

Mixture Model Estimation with Constraints:

Analysis of time-course gene expression with heterogeneous data

Ivan G. Costa Filho
Alexander Schliep

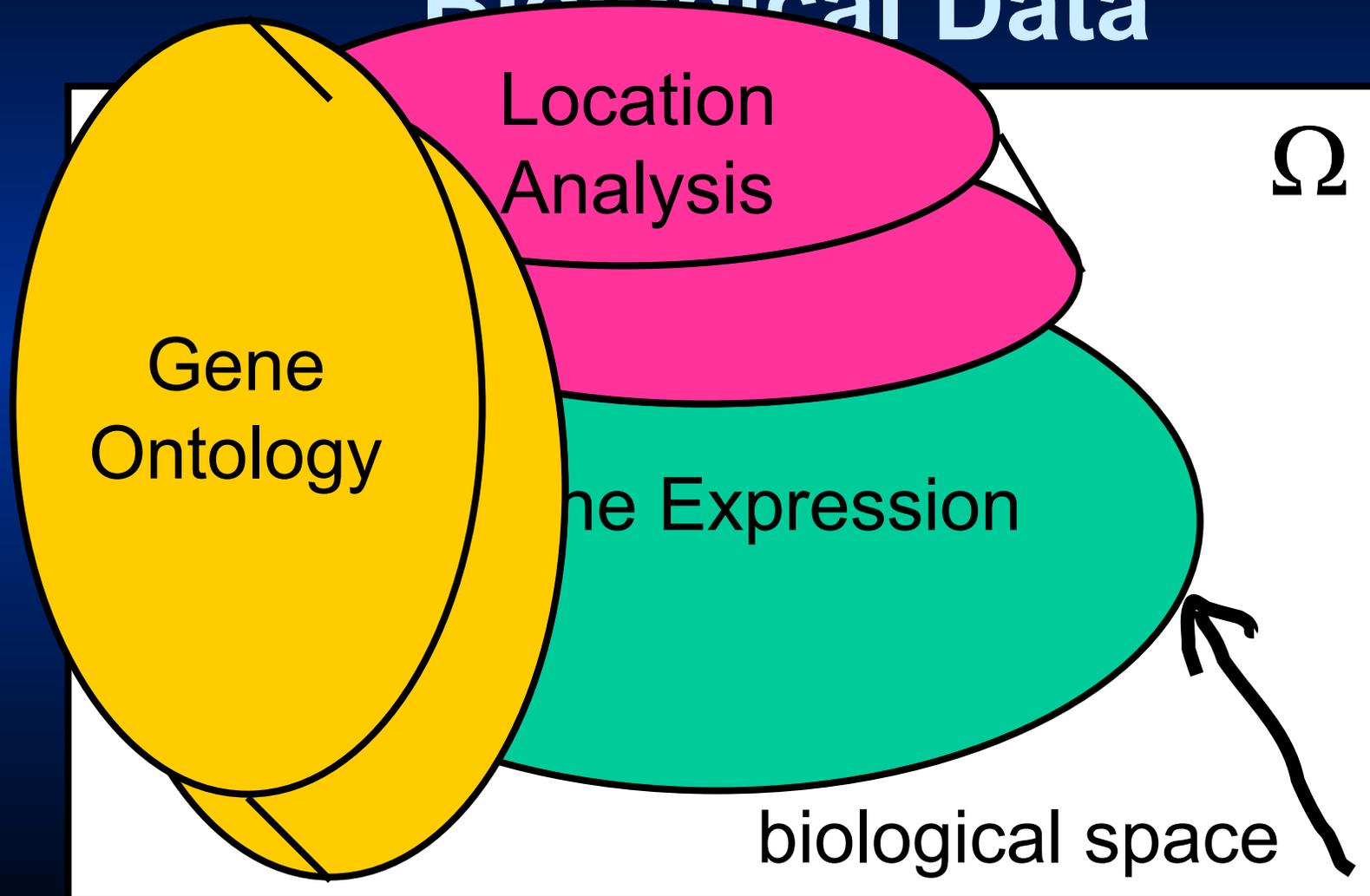


Computational Biology Department
Max-Planck-Institute for Molecular Genetics, Berlin

Motivation

Use additional **large scale biological data**
to **improve** clustering of gene
expression time-courses

Challenges of Heterogeneous Biological Data



Our approach

Semi-supervised

- Encode GO and location analysis as **soft pairwise constraints**
- Time-courses are modeled with a multivariate Gaussians
- Mixture estimation with constraints (Lange *et al.*, 2005, Lu and Leen, 2005)

Mixture Estimation with Constraints (1)

Maximize the complete likelihood:

$$P[X, Y|W, \Theta] = P[X|Y, \Theta] P[Y|W, \Theta]$$

where X is the observable data, Y the hidden data, Θ the model parameters, $W = \{W^+, W^-\}$ the pairwise constraints

The prior can be decomposed at:

$$P[Y | \Theta, W] = P[Y | \Theta] P[W^+ | Y, \Theta] P[W^- | Y, \Theta]$$

Mixture Estimation with Constraints (2)

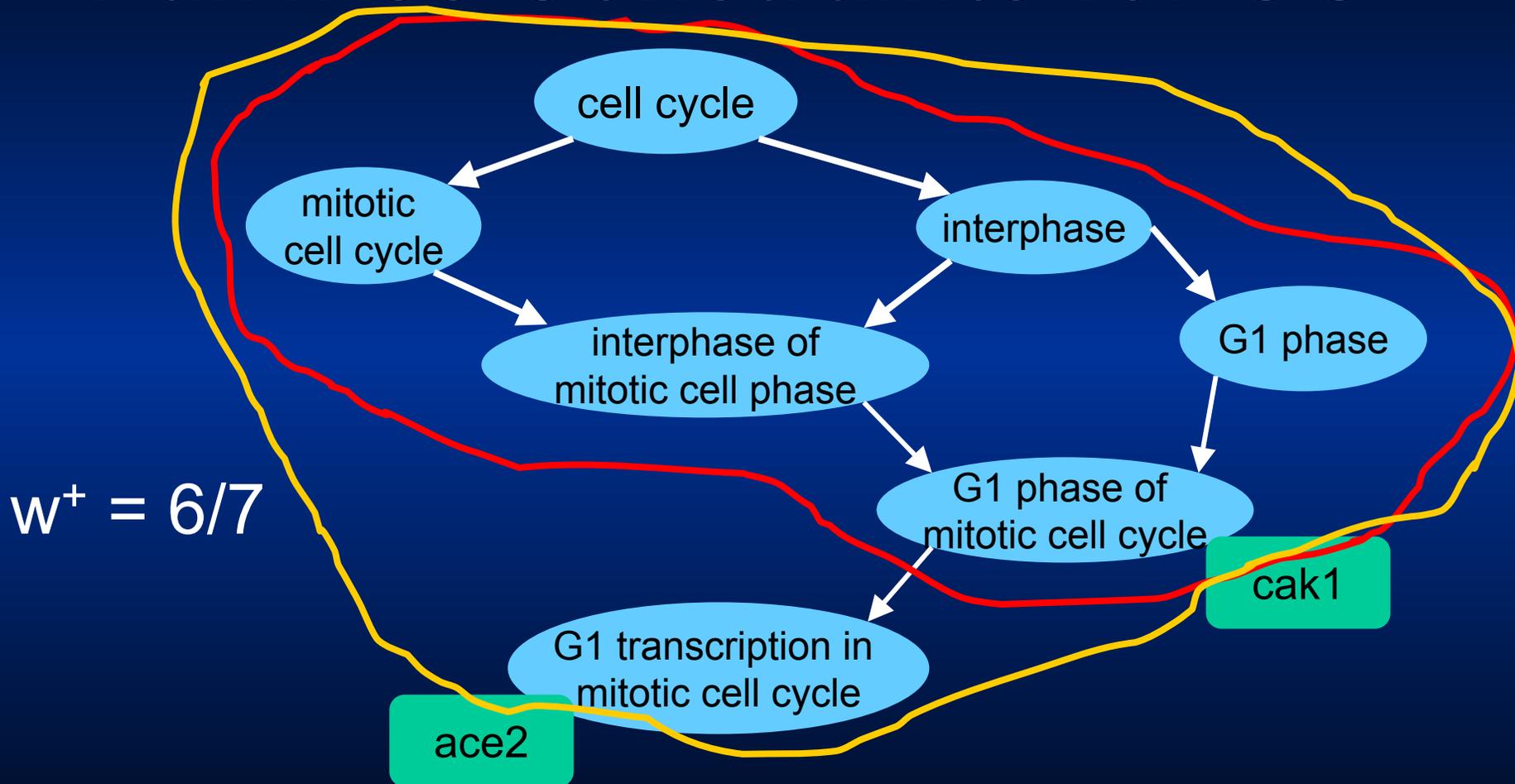
$$P[W^+ | Y, \Theta] \approx \exp \sum_i \sum_{j \neq i} -w_{ij}^+ 1\{y_i \neq y_j\} \lambda^+$$

The posterior assignments are approximated by means of Gibbs sampling (Lu and Leen, 2005)

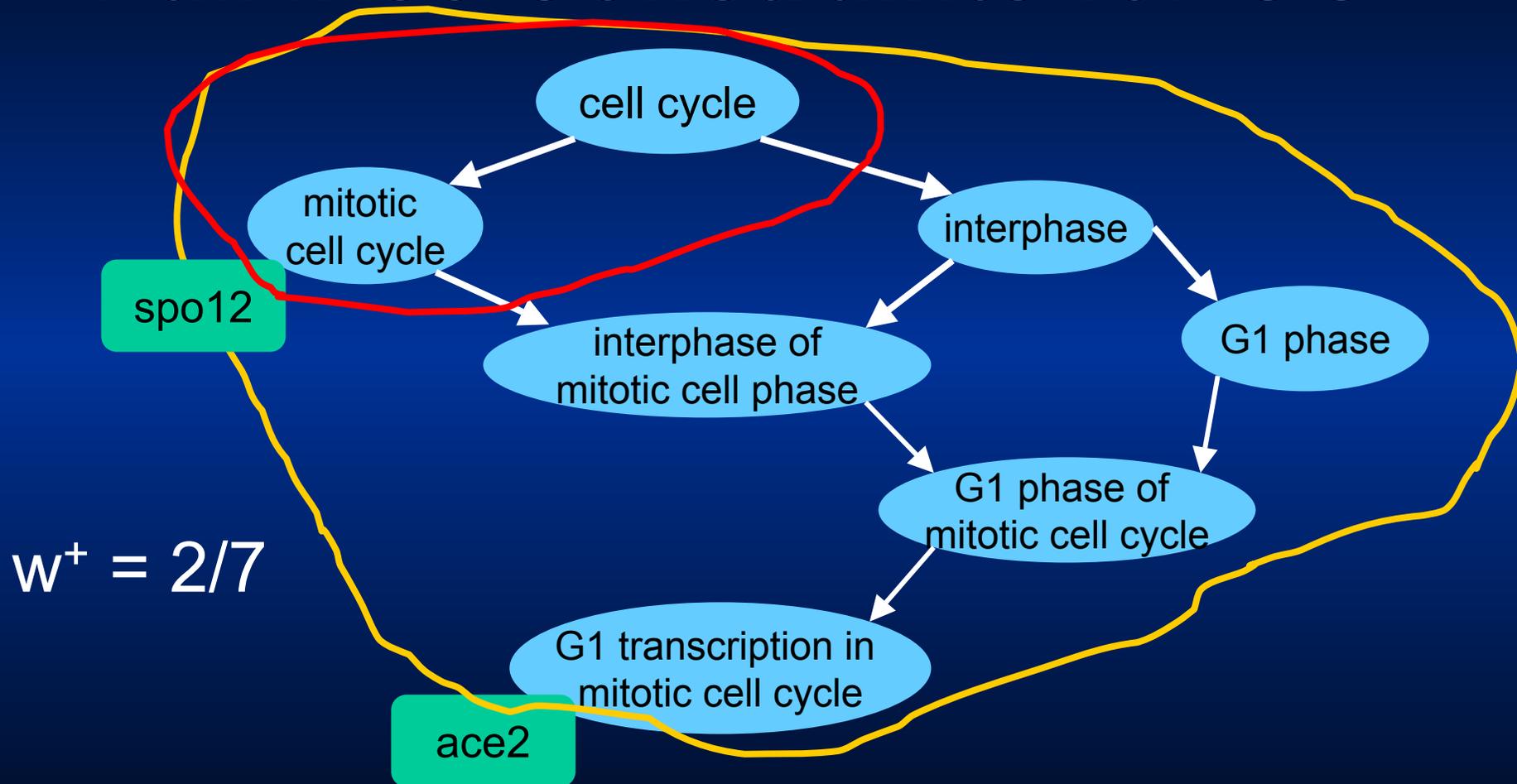
Data

- Gene expression data
 - time-courses of 384 genes during mitotic cell division in Yeast (Cho, 1998)
 - expert classification into five cell-cycle phases
- Constraints
 - Gene Ontology (GO)
 - transcription factor location analysis (Lee, 2002)
 - true labels

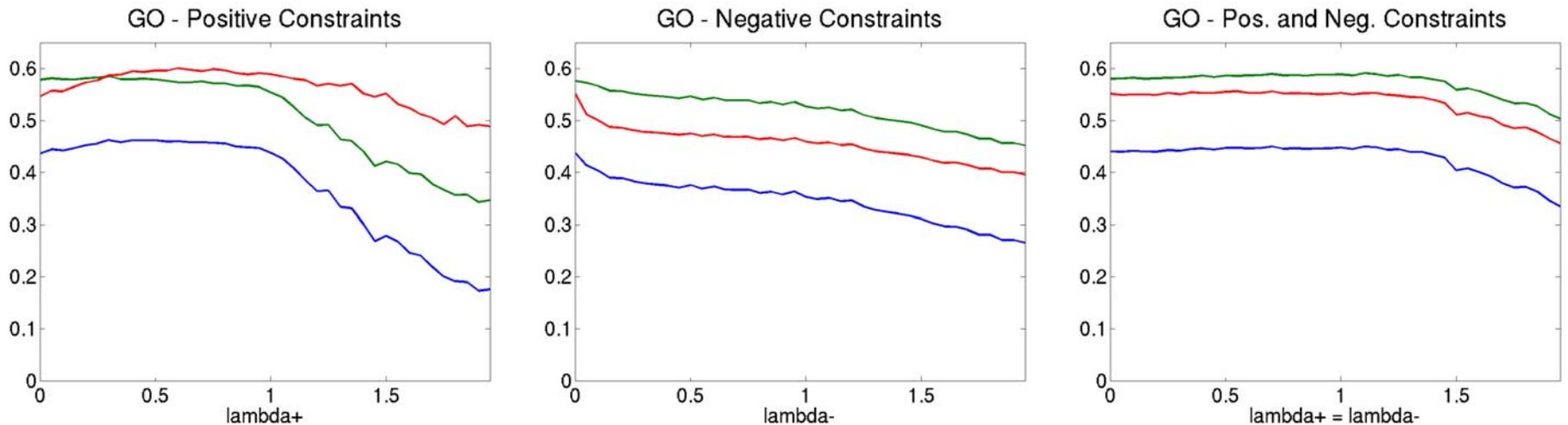
Pairwise Constraints for GO



Pairwise Constraints for GO



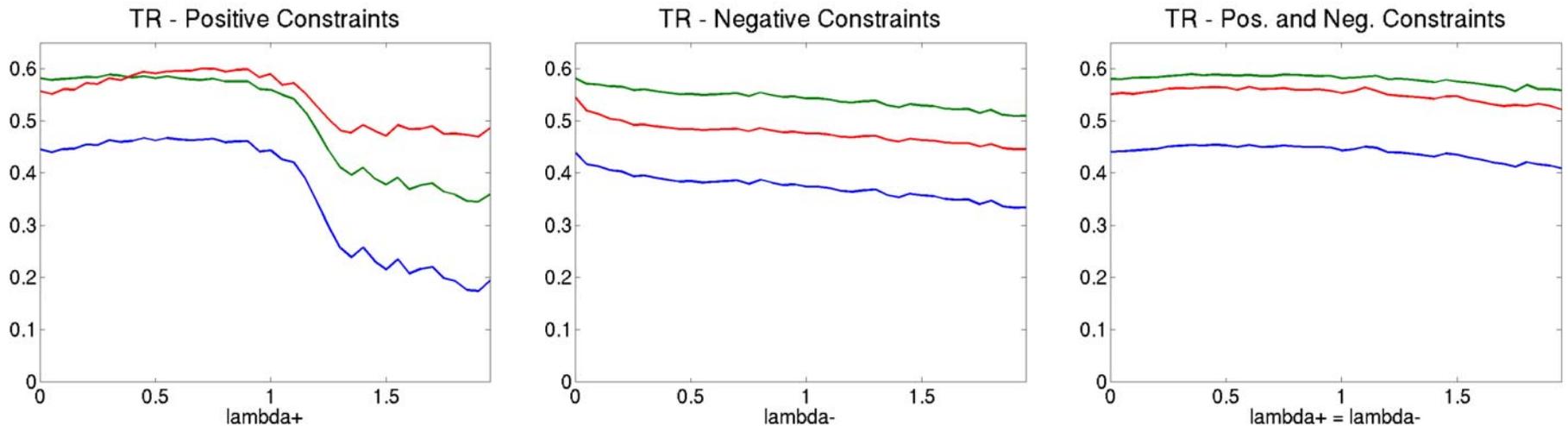
Constraints from GO



— corrected Rand
— specificity
— sensitivity

51% of gene
pairs constrained

Constraints from Location Analysis



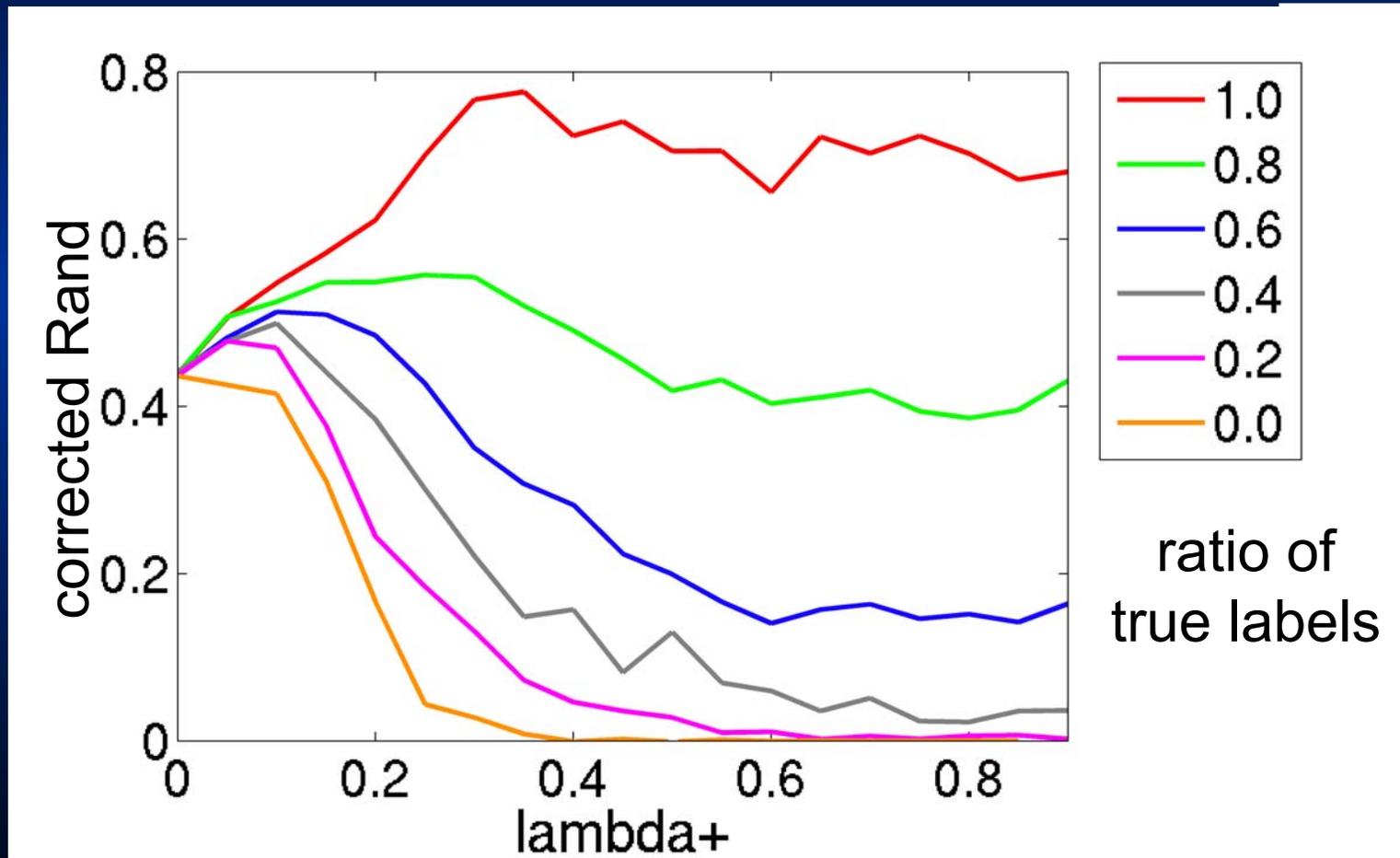
— corrected Rand
— specificity
— sensitivity

40% of gene
pairs constrained

Possible Explanations

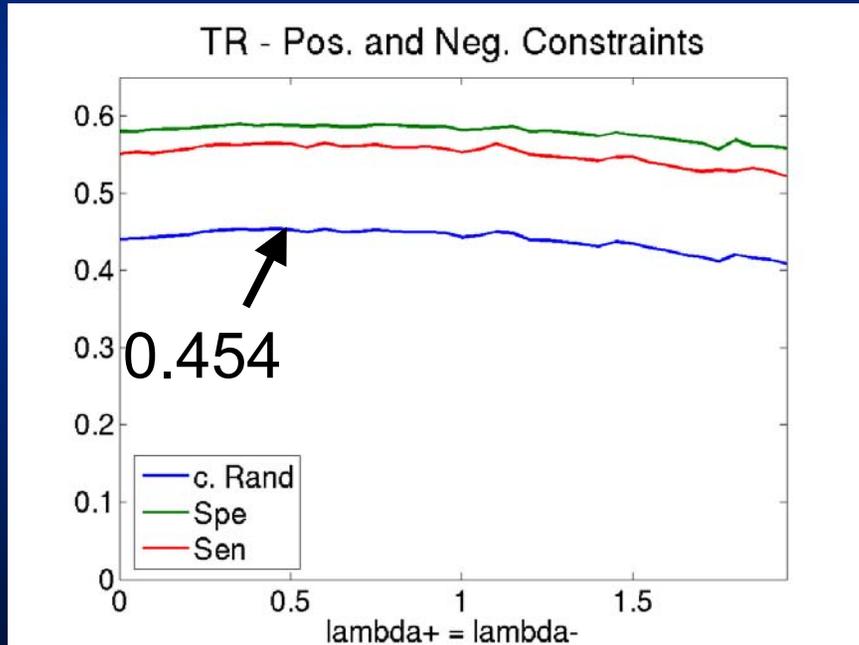
- Non-specific information content
- Noise in the data
- ...

Constraints from True and Random Labels

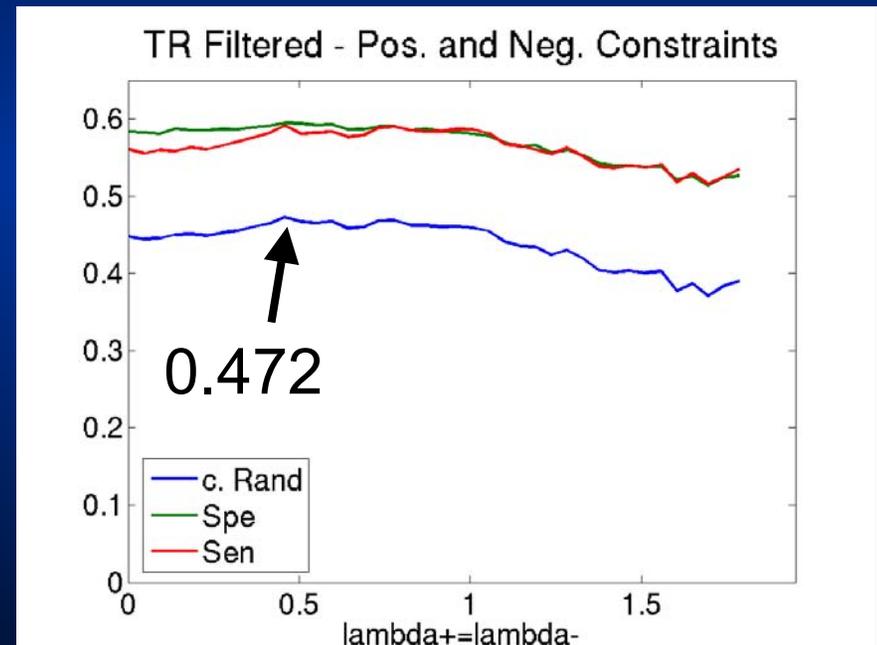


Filtered Constraints from Location Data

Non-filtered



Filtered



Summary and Outlook

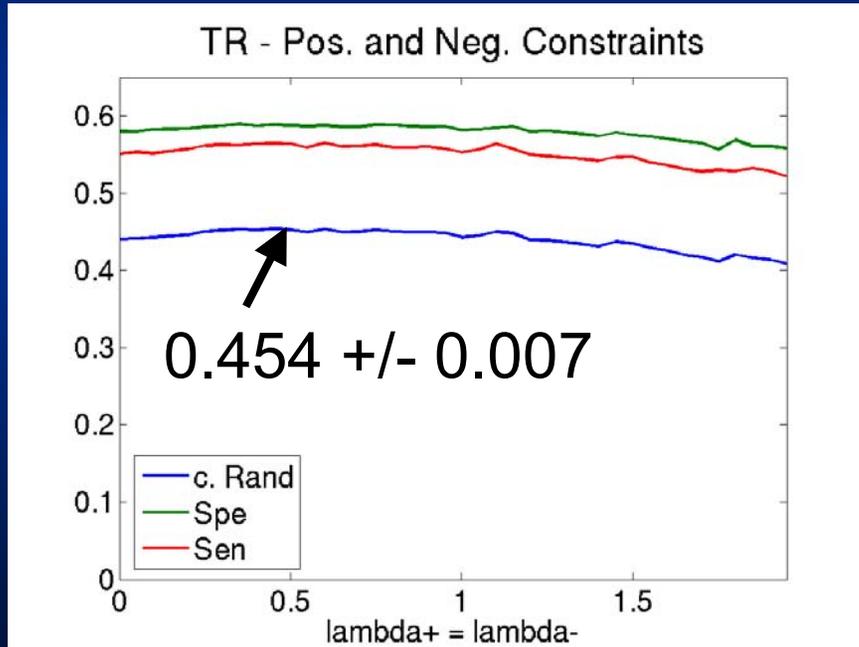
- A naïve analysis of heterogeneous data brings little improvements
- Data semantics?
 - GO does not yield further improvement
- Discovering relevant information?
 - learning of relevant constraints
 - find a right balance between $\lambda+$ and $\lambda-$

Thanks.

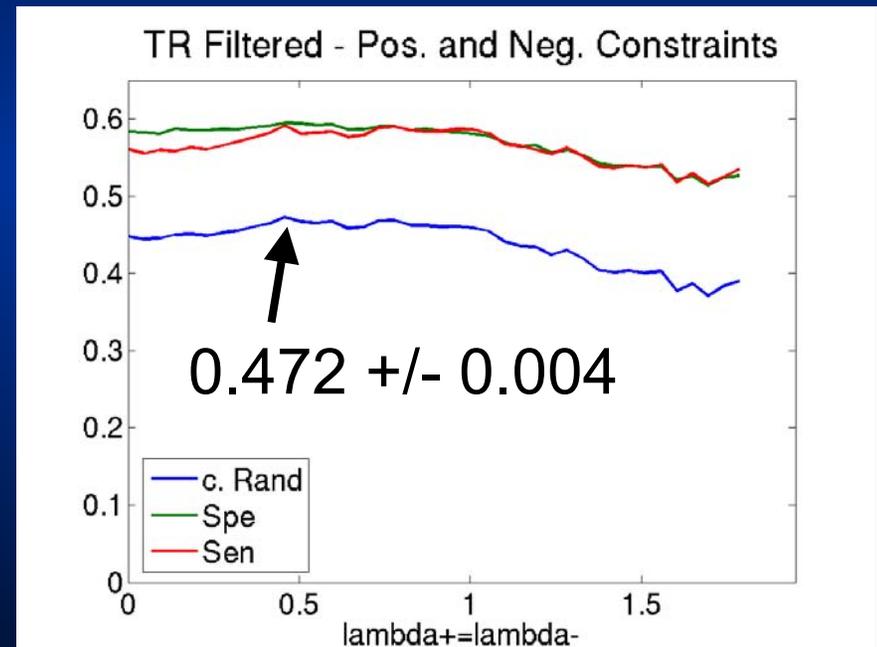
<http://algorithmics.molgen.mpg.de>

Filtered Constraints from Location Data

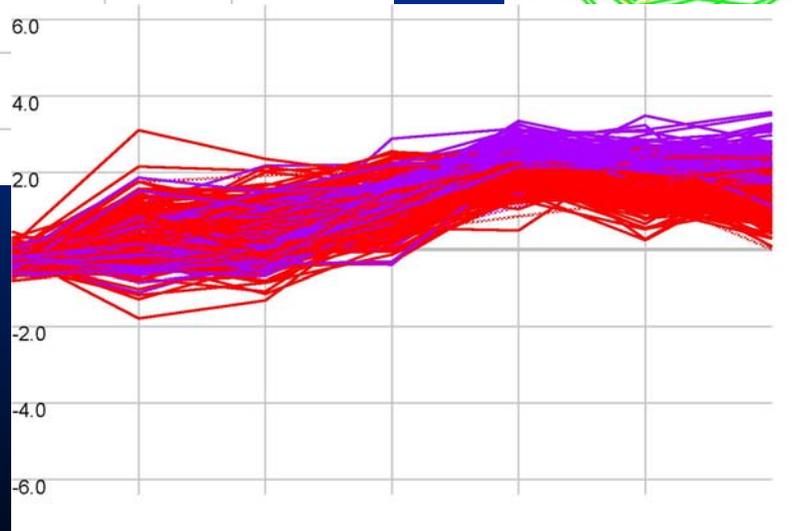
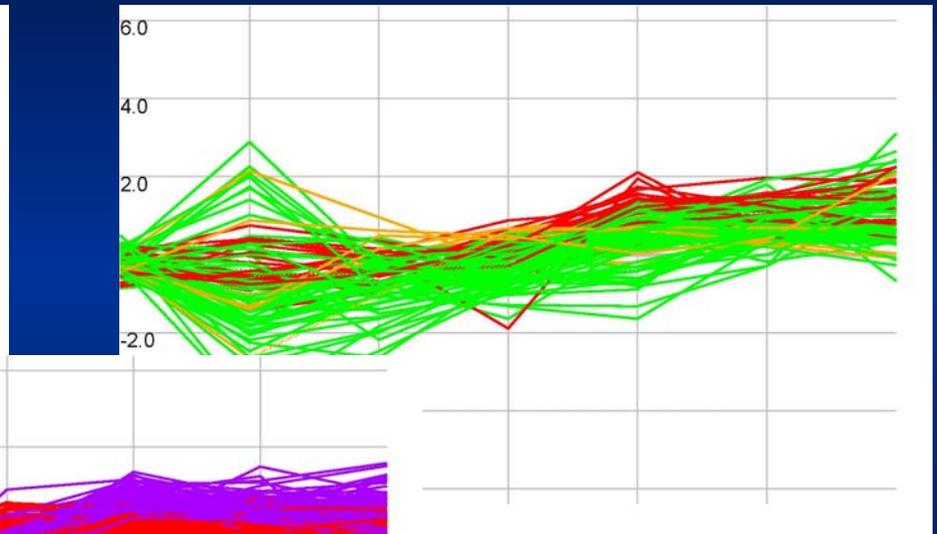
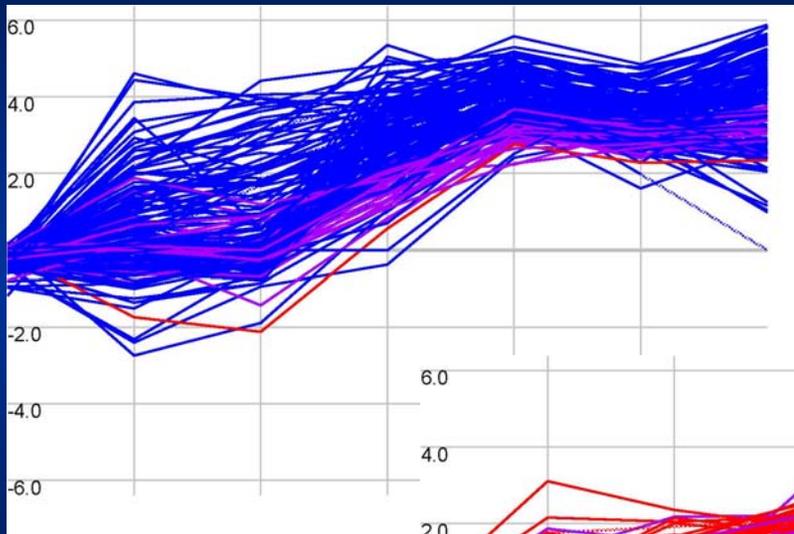
Non-filtered

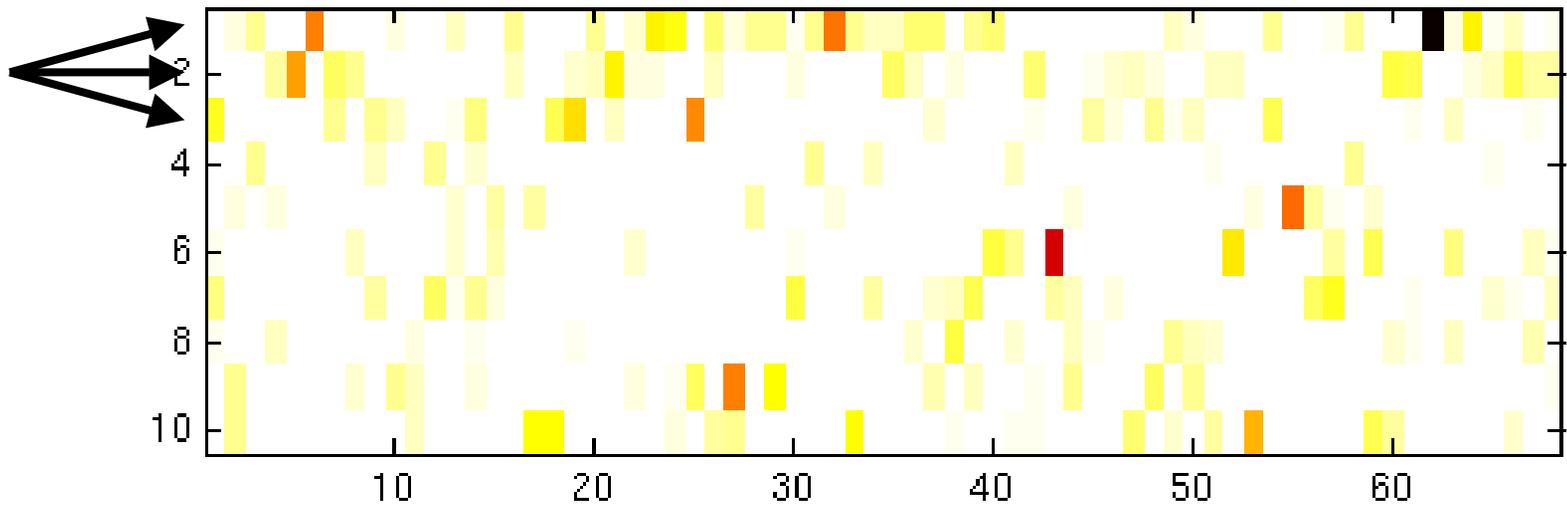


Filtered

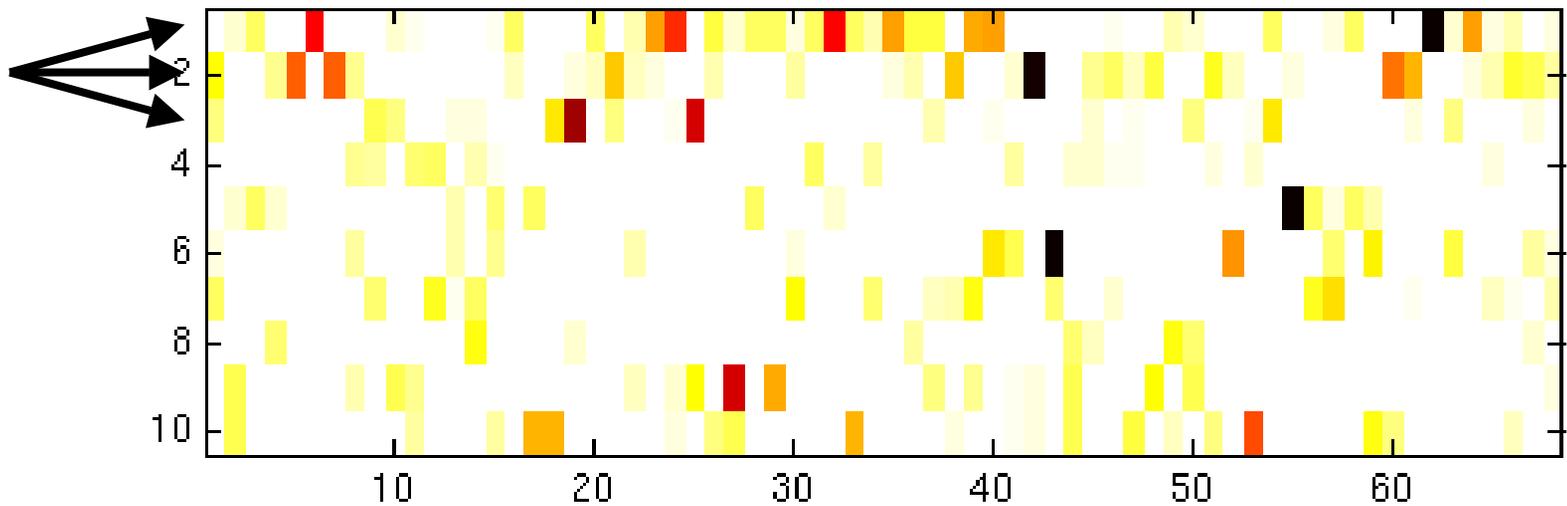


Examples - Sporulation data





mixture estimation



mixt. est. with constraints

Pairwise Constraints for GO

We define the constraints for all pair of annotated genes:

$$w_{ij}^+ = \frac{\#\{t_m \in \text{Dag}(g_i) \cap \text{Dag}(g_j)\}}{\#\{t_m \in \text{Dag}(g_i) \cup \text{Dag}(g_j)\}}$$

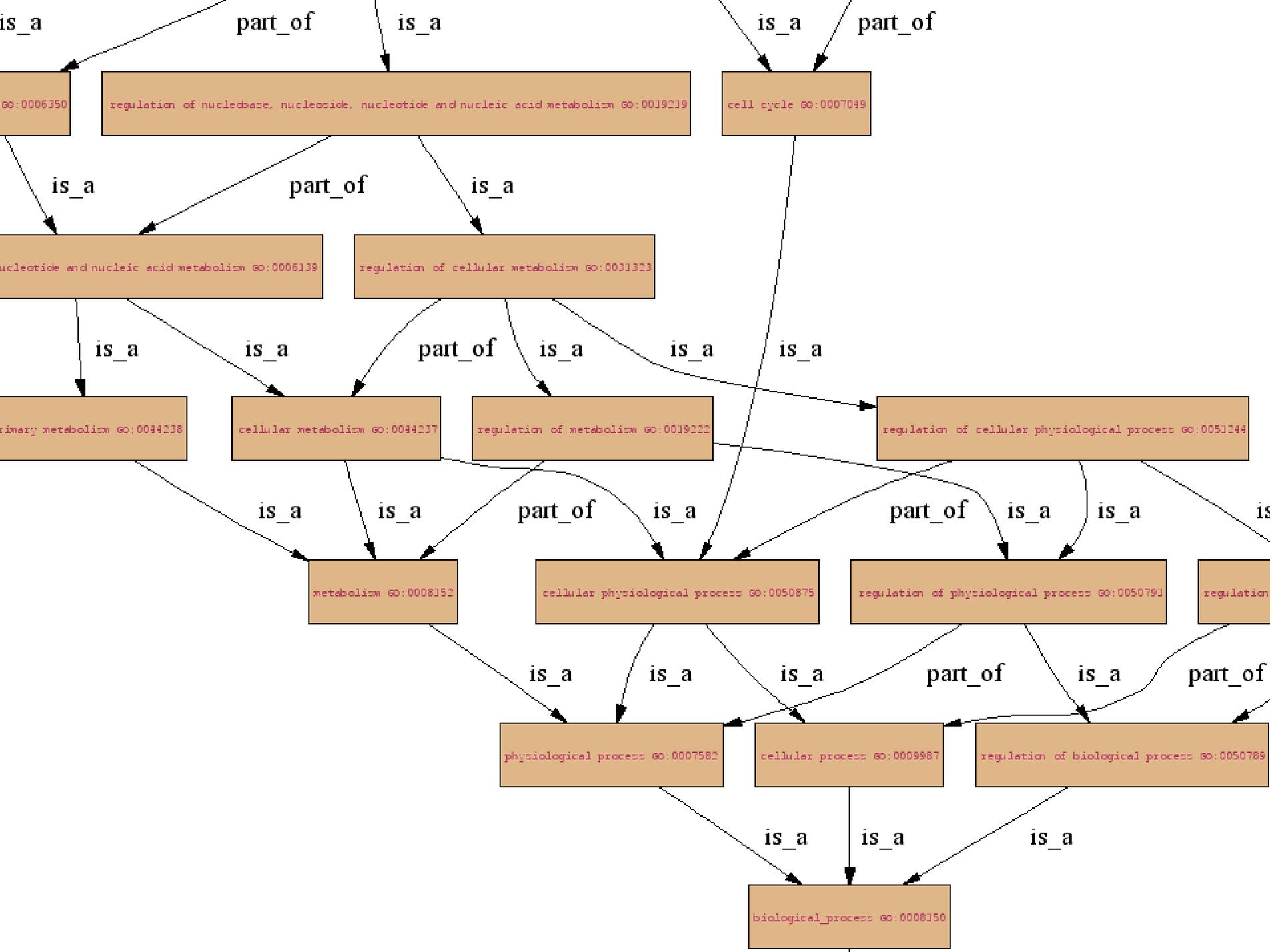
$$w_{ij}^- = \frac{\#\{t_m \in \text{Dag}(g_i) \cup \text{Dag}(g_j)\} - \#\{t_m \in \text{Dag}(g_i) \cap \text{Dag}(g_j)\}}{\#\{t_m \in \text{Dag}(g_i) \cup \text{Dag}(g_j)\}}$$

t_m are the terms in GO

$\text{Dag}(g_i)$ is the set of terms annotating g_i

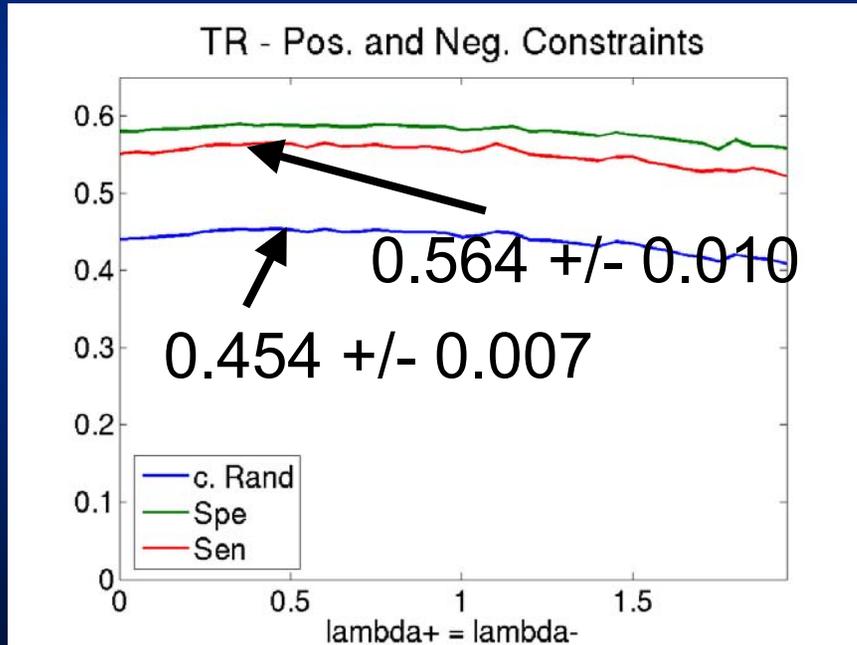
More is less.

Sometimes ...

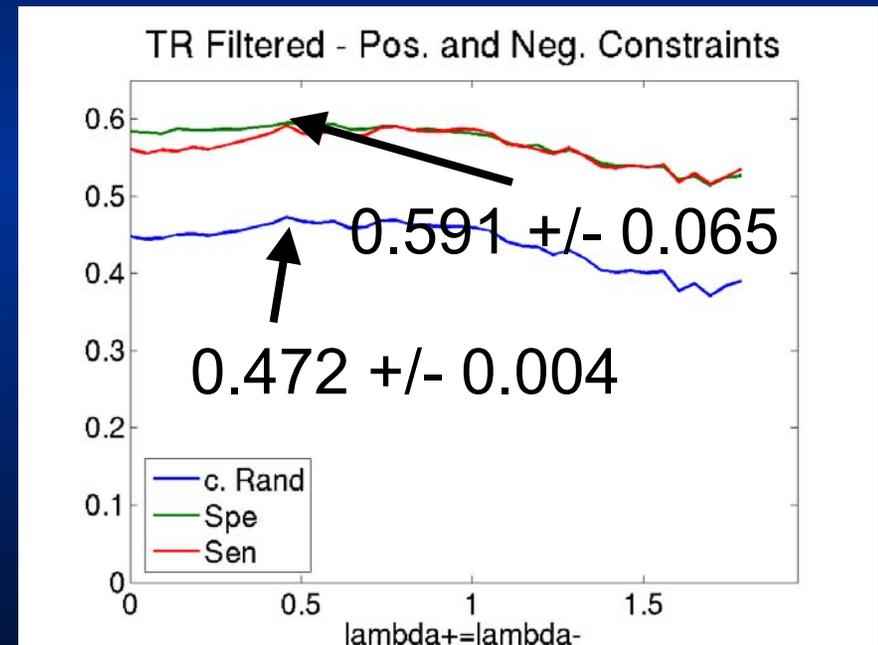


Expert Filter of Constraints Transcription Factors

Non-filtered

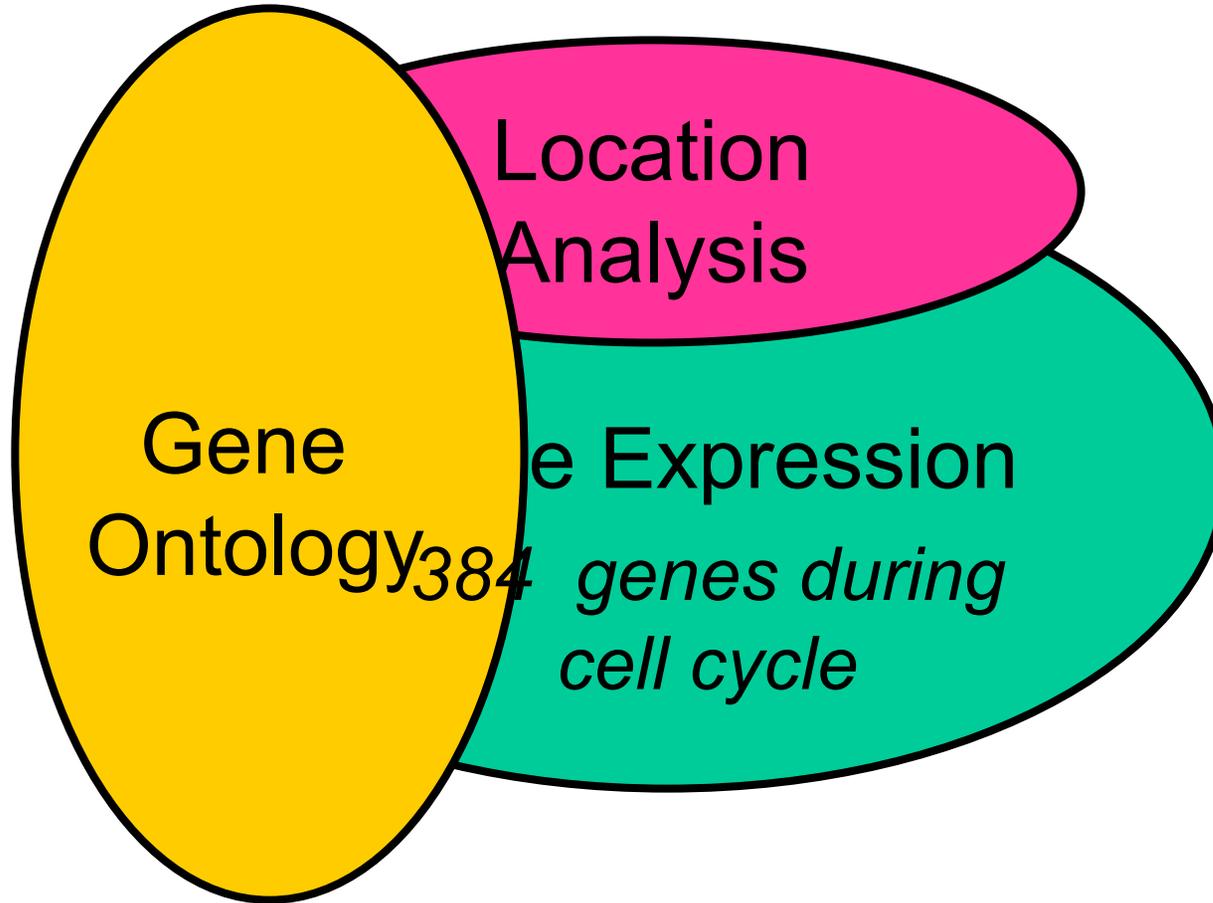


Filtered



Data

Ω



Mixture Estimation with Constraints (3)

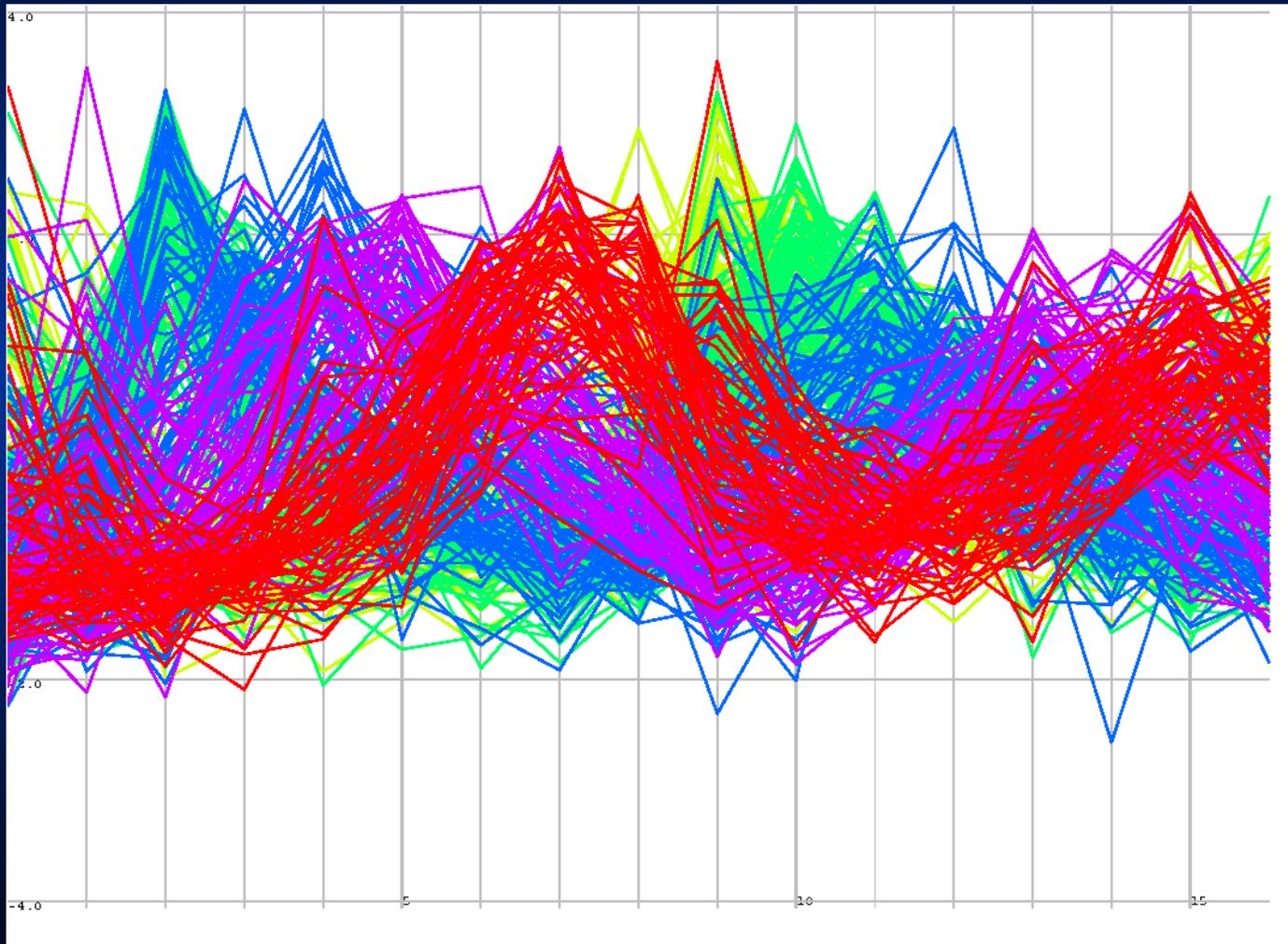
The posterior $P[y_i=k | Y_i', X, W, \theta_k]$ is approximated by means of Gibbs sampling :

$$\alpha_k P[x_i | \theta_k] \exp \sum_{j \neq i} -\lambda^+ w_{ij}^+ (1 - P[y_j = k | Y_i', X, W, \theta_k])$$

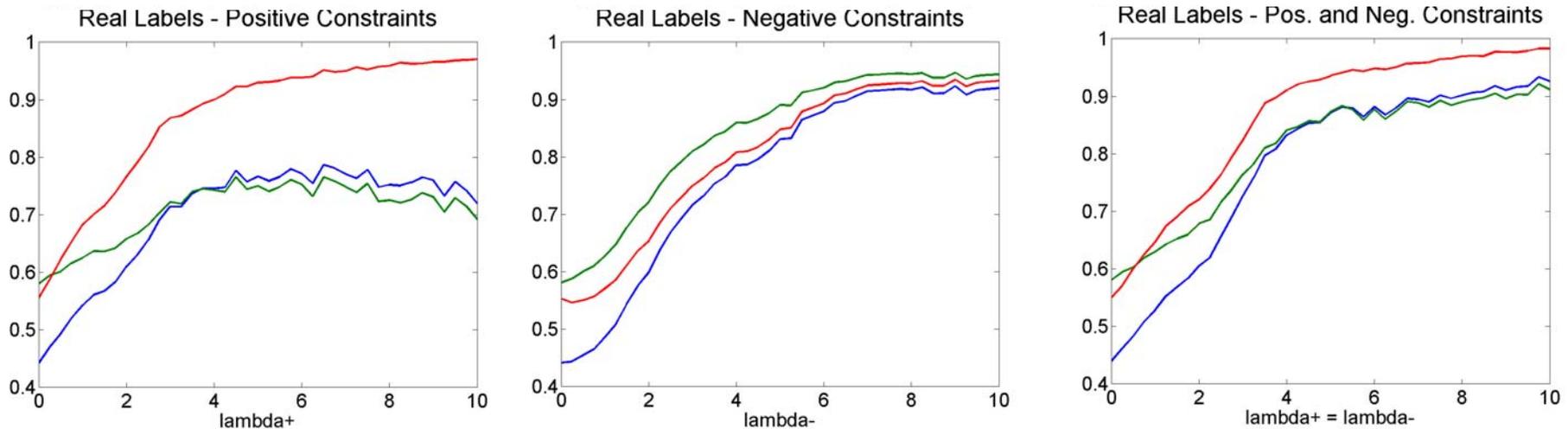
where

$$Y_i' = \{y_1^t, \dots, y_{i-1}^t, y_{i+1}^{t-1}, \dots, y_N^{t-1}\}$$

Yeast Cell Cycle



Constraints from True Labels



— corrected Rand
— specificity
— sensitivity

5% of gene
pairs constrained