Differentially Coexpressed Genes Dennis Kostka & Rainer Spang

Courses in Practical DNA Microarray Analysis



Nationales Genomforschungsnetz

This is joint work with

- Dennis Kostka -



Differential Expression

Molecular disease mechanisms constitute abnormalities in the coregulation of genes

Alterations in gene regulation typically result in up and down regulation of genes



BUT

Not all changes in coregulation show up as patterns of up or down regulated genes

Differential Coexpression



A Regulatory Mechanism is Breaking Down



Questions



How can we find differential coexpression patterns ?

Do these pattern exist in real data ?

Are they biologically meaningful ?

Did we really need a new method to find them ?

How can we find differential coexpression patterns ?

How did we find differential expression patterns ? By screening one gene after the other

Problem:

Differential expression is a property of a single gene, differential coexpression is a property of a set of genes

- ... we need to screen all subsets of genes on the chip
- ... this is hard and can only be done heuristically

The problem of finding differential coexpression is mainly a problem of efficient search

A: Decide on a score for differential coexpression

B: Greedy stochastic downhill search

- 1. Choose a random set of genes and score it
- 2. Randomly select a neighboring set (no more than k different genes) and calculate its score
- 3. If the score of the new set is lower, change to the new set otherwise keep the old set
- 4. Iterate until you find a local minimum of the score

The computational costs for scoring a candidate set are critical for the practicality of the algorithm

The algorithm is stochastic. Restarting it several times can result in different local minima corresponding to different differential coregulation patterns

A score for differential coexpression of several genes



$$S'(I,J) = \frac{1}{(|I|-1)(|J|-1)} \sum_{I,J} (a_{ij} - a_{i\bullet} - a_{\bullet j} + a_{\bullet \bullet})^2$$

The trick (borrowed from Cheng and Church):

Calculating S is computationally expensive, but it is very cheap to decide whether adding or replacing genes leads to a higher or lower score, much cheaper than for the correlation coefficient

Neighborhood structure: Neighboring sets differ only by a single gene. Given a group of genes I we wish to exclude gene k:

$$S(I) \propto rac{\mathsf{mean}_{I,G1}(\mathsf{res})}{\mathsf{mean}_{I,G2}(\mathsf{res})} = rac{A_k^{(1)} + B_k^{(1)}}{A_k^{(2)} + B_k^{(2)}}$$

and modulo refitting of the parameters:

$$S(I \setminus k) < S(I)$$
 iff $B_k^{(1)} / B_k^{(2)} > S(I)$

Some Fine-tuning

We include / exclude a β -fraction of the genes that meet the criterion for a reduced score

To tune the size of the finally found gene sets we introduce a tuning parameter α .

The final criterion for including or excluding a gene now reads:

 $C_k(\alpha) = B_k^{(1)} / B_k^{(2)} \pm \{\alpha \cdot S(I) + (1 - \alpha) \cdot 1 / |I|\} > 0$

Do these pattern exist in real data ?

Acute Lymphoblastic Leukemia

- About 1/3 of all pediatric cancers
- Different cytogenetic risk groups (e.g. 70% overall cure rate vs. 30% in phil+)
- We compared cytogenetically normal children to those with the phil+ translocation

Yeoh EJ, RossMEet al. (2002) Classication, subtype discovery, and prediction of outcome in pediatric acute lymphoblastic leukemia by gene expression proling, Cancer Cell, **1(2)**, 133-43.

Differential coexpression in phil+ leukemia



Are the patterns biologically meaningful ?

Pattern 1 (34 genes)				
HINT1 (3094) RHEB (6009) KIAA1579 (55225) RNP24 (10959) GNAS (2778) ATP5J2 (9551) COX6A1 (1337)	PSMA2 (5683) ARHGDIB (397) SRP19 (6728) ECHS1 (1892) GNAS (2778) PTOV1 (53635) STARD7 (56910)	PSMB1 (5689) NCL (4691) RNF4 (6047) C9orf10 (23196) OTOR (56914) YWHAQ (10971) CALM2 (805)	FNTA (2339) EXO70 (23265 HCP15 (157317) REA (11331) COPE (11316) CBFB (865) PPP2R1A (5518)	CG018 (90634) KARS (3735) HNRPA2B1 (3181) UGP2 (7360) ATP5A1 (498) HLA-Z (267017)
Pattern 2 (21 genes)				
PSMB2 (5690) ACTB (60) CAPNS1 (826) PPP1CC (5501) GDI2 (2665)	UBC (7316) ENG (2022) PFN1 (5216) PSMC5 (5705)	ARHA (387) GG2-1 (25816) SF3B2 (10992) ARPC5 (10092)	PSMA4 (5685) ACTR3 (10096) GNAI2 (2771) MRLC2 (103910)	ARPC2 (10109) PITPNB (23760) ARHA (387) AES (166)

- Proteasome-Ubiquitin Pathway (for several cancers including CLL, inhibition can induce apoptosis)
- Involved in degradation of p27 (prognostic factor in B-cell lymphoma)

Most others: protein synthesis, protein transport, protein degradation

Did we really need a new method to find the patterns?

Screening for differential expression:

The genes in the two patterns have ranks between 106-6114

Hierarchical clustering:



Thank You