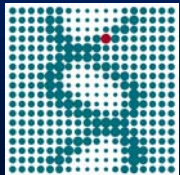


Clustering of Gene Expression Time Courses: Methods and Validity of Solutions

Ivan G. Costa Filho
Alexander Schliep



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

Goal

(1) Clustering of gene
expression **time-courses**

(2) ***Validating clusterings***
of genes

(3) Methods for clustering
heterogeneous data

gene expression + functional annotation, regulatory region
information, protein-protein interactions ...

(1) Clustering Method

Time-course models

Biological Truism

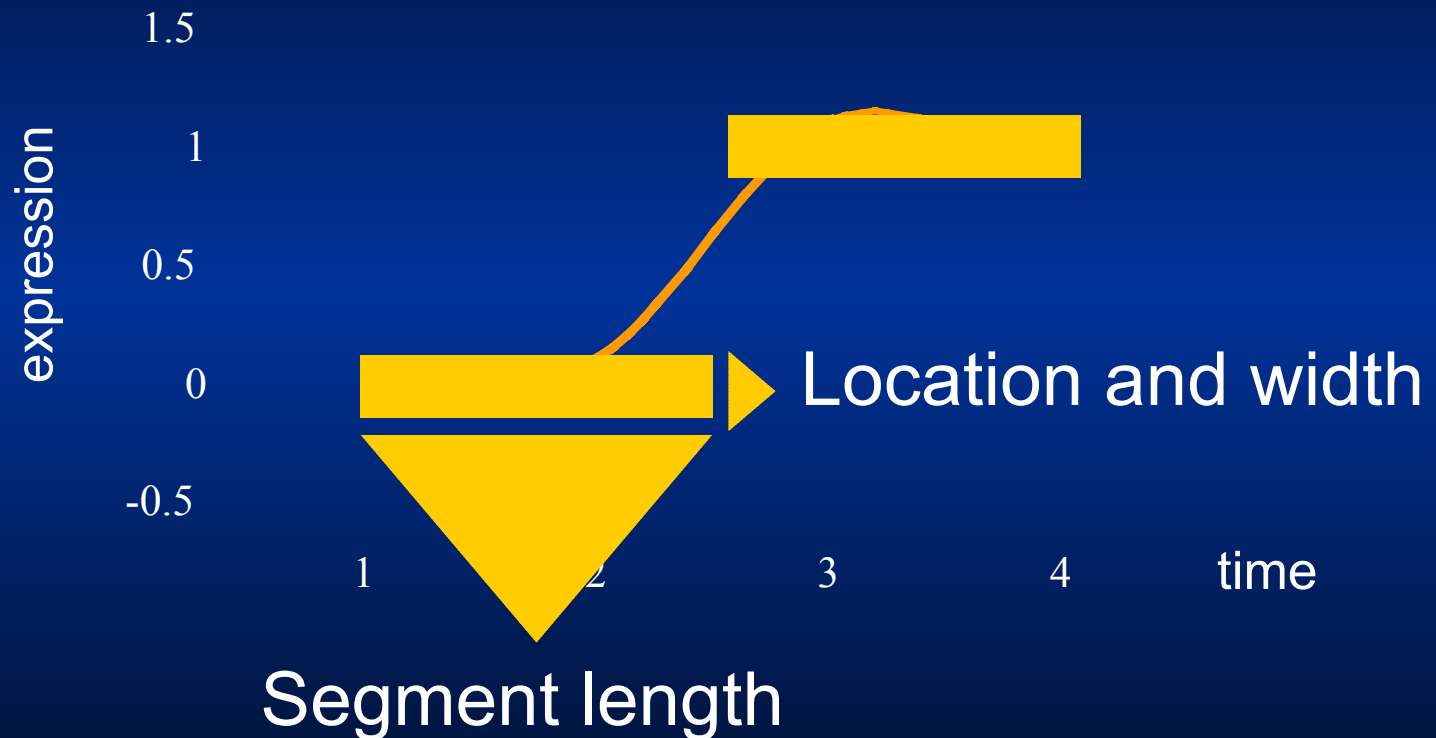
- Many genes have
 - multiple functions
 - are involved in several regulatory networks
- *Unique assignment to groups dubious*

Method Outline

1. Define a class of statistical models for time-courses
2. Combine them in a mixture model
3. Decode the mixture to infer groups

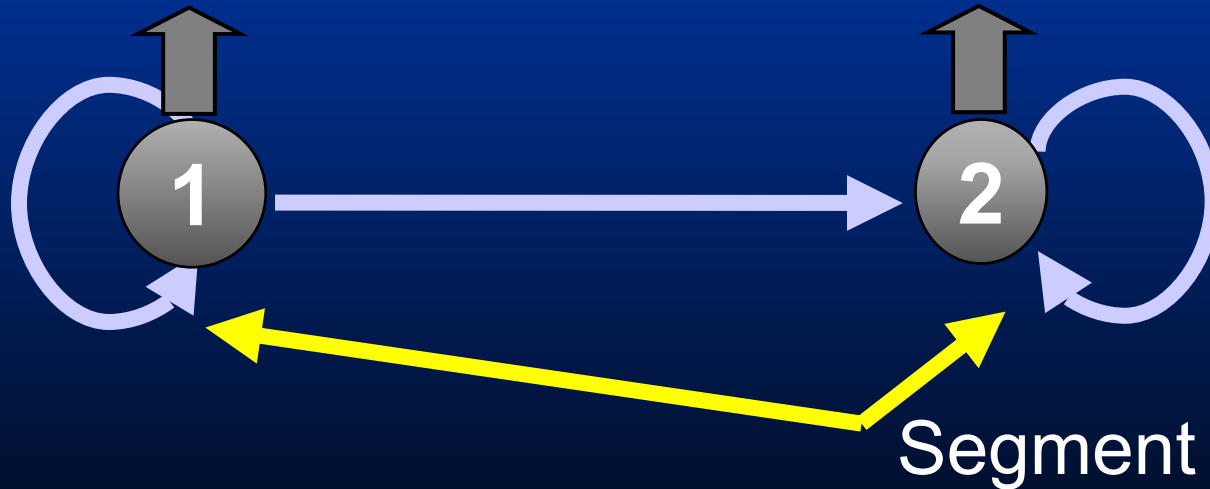
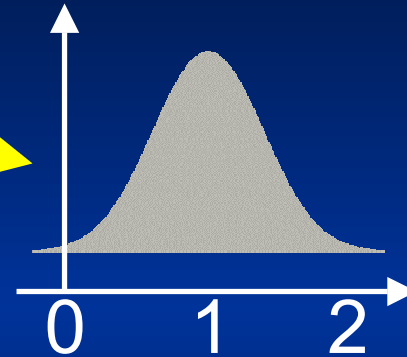
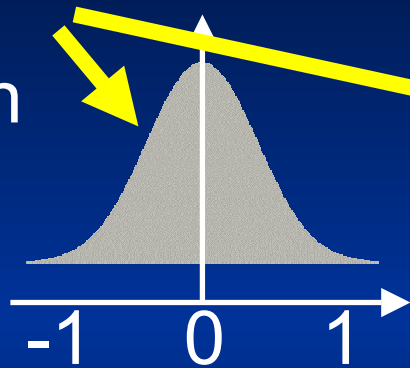
Time-course models

Example: Up-regulation



Prototype: Up-regulation

Location
and width



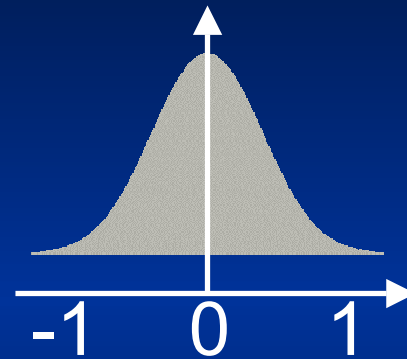
Segment length

Hidden Markov Model (HMM)

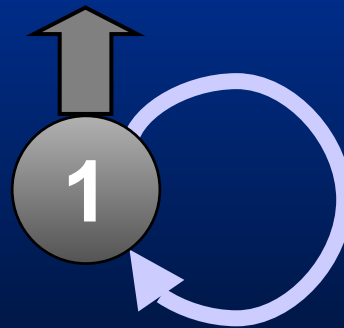
Example: Constant



Prototype: Constant



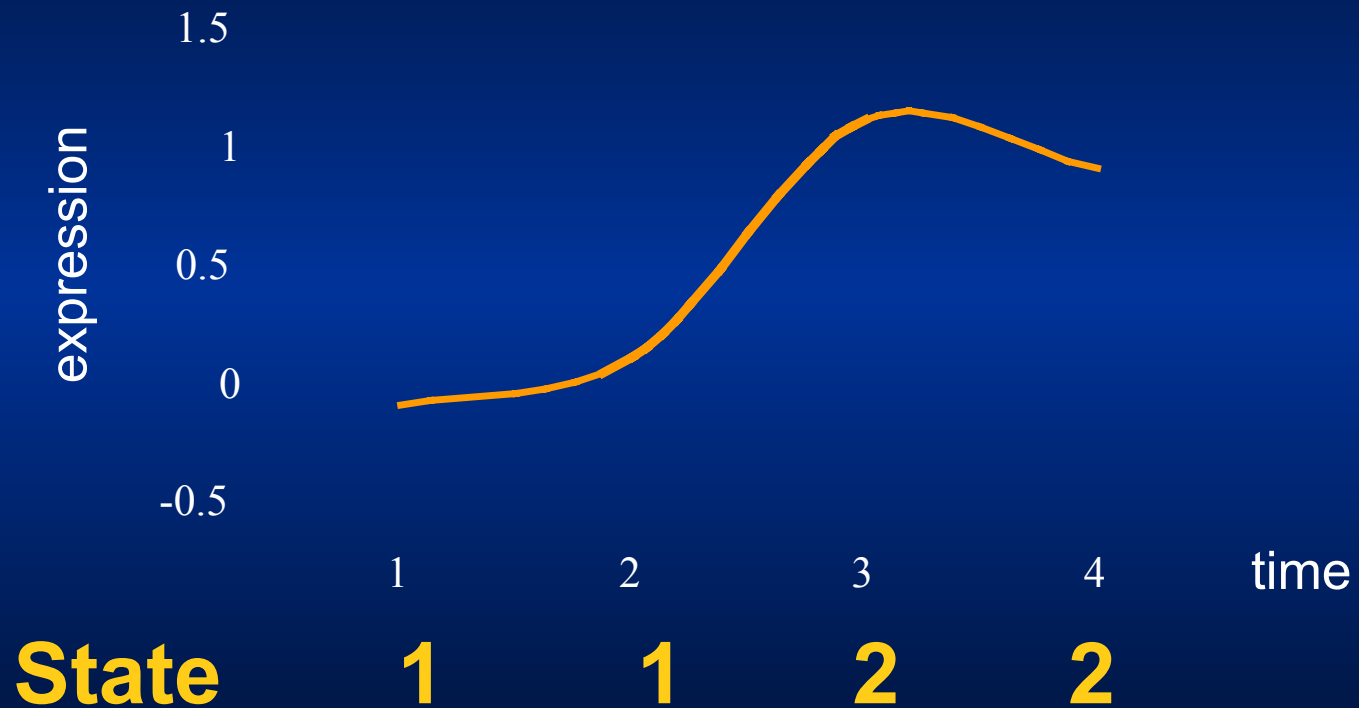
Location and width



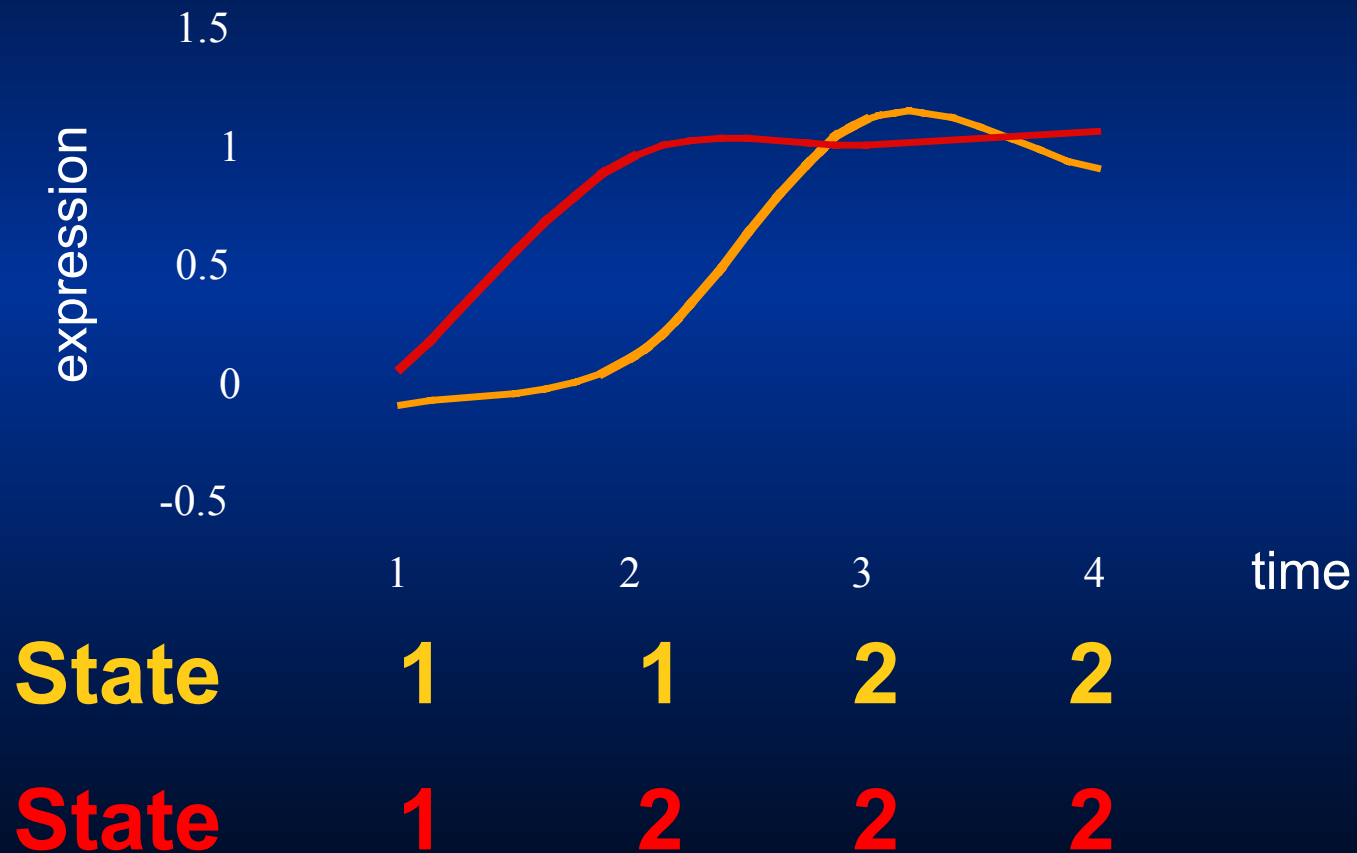
Segment length

Hidden Markov Model (HMM)

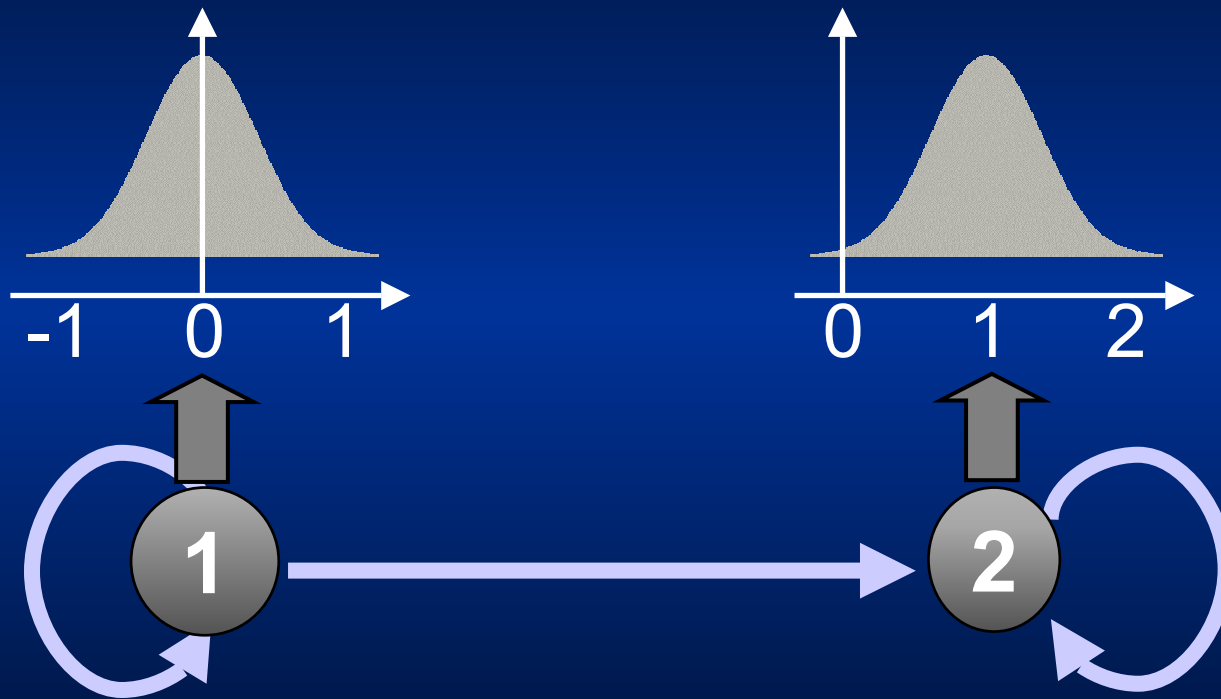
Viterbi Path: Up-regulation



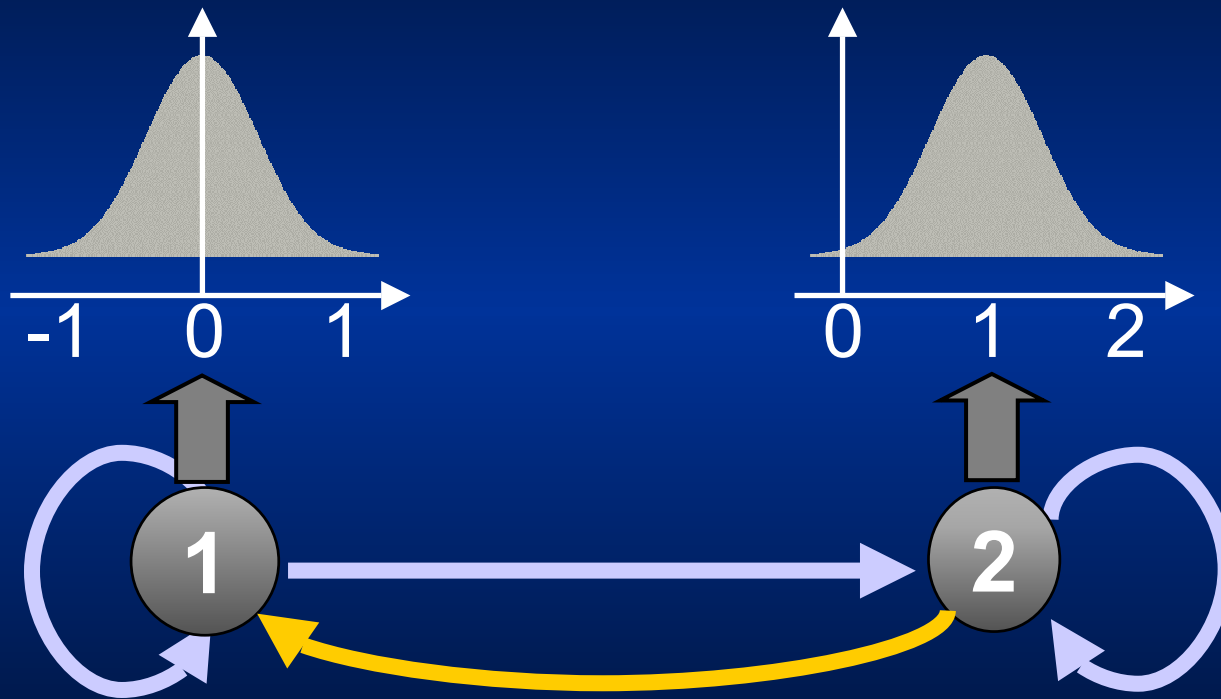
Viterbi Path: Up-regulation



Prototype: Cyclic



Prototype: Cyclic



Perspective

- #states = #time-points:
 - linear HMM " multi-variate Gaussian
 - covariance matrix $\text{diag}(\sigma_1, \dots, \sigma_t)$
- Typically #states \ll #time-points

Mixture of HMMs

Mixture models

- Mixture components: HMMs $\lambda_1, \lambda_2, \dots, \lambda_k$
- Mixture model " weighted sum of λ_i

$P[\text{gene} \mid \text{mixture}] =$

$$\alpha_1 P[\text{gene} \mid \lambda_1] + \alpha_2 P[\text{gene} \mid \lambda_2] + \dots$$

$\alpha_i \geq 0$, add to unity

HMM-based 'Clustering'

- Input:
 - genes profiles g_i
 - collection of k **HMMs**
- Initialization:
 - Assign the probability that a data-point belongs to each k **HMMs** randomly
- Iteration (until convergence of assignment):
 - Compute the **new HMM parameters (B-Welch)**
 - Re-assign g_i to a **HMM** proportionally to $P[\text{gene} \mid \lambda_1]$

Inference of groups

From Mixtures to Groups

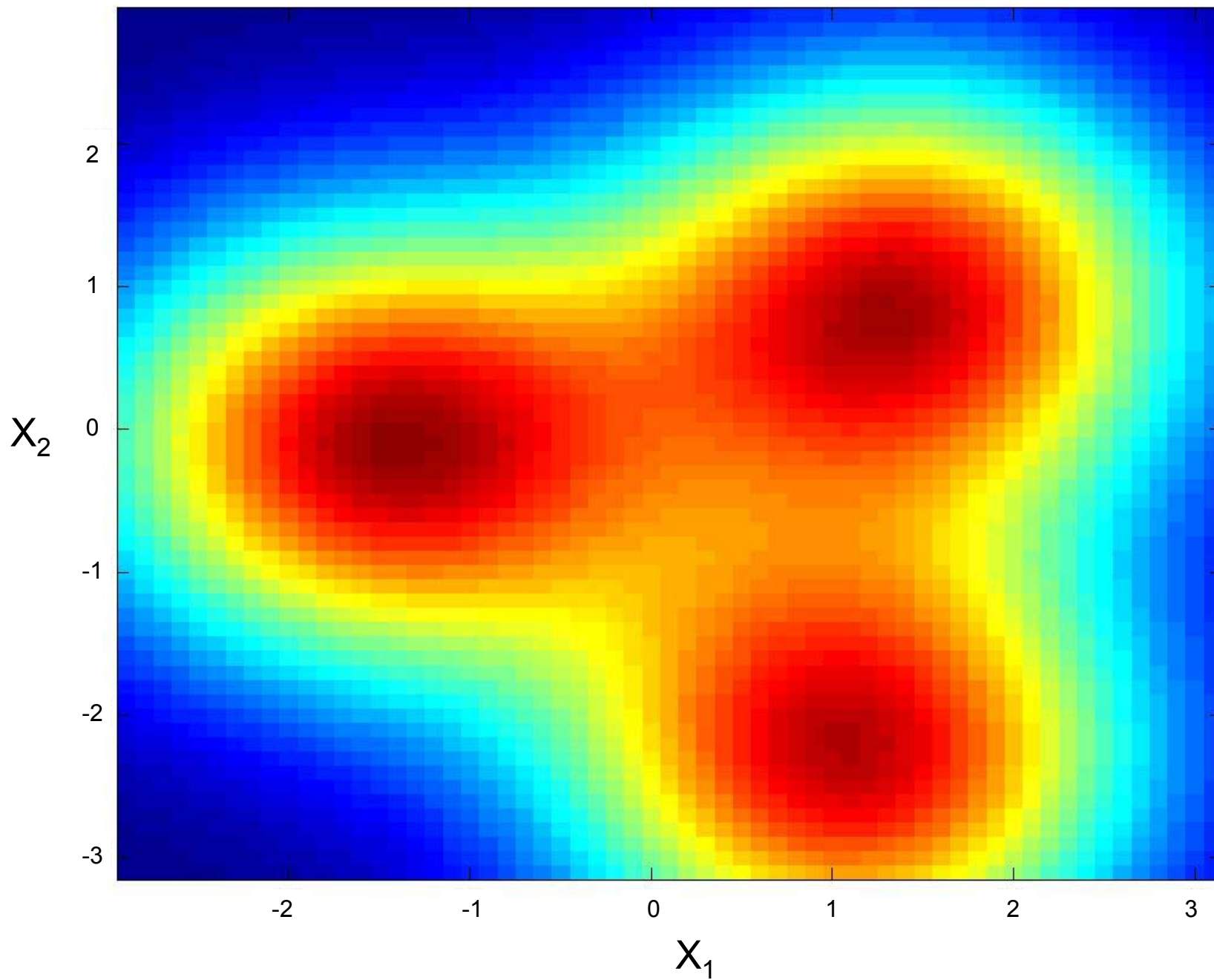
- Posterior probability of mixture component λ_i

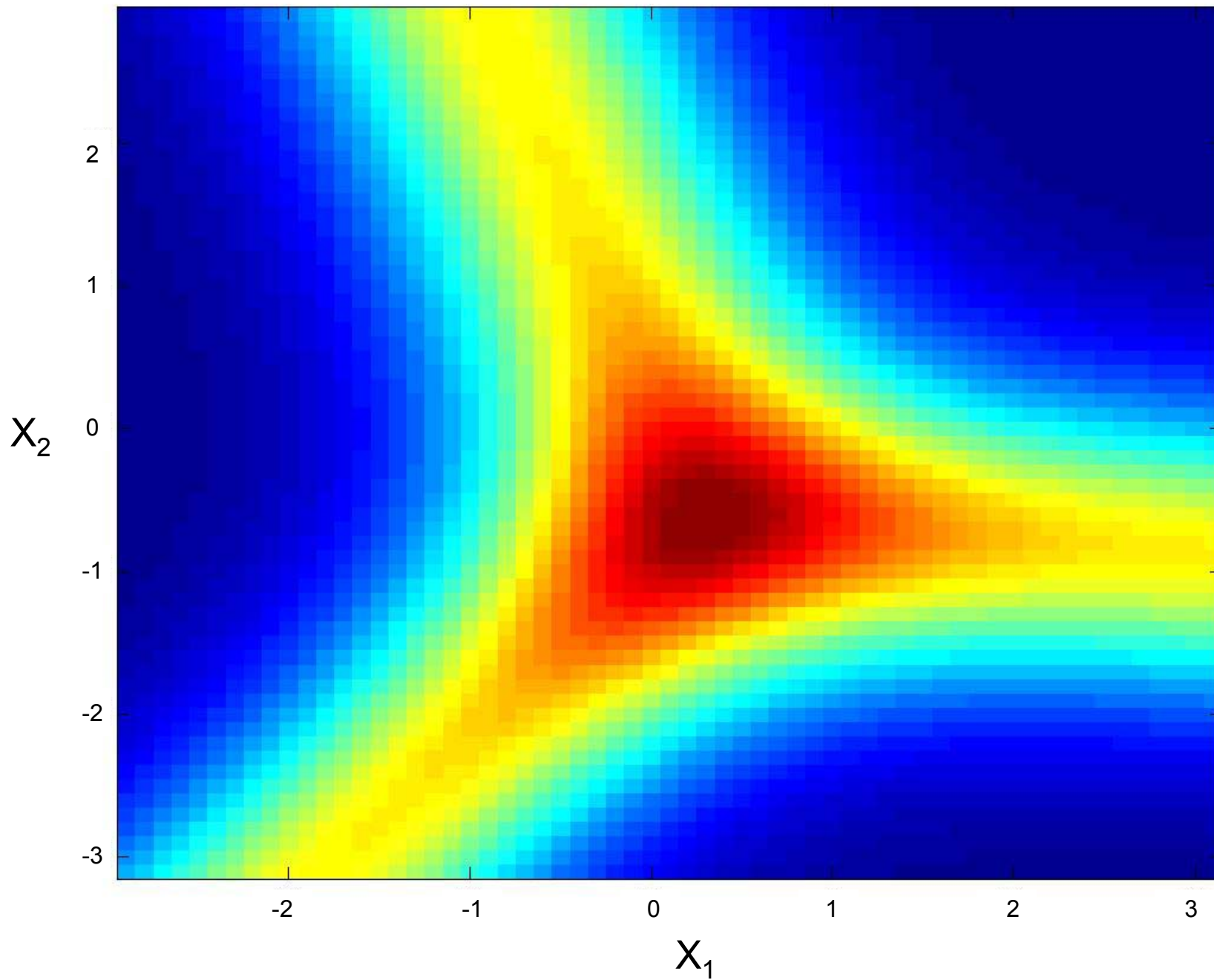
$$P[\lambda_i | gene]$$

- Shannon-entropy

$$H(\{ P[\lambda_i | gene] \}_{1 \leq i \leq k})$$

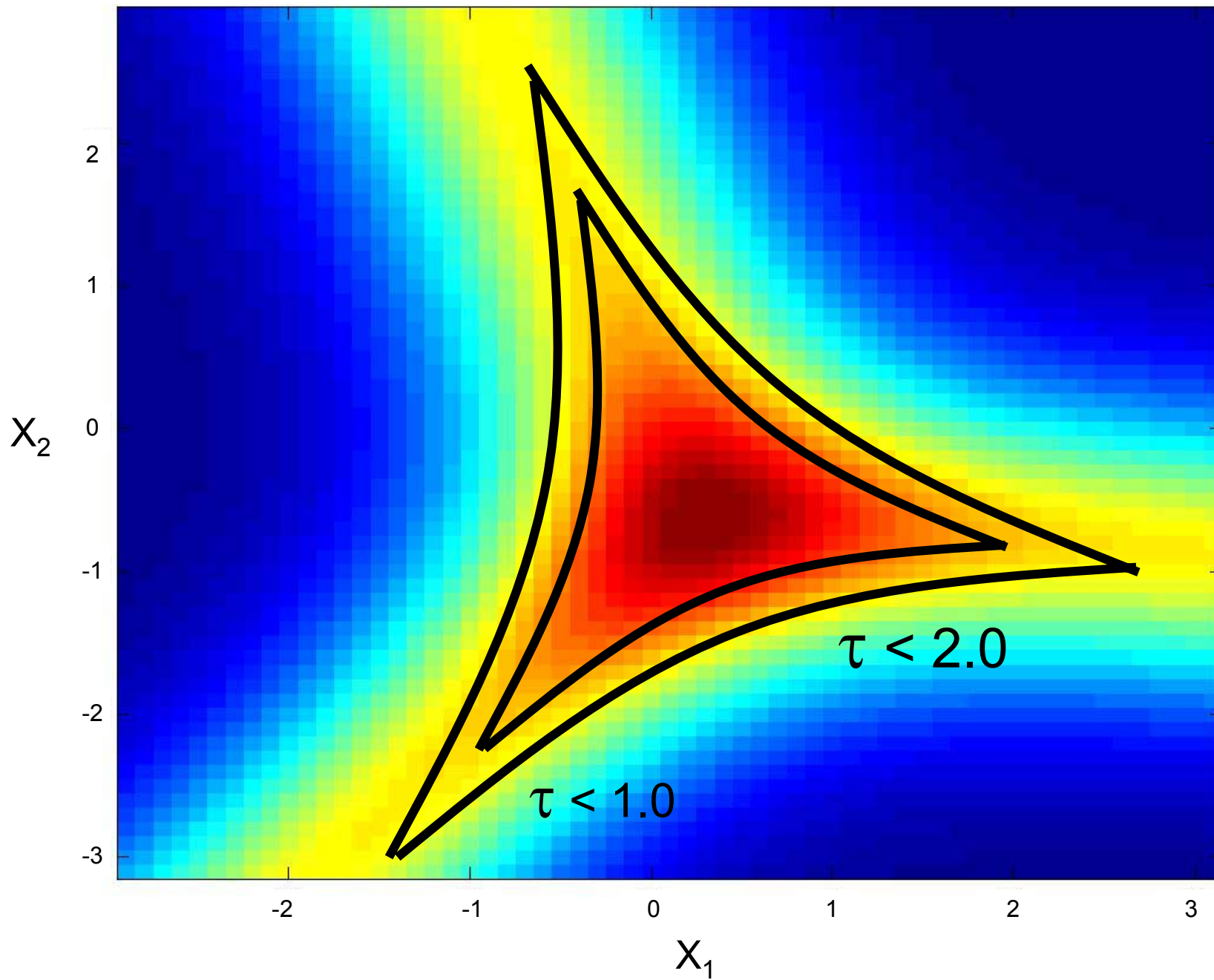
quantifies level of ambiguity in assignment





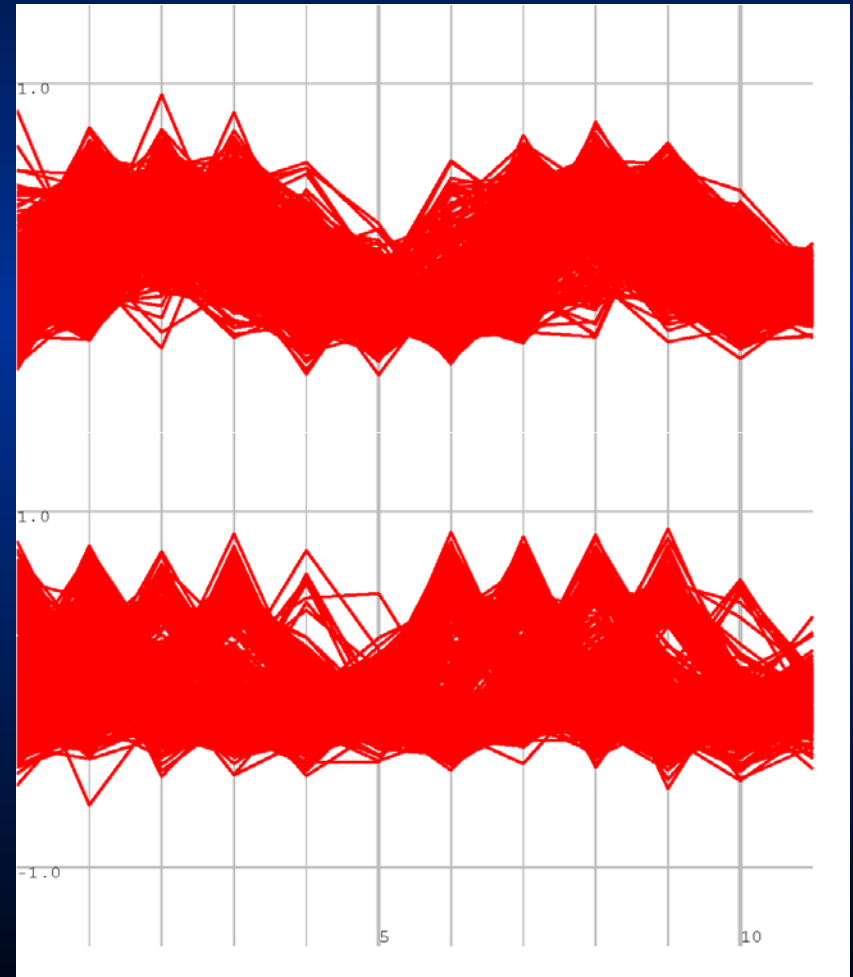
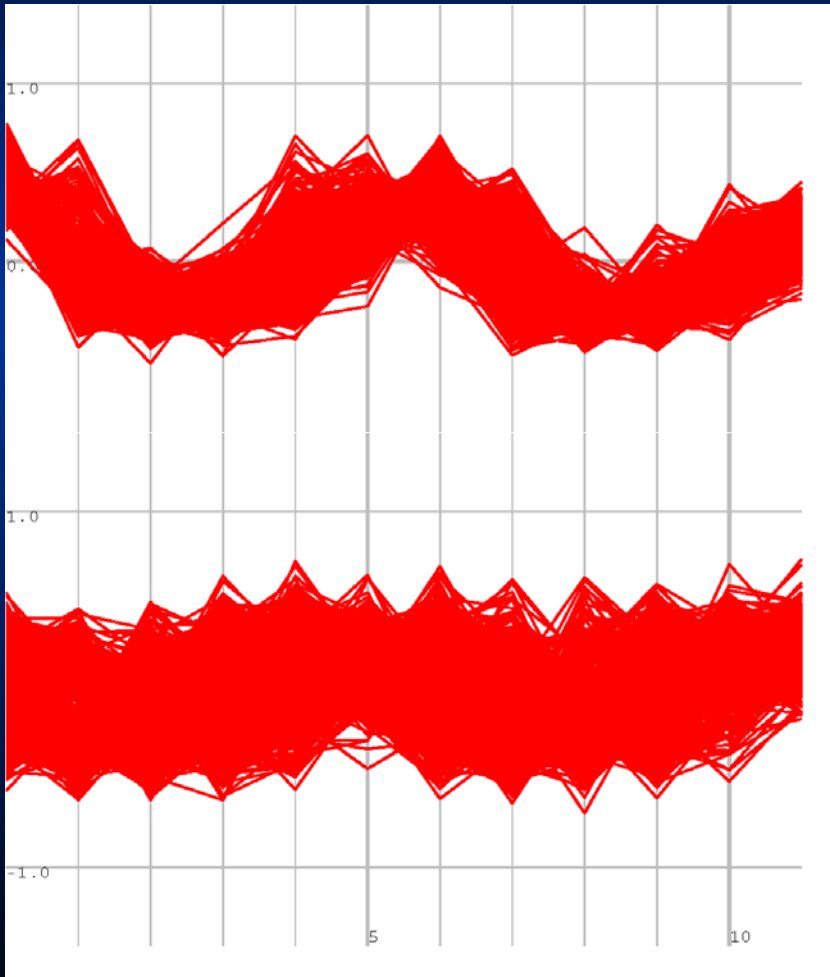
From Mixtures to Groups

- Choose entropy threshold τ
- If entropy of posterior is below τ
 - Assign gene to group i of maximal posterior
- Else:
 - leave gene unassigned



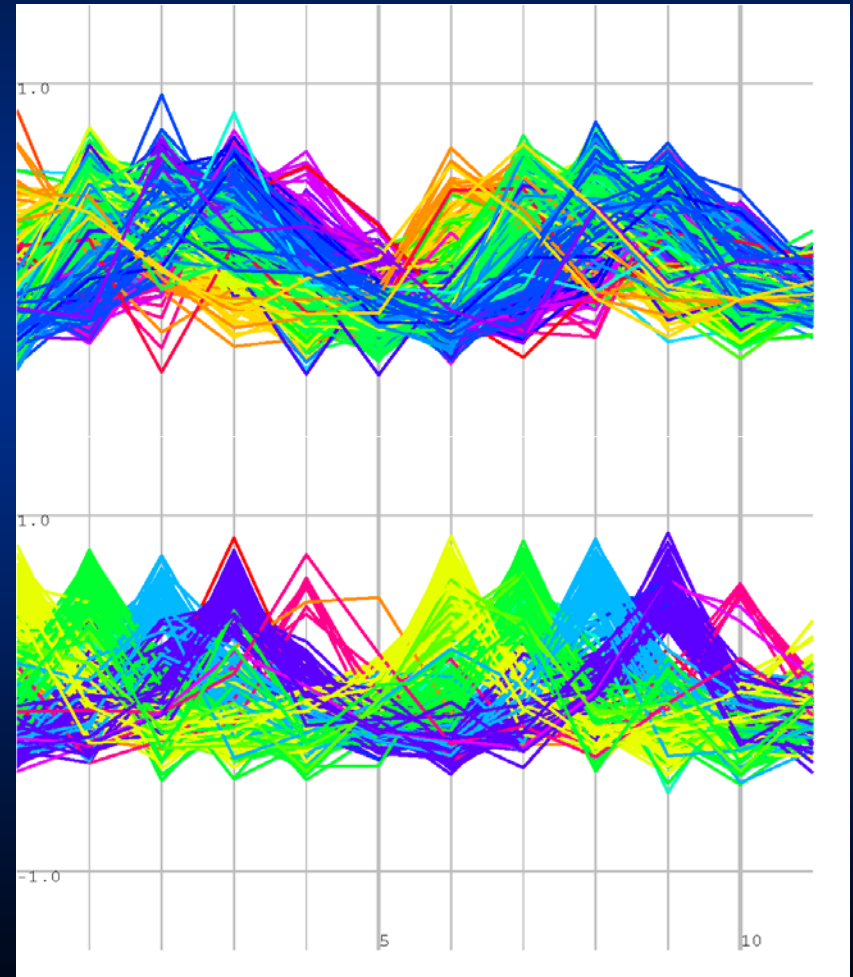
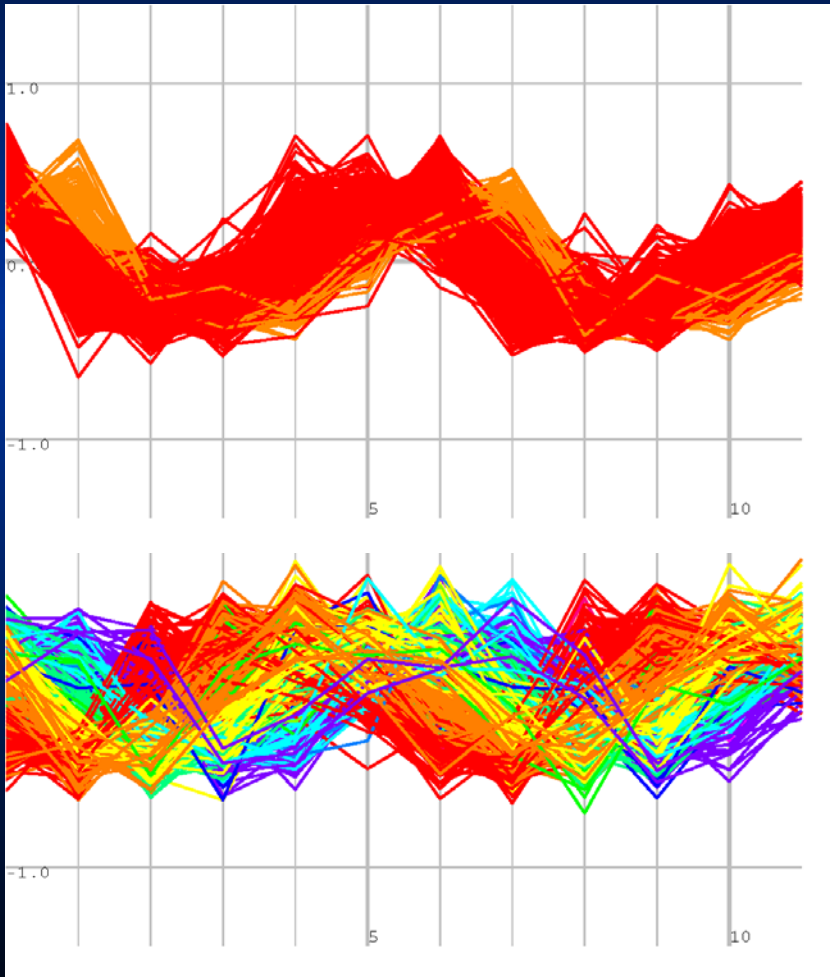
Results

Diurnal Time Courses



Results

Diurnal Time Courses



(2) Cluster Validation

Motivation

- Which clustering method to use?
- How good is a clustering?
- GO enrichment gives no comparative basis!

Use **gene annotation (GO)** as external data to validate **mixtures** (or clusterings) from gene expression

External Indices

In clustering, external indices look for the number of genes pairs that

	<i>Same categ.</i>	<i>Distinct categ.</i>
<i>Same cluster</i>	True Positive	False Positive
<i>Distinct clusters</i>	False Negative	True Negative

External Indices

$$\textit{corr.Rand} = \frac{\#TP + \#TN - n_c}{\#Pairs - n_c}$$

External Indices for Mixture Model

Given mixture models U and V , consider the posteriors of the mixture components :

$$\{P[u_i | g]\}_{1 < i < C} \quad \{P[v_j | g]\}_{1 < j < R}$$

$g_k \equiv g_l \leftrightarrow$ the event of co-occurrence of g_k and g_l

$$P[g_k \equiv g_l \text{ given } U] = \sum_{j=1}^C P[u_j | g_k] \cdot P[u_j | g_l]$$

External Indices for Mixture Model

$$TP = \sum_{k=1}^N \sum_{l=k+1}^N P[g_k \equiv g_l \text{ given } U] \cdot P[g_k \equiv g_l \text{ given } V]$$

$$TN = \sum_{k=1}^N \sum_{l=k+1}^N P[g_k \equiv g_l \text{ given } U]^c \cdot P[g_k \equiv g_l \text{ given } V]^c$$

Biological Data Experiments

- Gene expression data:
 - Yeast during sporulation (7 time points)
 - 1027 genes after 2 fold filtering
- ‘Clustering’ Methods:
 - Hierarchical clustering (Pearson correlation)
 - K-means (Pearson correlation)
 - Mixture of HMMs
 - Mixture of Multivariate Normals (full covariance)

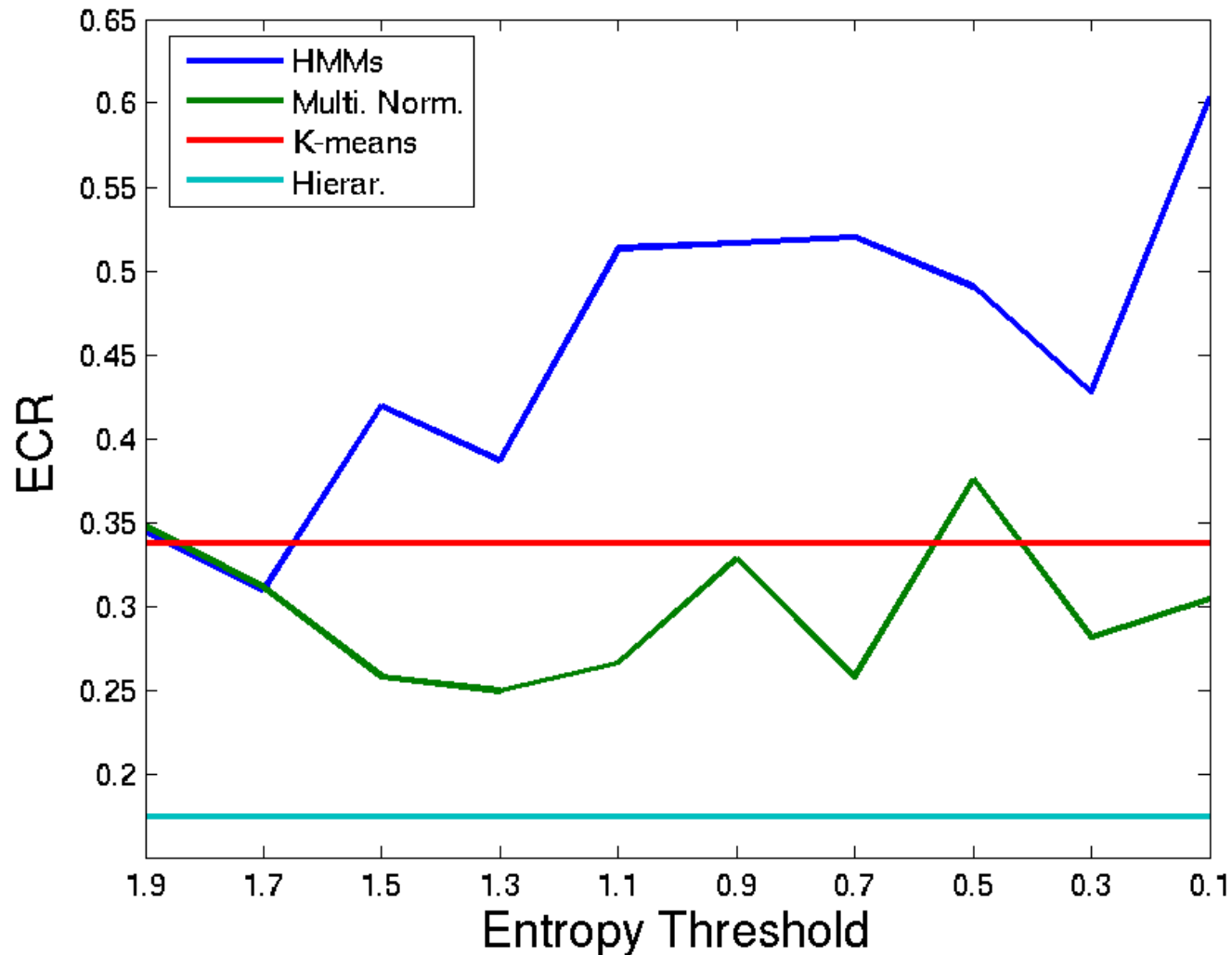
Biological Data Representation

- Mixture Model representation
 - Use each GO term as a component in the mixture
 - Maximum likelihood estimator of a multinomial distribution

$$P[t_i | g] = \begin{cases} 1/\#\{j | g \in t_j\}, & \text{if } g \in t_i \\ 0, & \text{otherwise} \end{cases}$$

Results

Methods x 'All' GO Terms

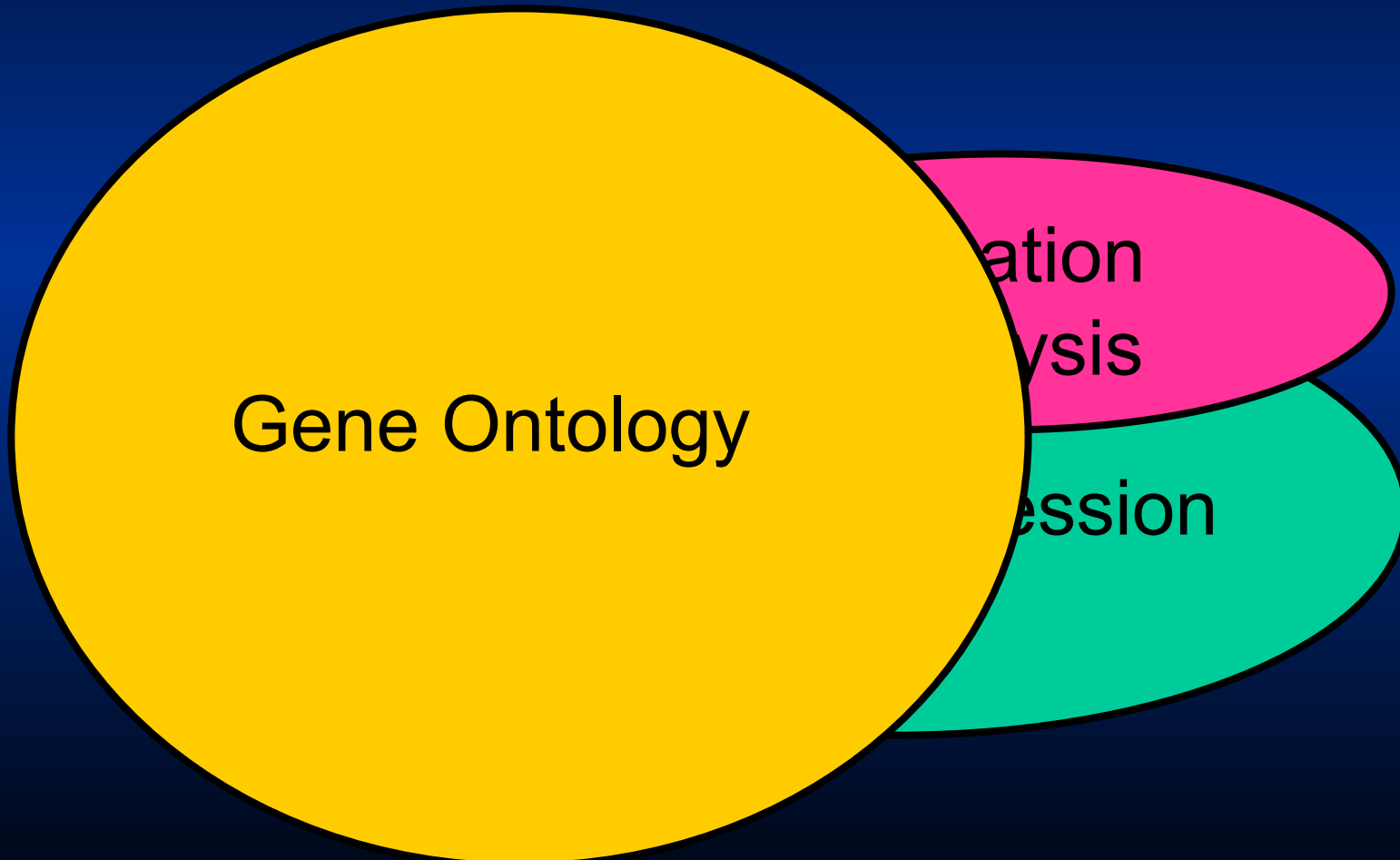


(3) Clustering of Heterogeneous data

Motivation

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

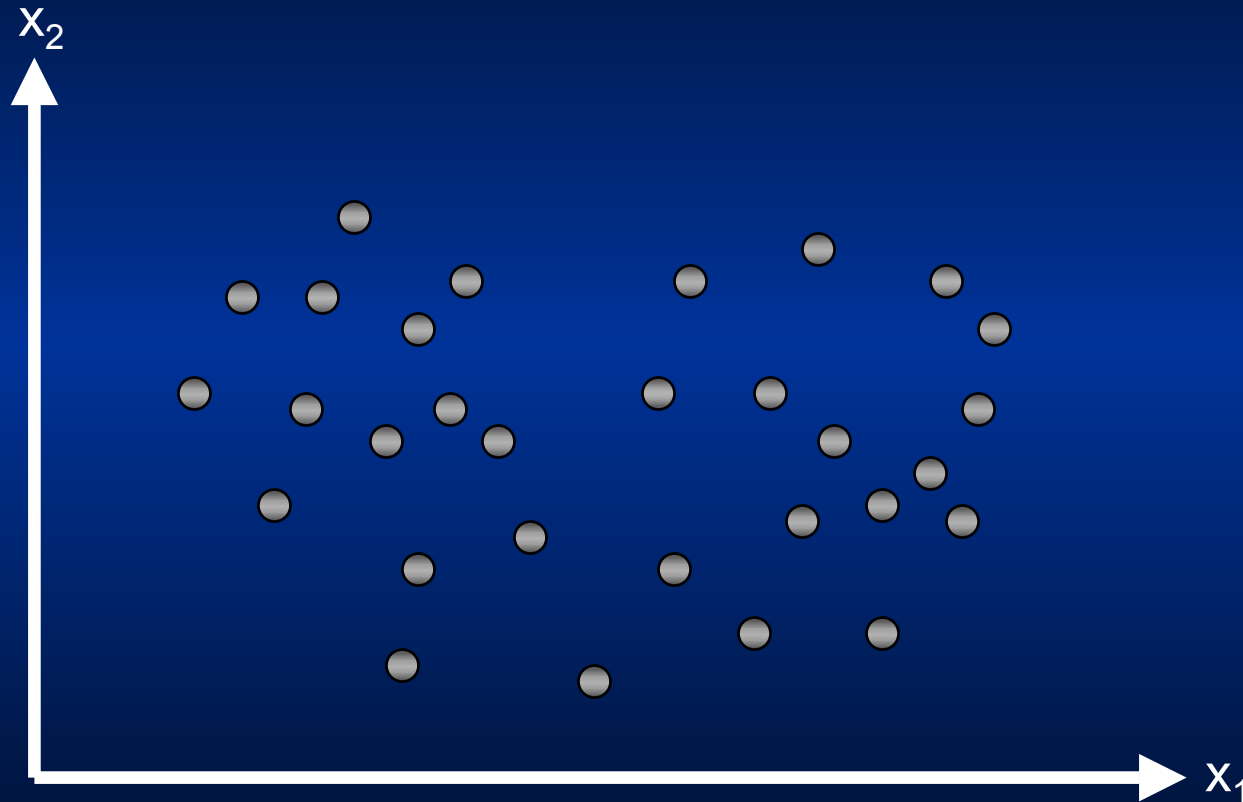
Challenges of Heterogeneous Biological Data



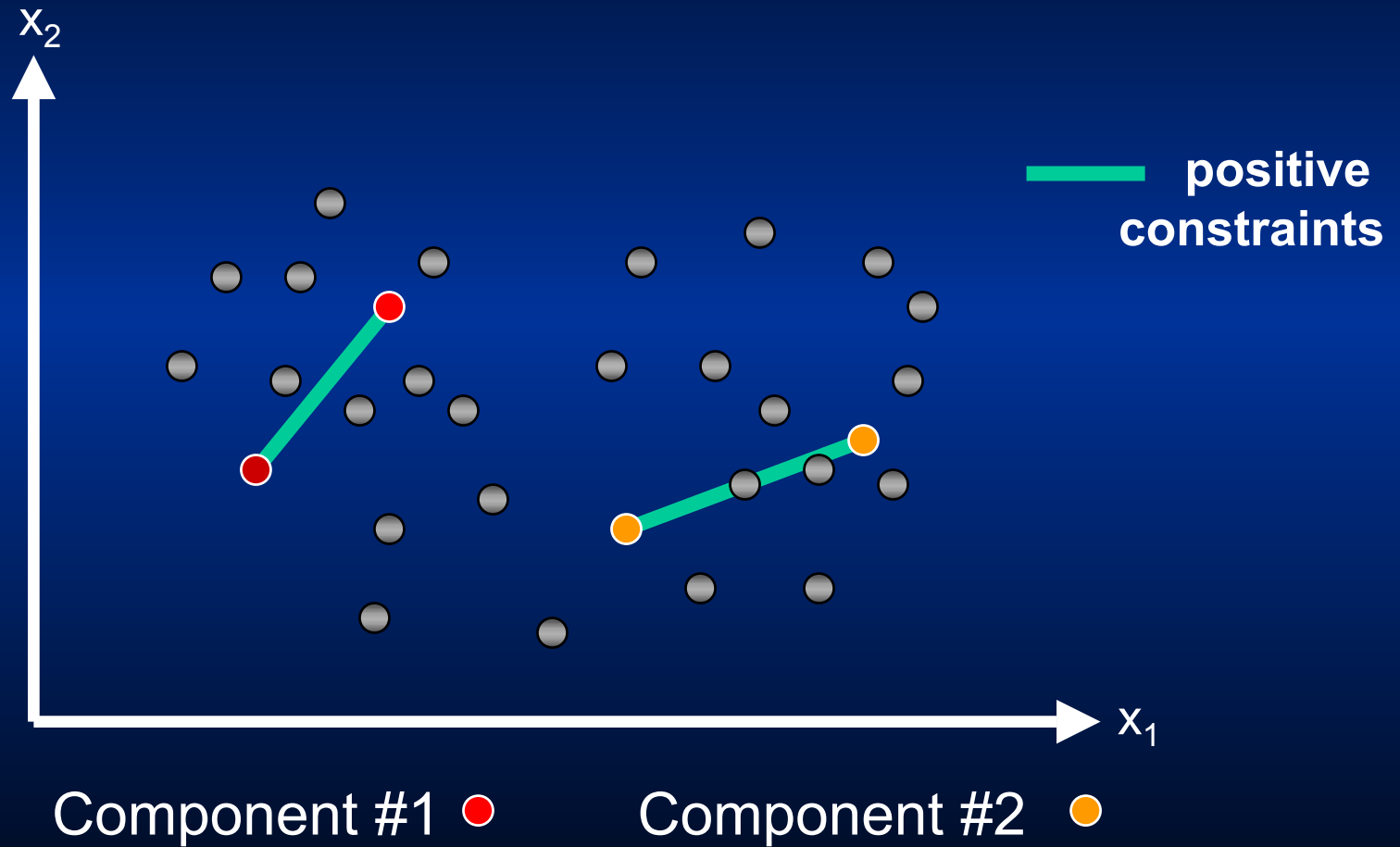
Our Approach

- Semi-supervised learning
 - Encode location analysis as **soft pairwise constraints**
 - Mixture estimation with constraints (Lange *et al.*, 2005, Lu and Leen, 2005)

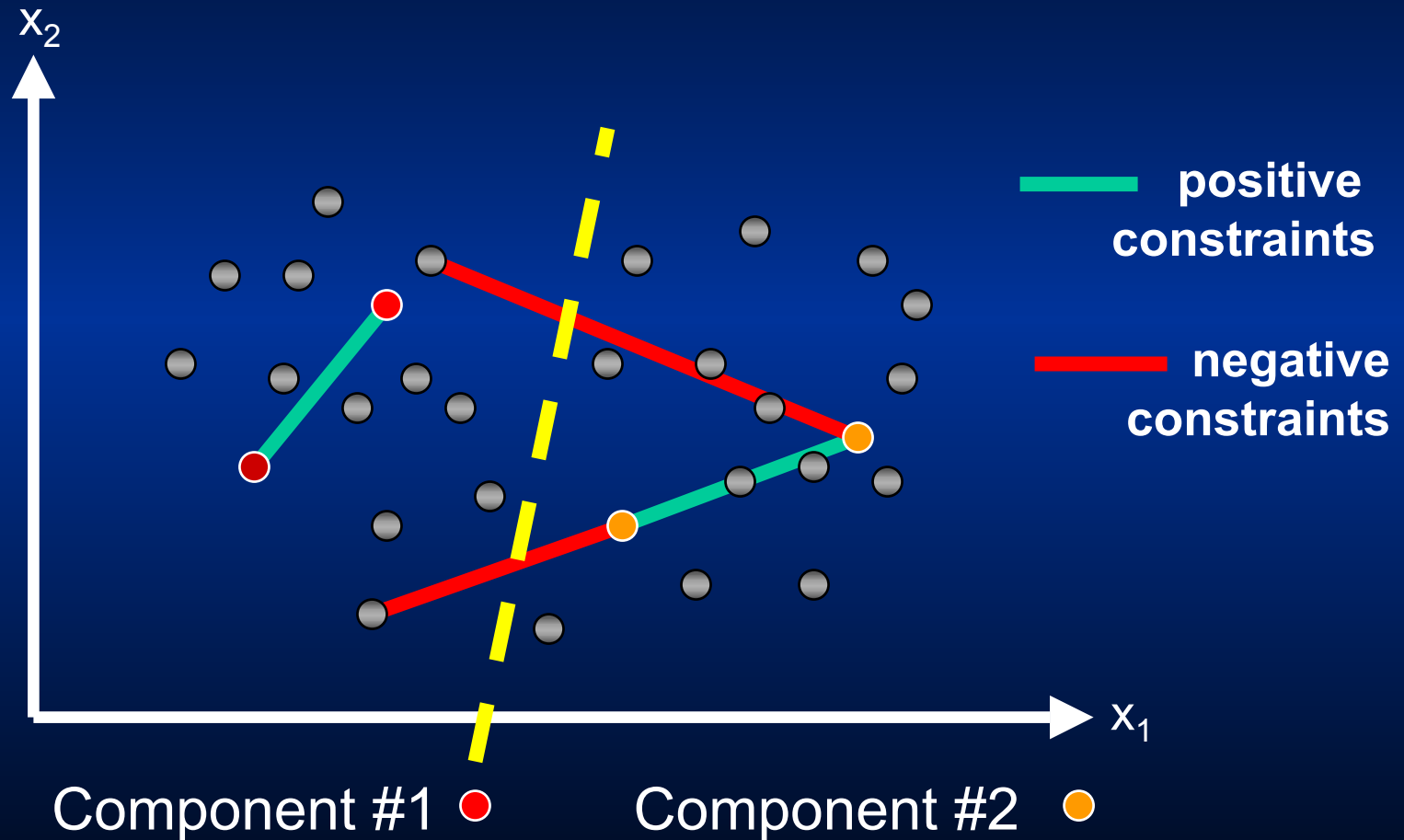
Semi-Supervised Learning



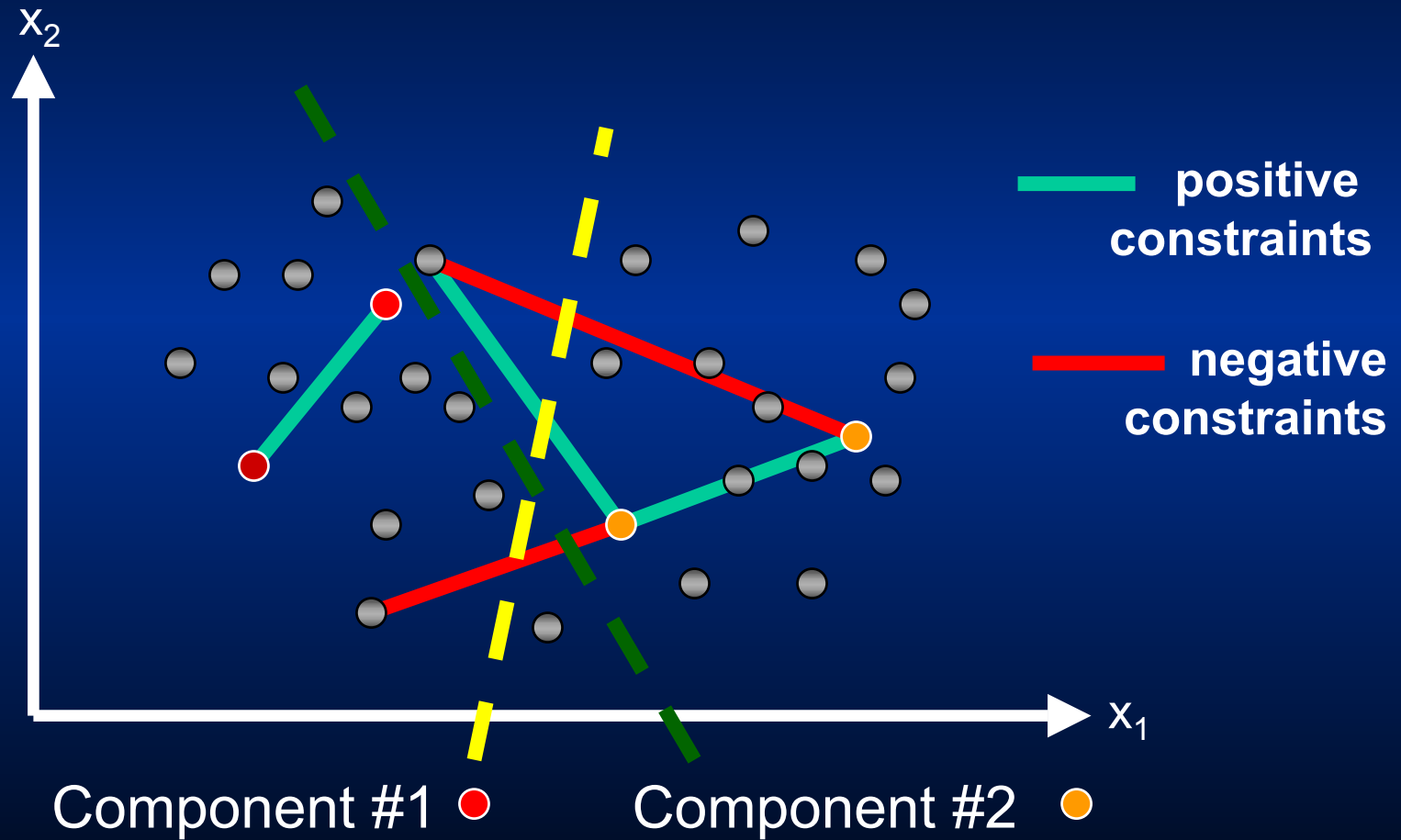
Semi-Supervised Learning



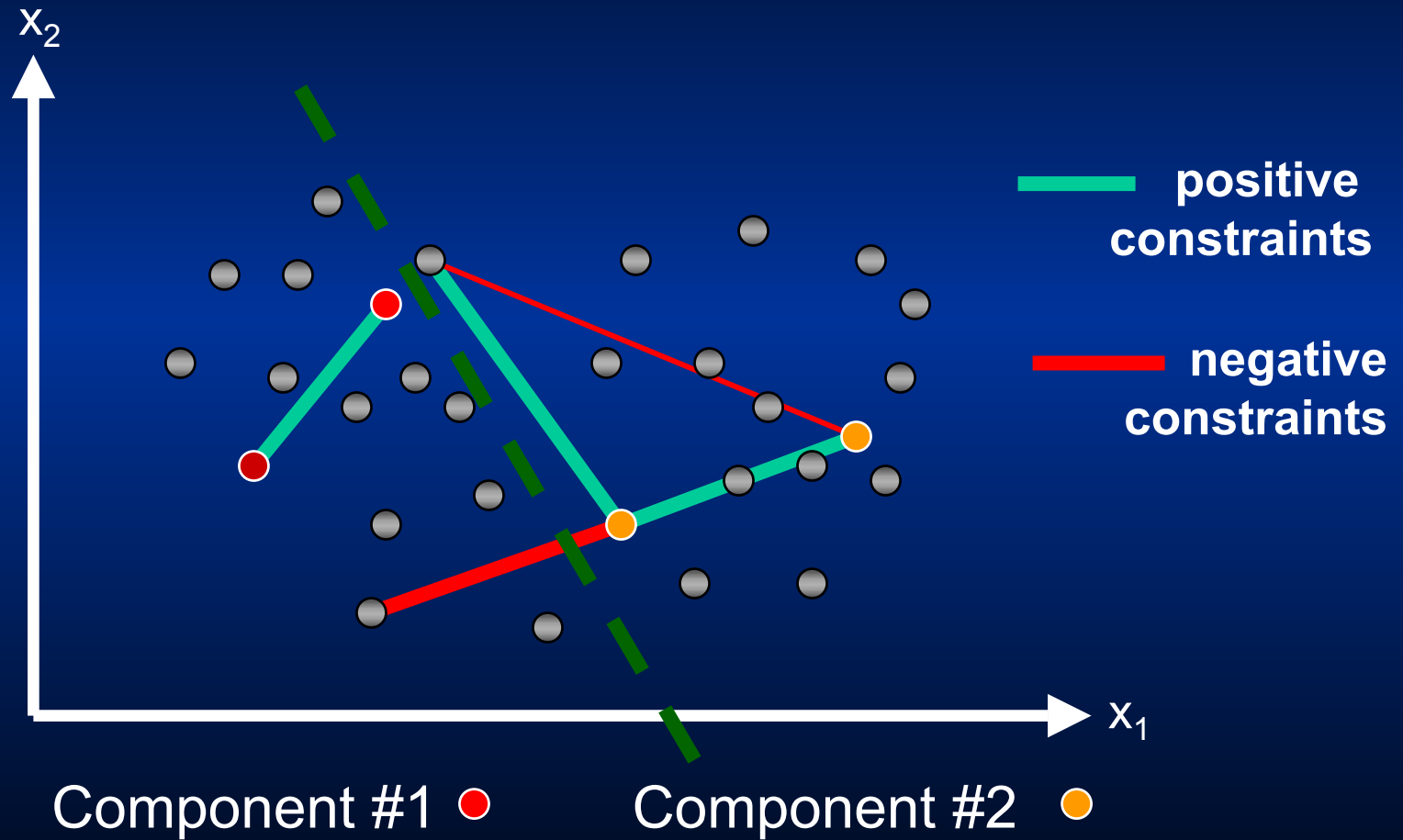
Semi-Supervised Learning



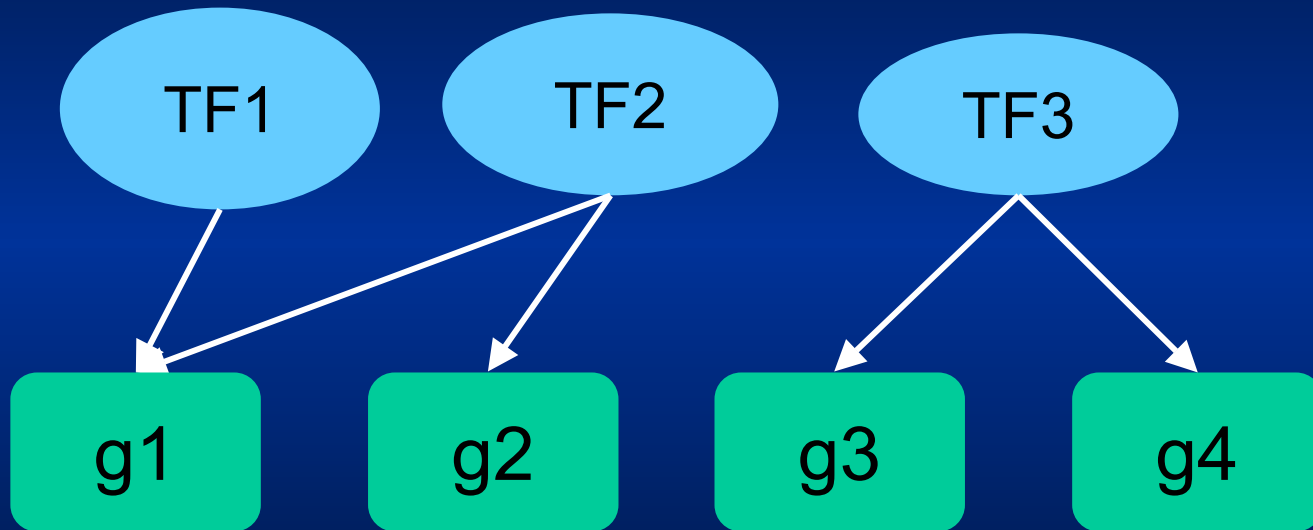
Semi-Supervised Learning



Semi-Supervised Learning



Pairwise Constraints Location Analysis



$$w^+_{(1,2)} = 0.5$$
$$w^-_{(1,2)} = 0.5$$

$$w^+_{(2,3)} = 0.0$$
$$w^-_{(2,3)} = 1.0$$

$$w^+_{(2,3)} = 1.0$$
$$w^-_{(2,3)} = 0.0$$

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)

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

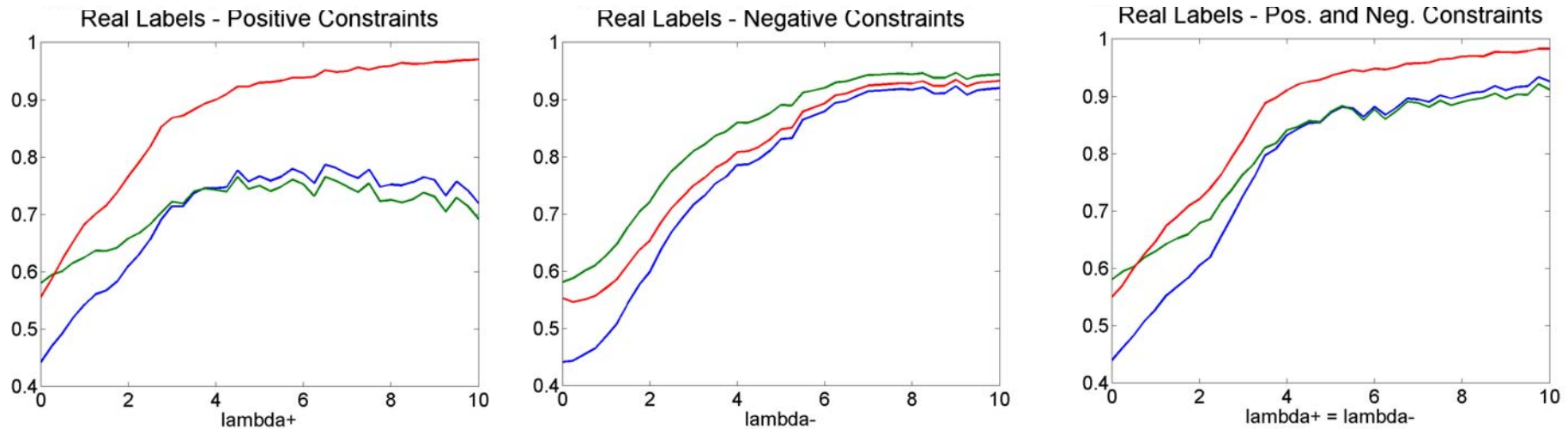
The posterior assignments are approximated by means of Gibbs sampling (Lu and Leen, 2005)

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
 - transcription factor location analysis (Lee, 2002)
 - true labels

Results

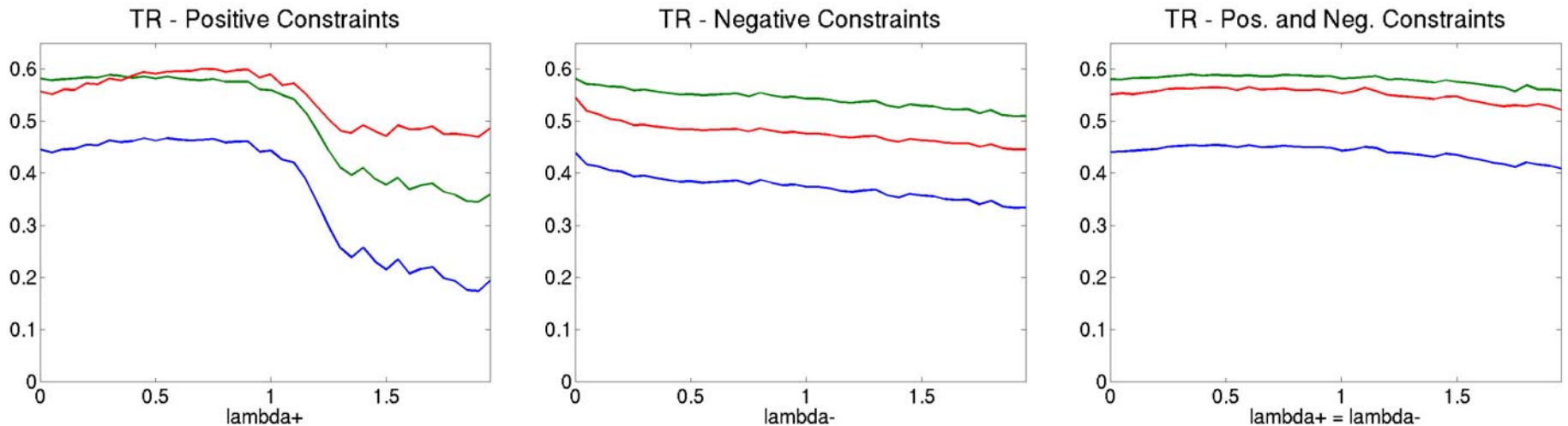
Constraints from True Labels



— corrected Rand
— specificity
— sensitivity

5% of gene
pairs constrained

Constraints from Location Analysis



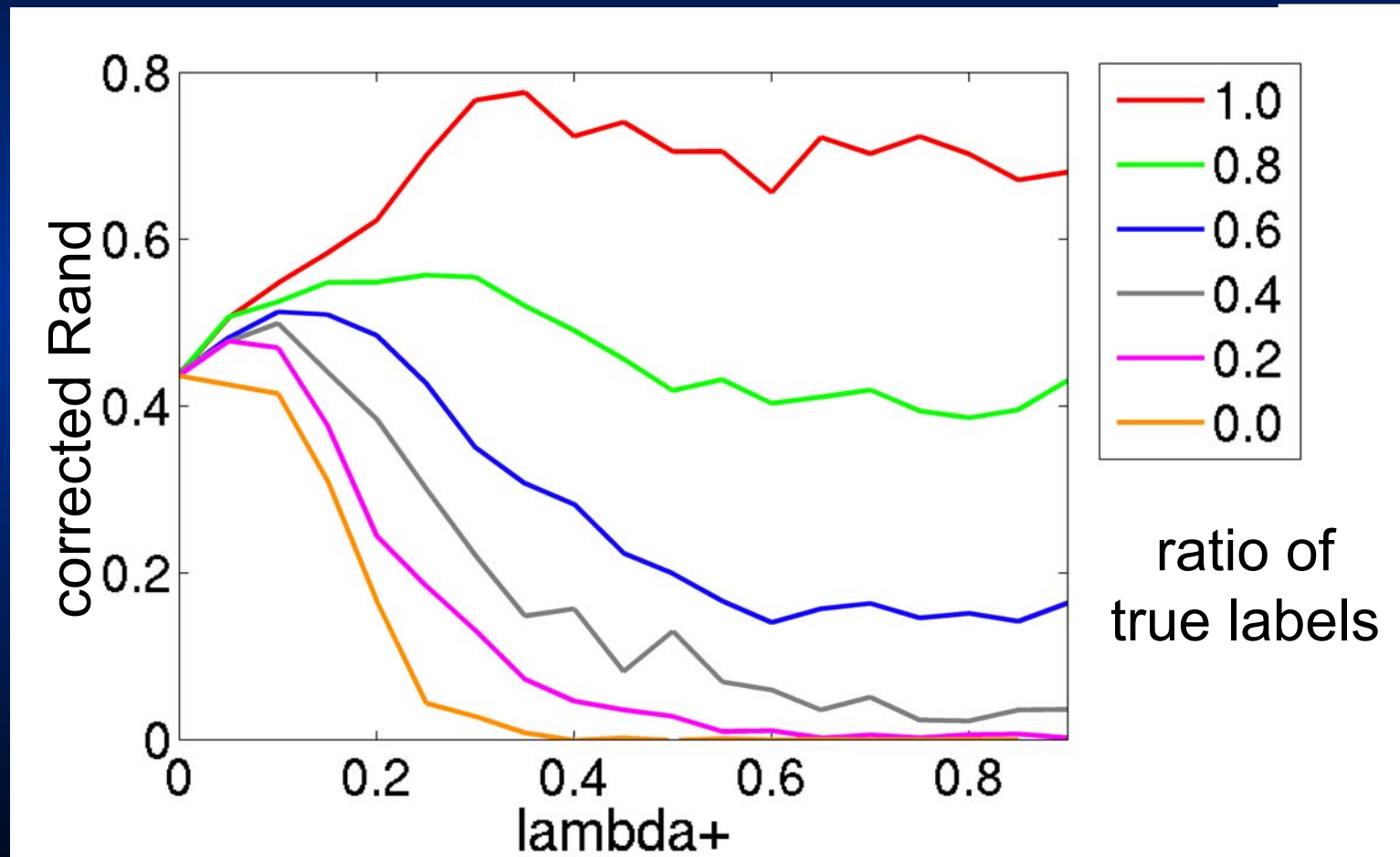
— corrected Rand
— specificity
— sensitivity

40% of gene pairs constrained

Possible Explanations

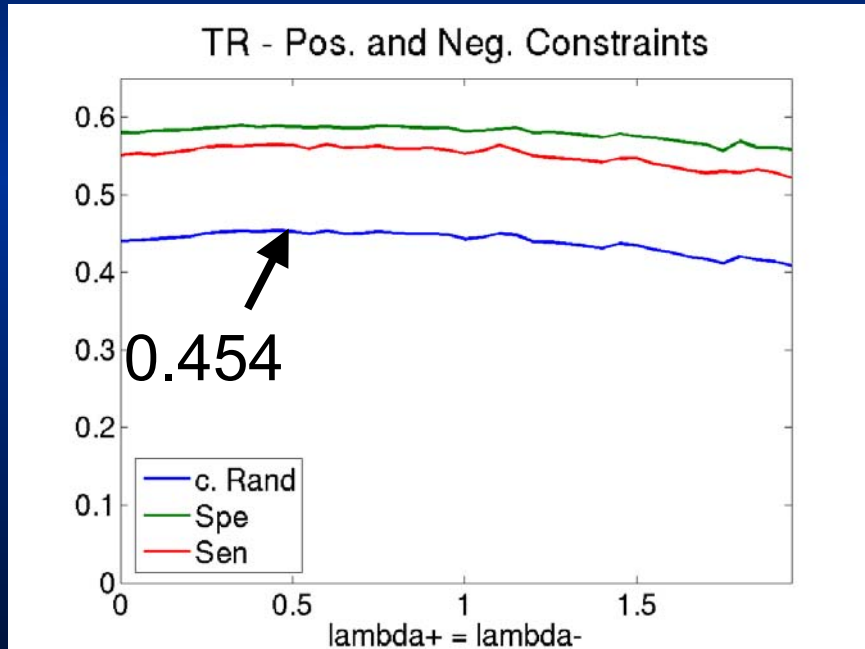
- Non-specific information content
- Noise in the data
- ...

Constraints from True and Random Labels

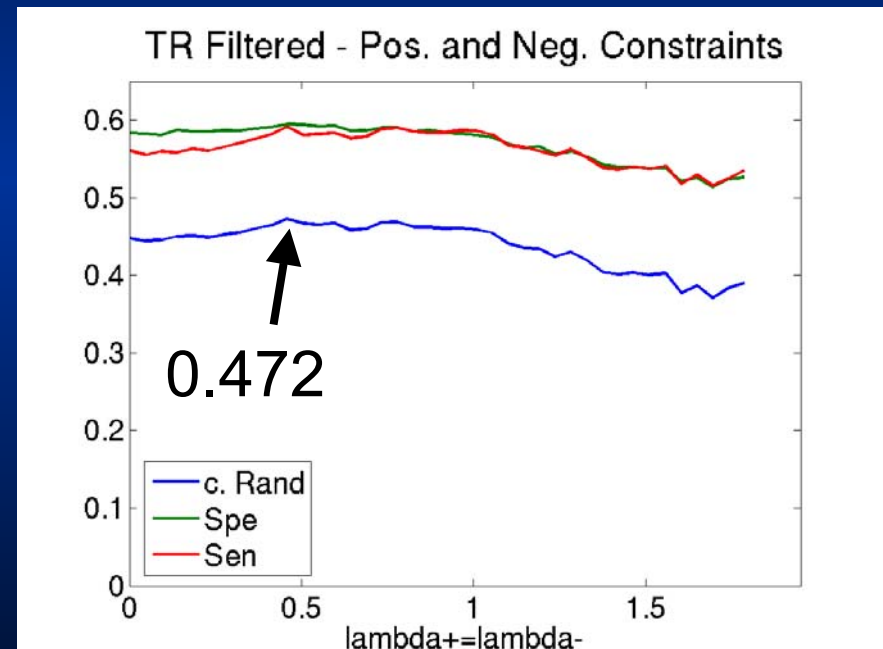


Filtered Constraints from Location Data

Non-filtered



Filtered

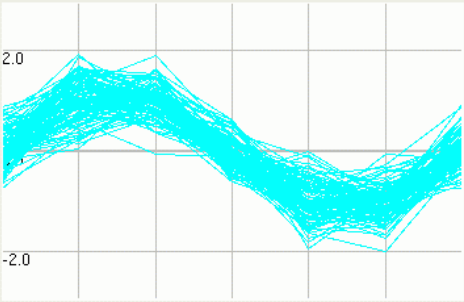


Software GQL

Shell No. 5 - Konsole
GQLCluster <2>

File Filters Models Cluster

Cluster Nr. 1 of size 104 (prior=0.132)

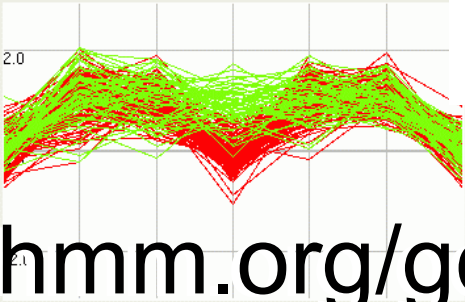


2.0
-2.0

0.05 +/- 0.13 1.05 +/- 0.12 -0.07 +/- 0.10 -1.03 +/- 0.12 -0.00 +/- 0.10
1.08 1.90 1.10 1.91 1.00

Show Details

Cluster Nr. 2 of size 175 (prior=0.184)



2.0
-2.0

0.03 +/- 0.12 1.07 +/- 0.14 0.48 +/- 0.35 1.07 +/- 0.12 -0.00 +/- 0.09
1.05 1.69 1.55 1.71 1.00

Show Details

Cluster Nr. 3 of size 104 (prior=0.132)

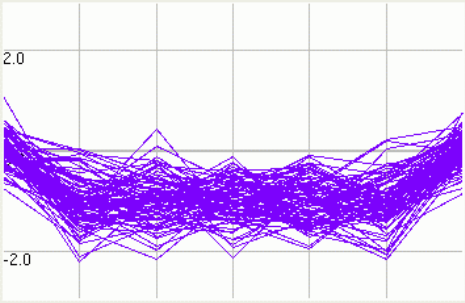
GQLCluster: Details for cluster 3

id=419 p=1.00	This is profile id = #419: Acc#acc419, its
id=417 p=1.00	This is profile id = #417: Acc#acc417, its
id=415 p=1.00	This is profile id = #415: Acc#acc415, its
id=401 p=1.00	This is profile id = #401: Acc#acc401, its
id=379 p=1.00	This is profile id = #379: Acc#acc379, its
id=372 p=1.00	This is profile id = #372: Acc#acc372, its
id=349 p=1.00	This is profile id = #349: Acc#acc349, its
id=396 p=1.00	This is profile id = #396: Acc#acc396, its
id=393 p=1.00	This is profile id = #393: Acc#acc393, its
id=355 p=1.00	This is profile id = #355: Acc#acc355, its
id=374 p=1.00	This is profile id = #374: Acc#acc374, its
id=359 p=1.00	This is profile id = #359: Acc#acc359, its
id=353 p=1.00	This is profile id = #353: Acc#acc353, its
id=409 p=1.00	This is profile id = #409: Acc#acc409, its
id=354 p=1.00	This is profile id = #354: Acc#acc354, its
id=424 p=1.00	This is profile id = #424: Acc#acc424, its
id=404 p=1.00	This is profile id = #404: Acc#acc404, its
id=394 p=1.00	This is profile id = #394: Acc#acc394, its
id=358 p=1.00	This is profile id = #358: Acc#acc358, its

-0.03 +/- 0.12 0.13
1.00

Show Details

Cluster Nr. 4 of size 88 (prior=0.111)



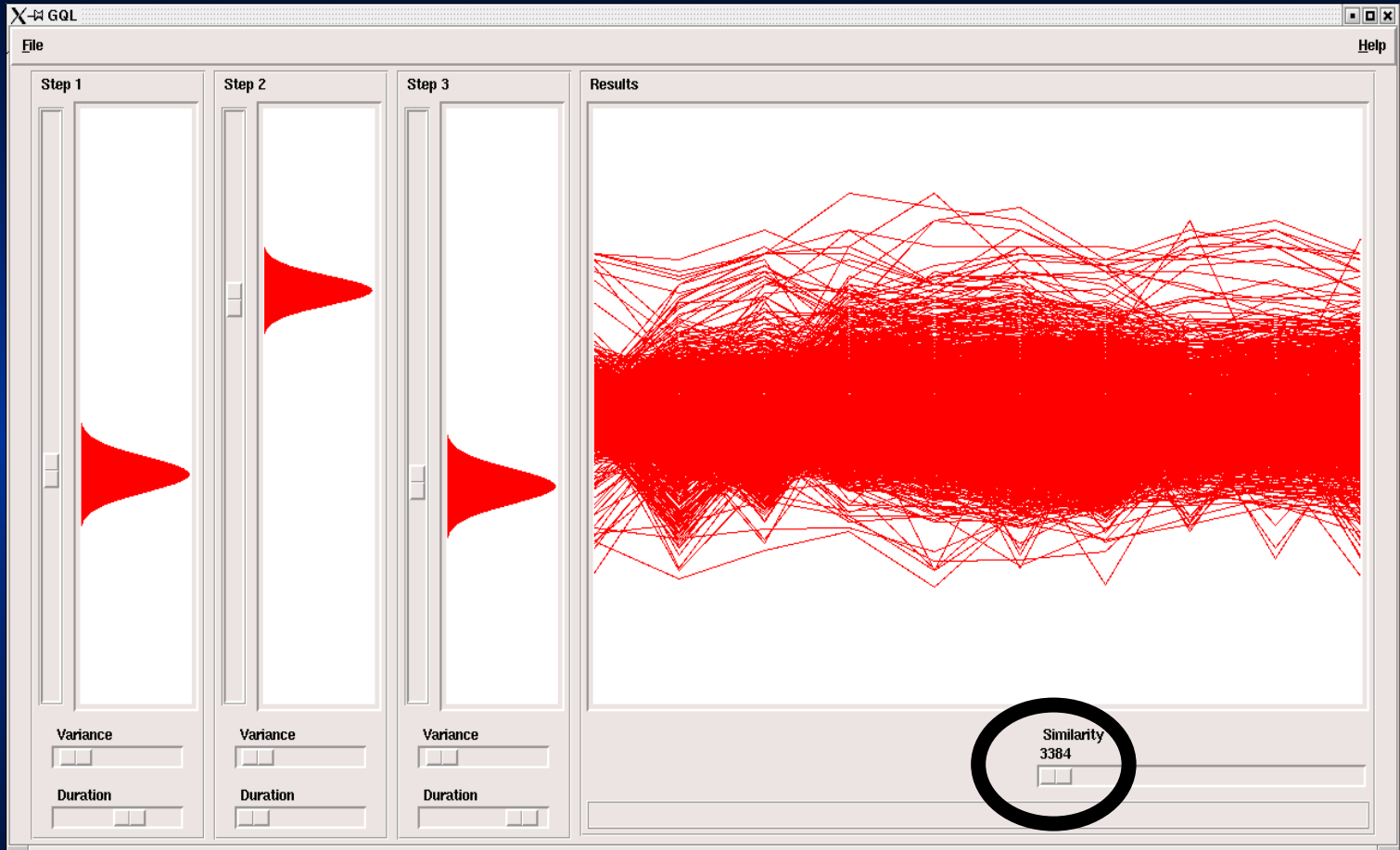
2.0
-2.0

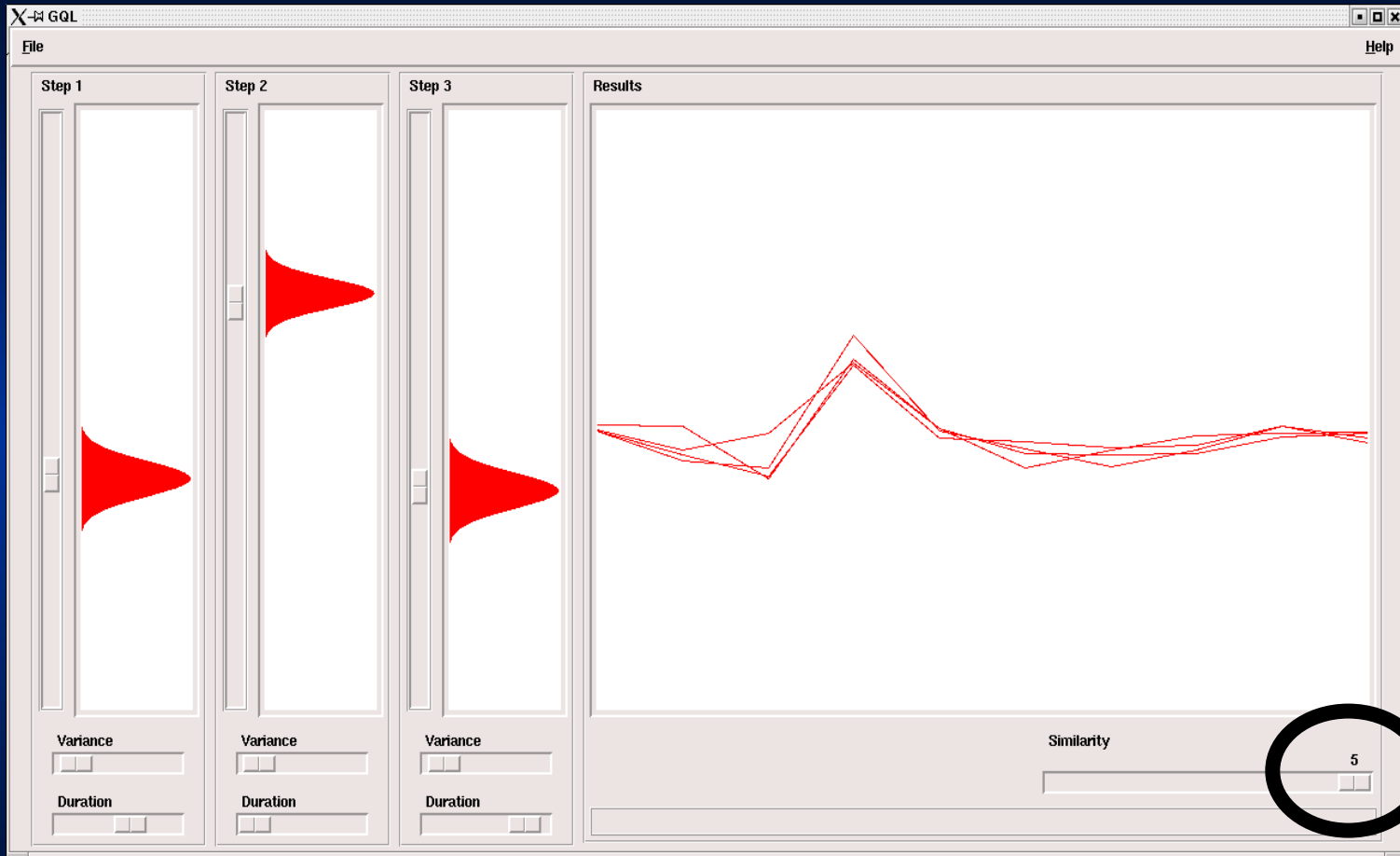
0.03 +/- 0.10 -1.01 +/- 0.23 -0.96 +/- 0.18 -1.04 +/- 0.19 0.00 +/- 0.10
1.00 1.65 1.66 1.65 1.04

Show Details

QuickTime™ and a TIFF (Uncompressed) decompressor are needed to see this picture.

<http://ghmm.org/gql>





Summary

- Clustering of Time-Series
 - Flexible: *cyclic & transient time-courses*
 - Interactive & robust
- Cluster Validation
 - Methodology for evaluating clustering given functional annotation
- Heterogeneous Analysis
 - Successful integration of location analysis

Outlook

- Cluster Validation
 - Perform a ‘extensive’ evaluation of clustering methods
- Heterogeneous Analysis
 - learning of relevant constraints
 - *In-situ hybridization, protein-protein interactions,*
...
- Clustering of Development Trees
 - In progress ...

Acknowledgements

- *A. Schönhuth* and *C. Steinhoff*
- GHMM (<http://ghmm.org>):
Wasinee Rungsaityotin, Benjamin Georgi, Ben Rich, Matthias Heinig, Alexander Riemer, Janne Grunau
- *T. Beissbarth* for help with *GO*

Thanks.

ghmm.org/gql

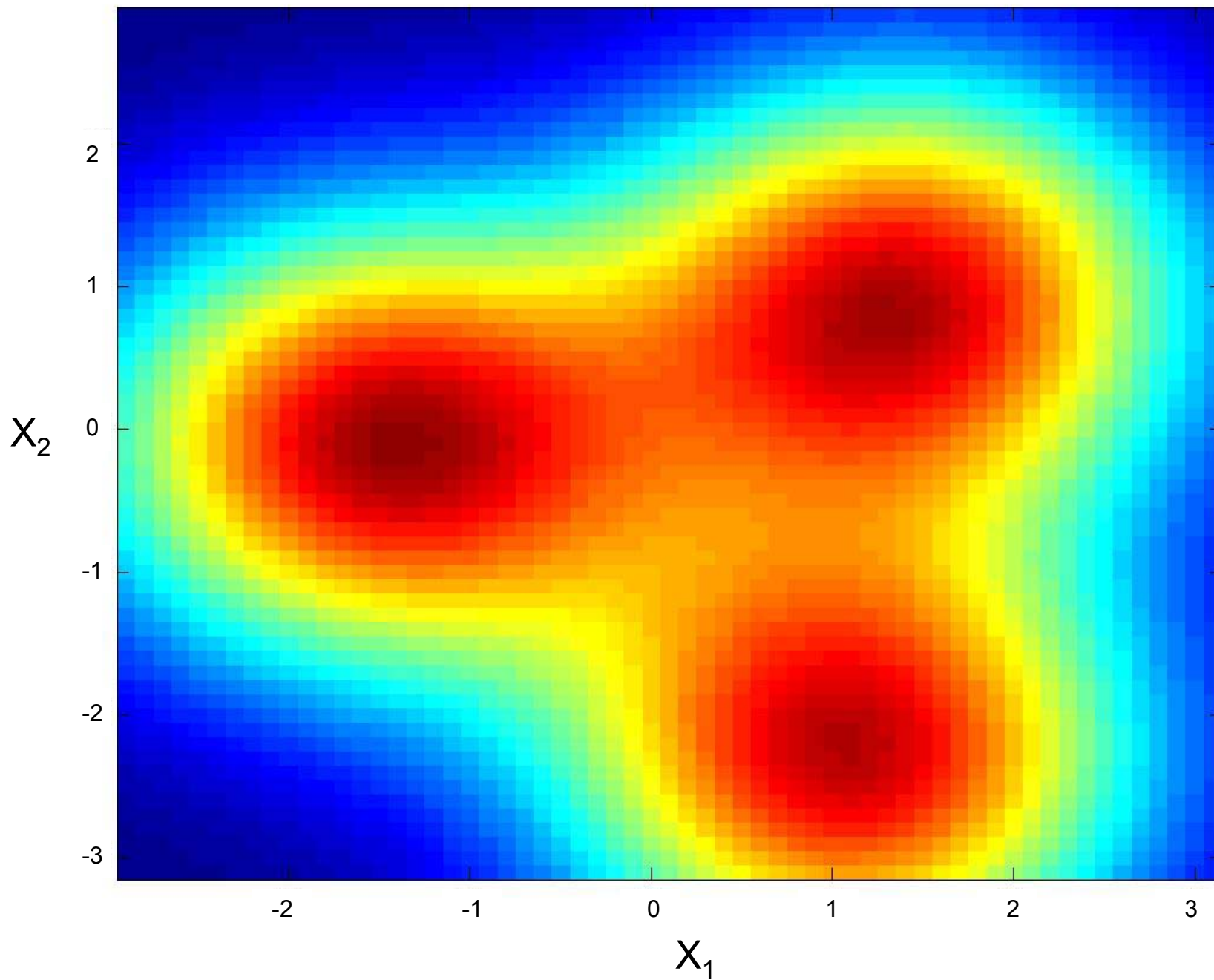
Motivational Buzzphrases

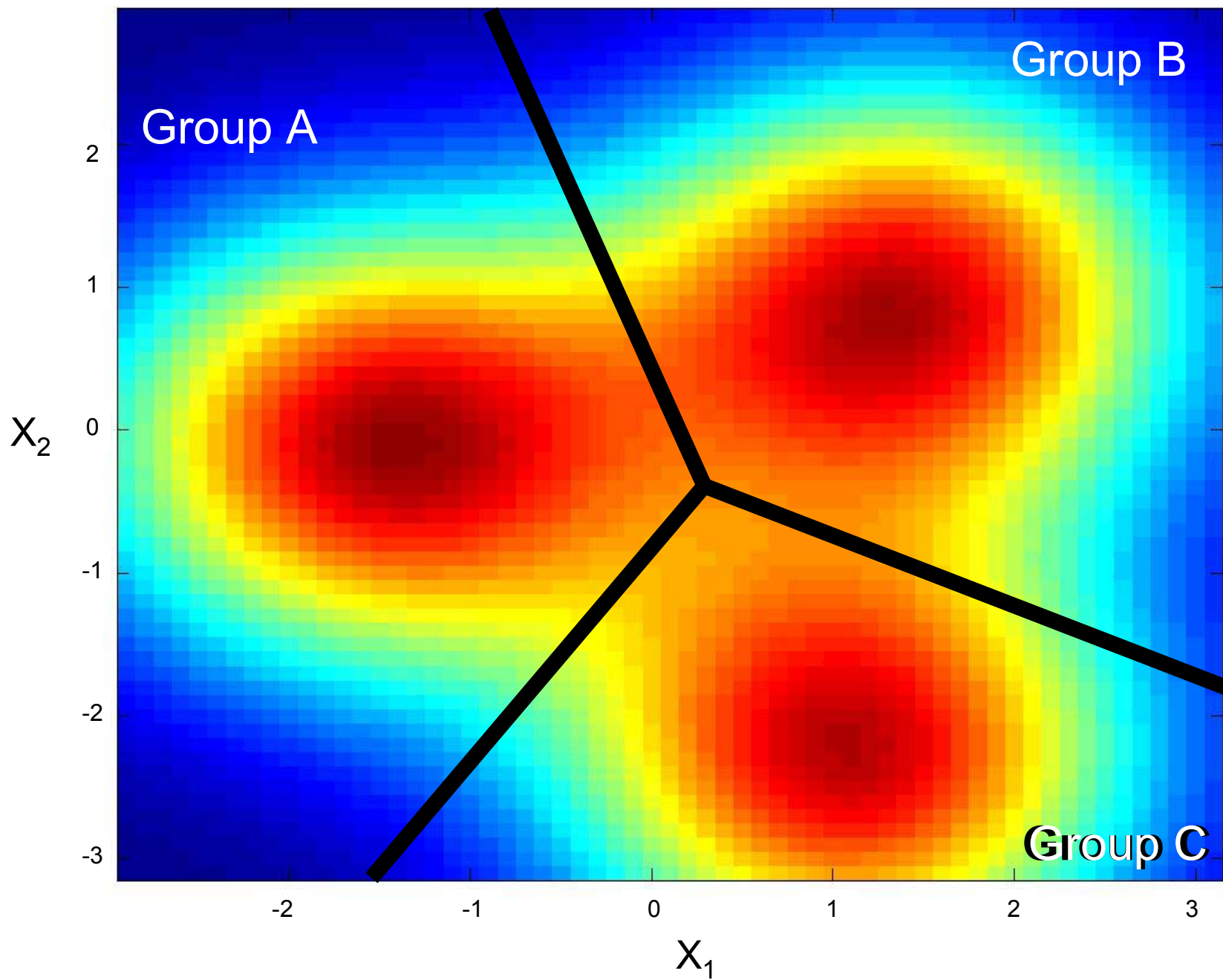
Model qualitatively

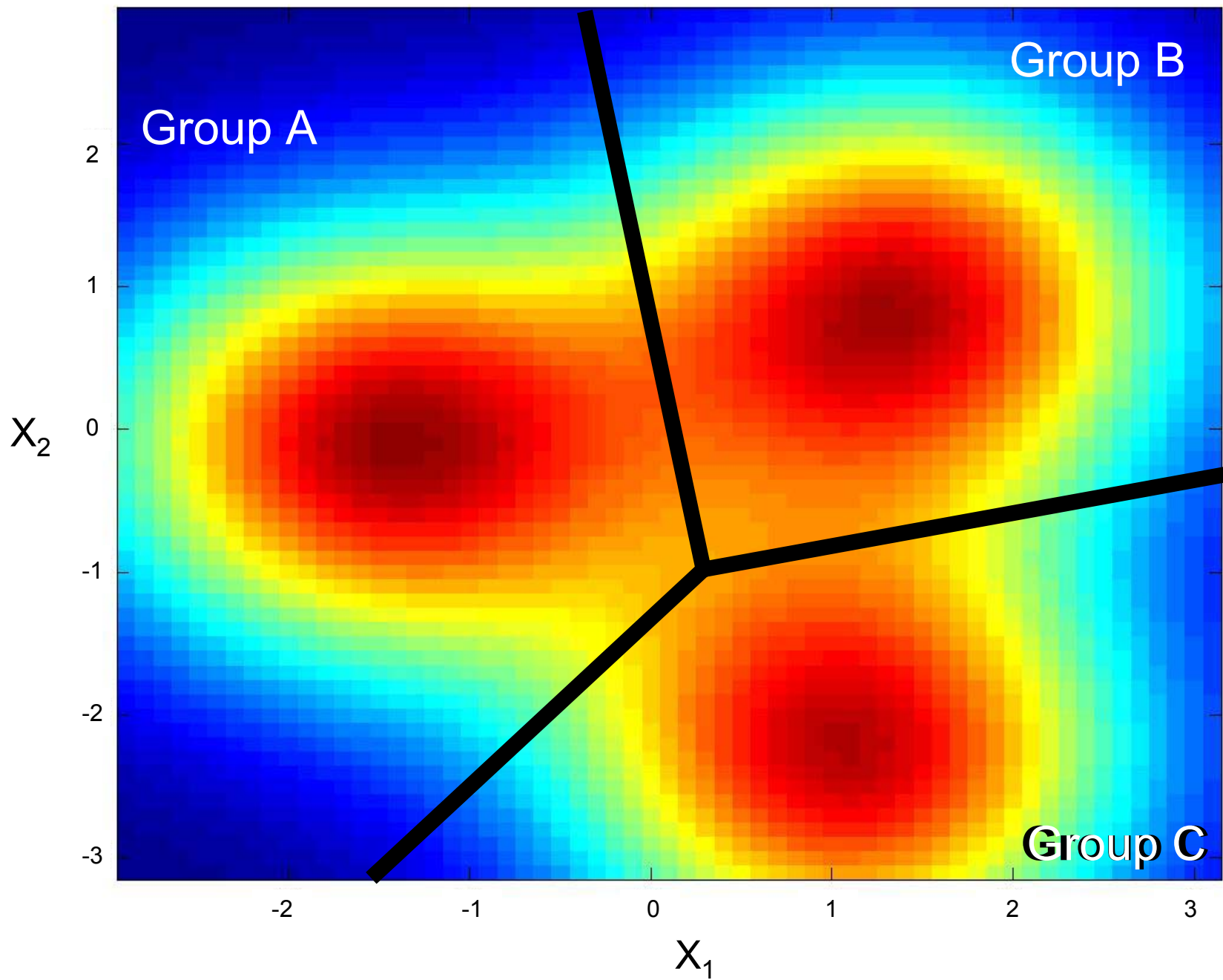
Embrace ambiguity

Assure robustness

Don't be ignorant







Yeast Cell Cycle

