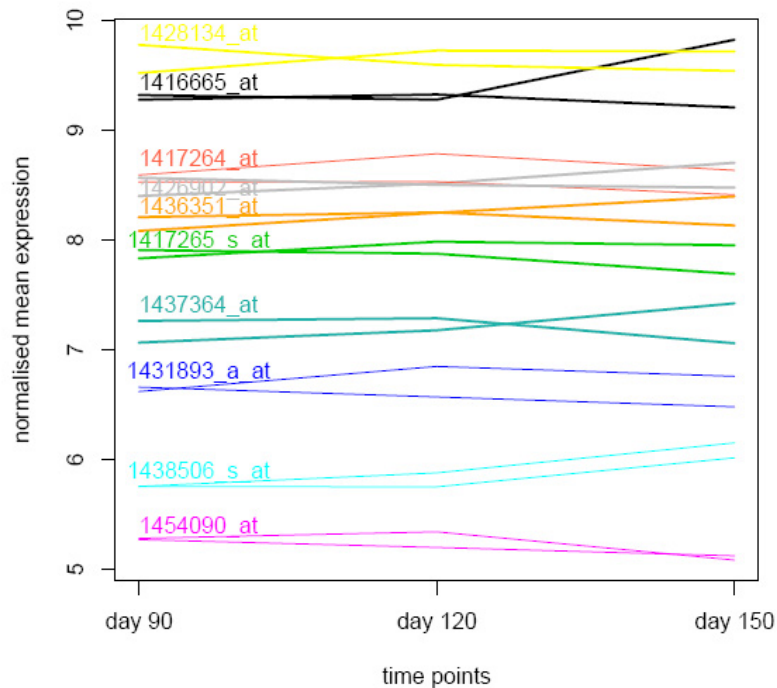
		days after inoculation		
		90	120	150
group	control	n = 3	n = 3	n = 3
	infected	n = 3	n = 3	n = 3

Is the time course of expression different for the two treatment groups in predefined biological processes?

→ Testing the interaction between time and treatment

Example V

```
> GlobalAncova(xx = expressions,  
               formula.full = ~ group * time,  
               formula.red   = ~ group + time,  
               model.dat     = phenodata,  
               test.genes    = ubiquinone.metabolism)
```



Summary: Two Perspectives on Gene Groups

Question 1: Two groups of genes have to be compared with respect to gene expression: Is the gene expression in gene group A different from the expression in gene group B?

Genes of group A

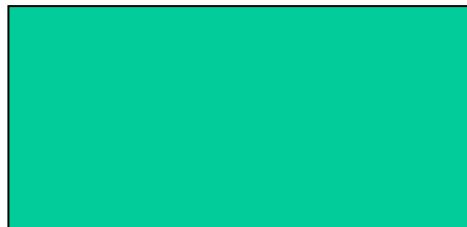
Genes of group B

Question 2: Is there differential gene expression between different biological entities not in terms of single genes but with respect to a defined group of genes?

Entity I

Entity II

Well defined
group of genes

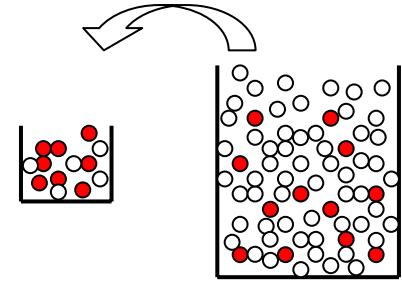


Group testing

Summary: Perspectives of Group Testing Methods

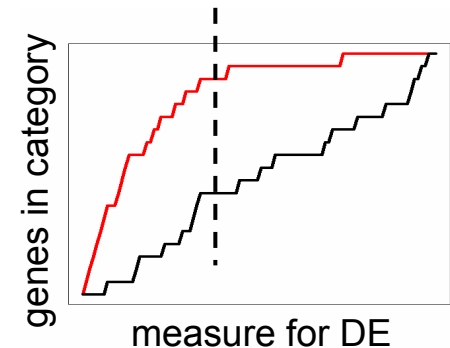
- **Fisher-test approach**

Are there more interesting genes in the category than expected by randomly drawing?



- **Gene set enrichment analysis**

Do the genes in the category have high ranks with respect to differential expression?



- **Global test / GlobalANCOVA / Category**

Can there be found differential expression in the category?



References

- Alexa A, Rahnenführer J, Lengauer T. Improved significance scoring of functional groups from gene expression data by decorrelating GO graph structure. *Bioinformatics* 2006; 22 (13): 1600-1607.
- Draghici S, Khatri P, Martins RP, Ostermeier GC, Krawetz SA. Global functional profiling of gene expression. *Genomics* 2003; 81: 98-104.
- R. Gentleman with contributions from S. Falcon (2006). Category: Category Analysis. R package version 2.0.0.
- Goeman JJ, van de Geer SA, de Kort F, van Houwelingen HC. A global test for groups of genes: testing association with a clinical outcome. *Bioinformatics*. 2004; 20: 93–99.
- Goeman JJ and Oosting J. (2006). Globaltest: testing association of a group of genes with a clinical variable. R package, version 4.4.0. (includes focus-level approach)
- Grossmann S, Bauer S, Robinson PN, Vingron M. An improved statistic for detecting over-representated Gene Ontology annotations in gene sets. *Research in Computational Molecular Biology: 10th Annual International Conference, RECOMB 2006, Venice, Italy, April 2-5, 2006. Proceedings: Lecture Notes in Computer Science 3909, Mar 2006, 85-98.*
- Mansmann U, Meister R. Testing differential gene expression in functional groups. *Methods Inf Med* 2005; 44(3).
- Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, Paulovich A, Pomeroy SL, Golub TR, Lander ES, Mesirov JP. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *PNAS* 2005; 102 (43): 15545-15550.