

Structured Analysis of Microarrays





Typical frustrations

Falsely predicted patients:

The predictive model seems to work very well for 80% of the patients and for the remaining 20% of patients you get wrong predictions although the pattern in the expression profiles seem to be quite obvious.

Genes that make no sense in this context:

If you look at the list of genes that drive the model the list does not tell you a unique biological story. You observe some genes that are expected to be there and many more genes that look like being randomly collected.

Implicit assumptions of standard approaches

The models aim to identify characteristic expression pattern for the whole patient groups (global molecular signatures) $% \left({\left[{{{\rm{T}}_{\rm{T}}} \right]_{\rm{T}}} \right)_{\rm{T}}} \right)$

The patient groups are seen as molecular homogenous groups This is often not the case

There might be many different molecular phenotypes with a poor treatment outcome

Genes are anonymous variables: x1,...xn

This ignores that for many genes we already know a lot about their function and the biological processes they are involved in

Our approach:

- 1. Sub-class finding instead of global class prediction
- 2. Use of functional annotations of genes

Molecular Symptoms













































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
- \ldots 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







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.







