## V14 – Gene Regulation

### Thu, Dec. 5, 2019

- Co-expression modules
  - Motifs in GRNs
- Master Regulatory Genes in GRNs



Observed modules

# **Module detection methods**

Module detection is a cornerstone also in the biological interpretation of large gene expression compendia.

Such modules are groups of genes with similar expression profiles, which also tend to be functionally related and co-regulated.

Approaches:

- (a) clustering
- (b) decomposition methods
- (c) biclustering local co-expression (also (b))
- (d) direct network inference
- (e) iterative network inference.

(d) and (e) also model the regulatory relationships between the genes (see VI3, here used: GENIE3 and CLR)

Saelens et al. Nature Commun. 9, 1090 (2018)

# **Limitations of clustering**

(1) Clustering methods look at co-expression among all samples.
 As transcriptional regulation is highly context specific, clustering potentially misses local co-expression effects which are present in only a subset of all biological samples.

(2) Most clustering methods are **unable to assign genes to multiple modules**. The issue of overlap between modules is especially problematic given the increasing evidence that gene regulation is highly combinatorial and that gene products can participate in multiple pathways.

(3) Clustering methods **ignore the regulatory relationships between genes**. As the variation in target gene expression can at least be partly explained by variation in transcription factor expression, including this information could therefore boost module detection.

## **Other module-detection methods**

Decomposition methods and biclustering try to handle local co-expression and overlap.

These methods differ from clustering because they allow that genes within a module do **not need to be co-expressed** in all biological samples, but that a sample can influence the expression of a module to a certain degree (decomposition methods) or not at all (biclustering methods).

Alternatively, direct network inference (direct NI) and iterative NI use the expression data to additionally model the **regulatory relationships** between the genes.

# **Module detection methods**



Saelens et al. Nature Commun. 9, 1090 (2018)

# Module detection methods: performance



ICA-based decomposition methods work best in detecting co-expression modules that overlap with known regulatory modules.

Saelens et al. Nature Commun. 9, 1090 (2018)

ICA method: https://www.ece.ucsb.edu/wcsl/courses/ECE594/594C\_FI0Madhow/comon94.pdf https://www.cs.helsinki.fi/u/ahyvarin/papers/NN00new.pdf

# Independent Component Analysis (ICA)

ICA decomposes the expression data matrix X into a number of "components" (k = 1,2,..K).

Each component is characterized by an **activation pattern** over genes  $(S_k)$  and another over samples  $(A_k)$ 

$$X = \sum_{k=1}^{K} S_k \otimes A_k + E$$

This is done in such a way that the gene activation patterns  $(S_1, S_2, ..., S_K)$  are as **statistically as independent as possible** while also the residual "error" matrix *E* is minimized.

K is the number of inferred independent components (IC) to which pathways and regulatory modules map.

# **ICA model of gene expression**

**EXPRESSION** 

Samples

pathway - B

pathway-C

pathway - A

Genes

pathway - D

In the ICA model, the gene expression matrix is decomposed into the product of a "source" matrix S and a "mixing" matrix A.

#### **Color codes** for heatmaps:

red, overexpression; green, underexpression; blue, upregulation; yellow, downregulation.

The **columns of S** describe the activation levels of genes in the various inferred independent components.

SOURCE (S)

Genes

Х

MIXING (A)

Samples

The **rows of A** give the activation levels of the independent components across tumor samples.

# **ICA model of gene expression**



In Principal Component Analysis (PCA), the PCs are constructed to be orthonormal to each other. In ICA, the ICs are constructed to be statistically independent from each other.

Module: co-active set of genes in one IC (e.g. they have a high z-score compared to the other genes in this IC). In IC-2, they are either on (right) or off (left).



The *IC*-*k*-submatrix is obtained by multiplying the *k*-th column of S,  $S_k$ , with the *k*-th row of A,  $A_k$ .

ICA algorithms seem to be much **more robust** than PCA-based methods: pathways that were found to be differentially activated through ICA in one cohort were also consistently differentially activated in the other cohorts.

Also, whereas using PCA no regulatory module was found to be differentially activated across 4 major breast cancer studies, the ICA algorithms found an average of 4 modules.

## **9.5 Network Motifs**



Shai S. Shen-Orr<sup>1</sup>, Ron Milo<sup>2</sup>, Shmoolik Mangan<sup>1</sup> & Uri Alon<sup>1,2</sup>

*Nature Genetics* **31** (2002) 64

RegulonDB + hand-curated literature evidence

- $\rightarrow$  break down network into motifs
  - $\rightarrow$  statistical significance of the motifs?
    - $\rightarrow$  behavior of the motifs <=> location in the network?

# **Detection of motifs**

Represent transcriptional network as a connectivity matrix M such that  $M_{ij} = 1$  if operon j encodes a TF that transcriptionally regulates operon i and  $M_{ij} = 0$  otherwise.

Scan all  $n \times n$  submatrices of M generated by choosing n nodes that lie in a connected graph, for n = 3 and n = 4.

Submatrices were enumerated efficiently by recursively searching for nonzero elements.



Connectivity matrix for causal regulation of transcription factor *j* (row) by transcription factor *i* (column). Dark fields indicate regulation. (Left) **Feed-forward loop** motif. TF 2 regulates TFs 3 and 6, and TF 3 again regulates TF 6. (Middle) Single-input multiple-output (**SIM**) motif. (Right) Densely-overlapping region.

For n = 3, the only significant motif is the **feedforward loop**.

Also significant are SIMs and densely overlapping regulation motifs..

Shen-Orr et al. Nature Gen. 31, 64 (2002)

# **Motif Statistics**

Compute a p-value for submatrices representing each type of connected subgraph by comparing # of times they appear in real network vs. in random network.

Structure	Appearances in real network	Appearances in randomized network (mean ± s.d.)	<i>P</i> value
Coherent feedforward loop	34	4.4 ± 3	<i>P</i> < 0.001
Incoherent feedforward loop	6	2.5 ± 2	<i>P</i> ~ 0.03
Operons controlled by SIM (>13 operons)	68	28 ± 7	<i>P</i> < 0.01
Pairs of operons regulated by same two transcription factors	203	57 ± 14	<i>P</i> < 0.001
Nodes that participate in cycles*	0	$0.18 \pm 0.6$	<i>P</i> ~ 0.8

Listed motifs are highly **overrepresented** compared to randomized networks

**No cycles**  $(X \rightarrow Y \rightarrow Z \rightarrow X)$  were identified,

but this was not statistically significant in

comparison to random networks

Bioinformatics 3 – WS 19/20

## **Generate Random Networks**

For a stringent comparison to randomized networks, one generates networks with precisely the same

- number of operons,
- interactions,
- TFs and
- number of incoming and outgoing edges for each node as in the real network (here the one from *E. coli*).

One starts with the real network and repeatedly swaps randomly chosen pairs of connections ( $X1 \rightarrow Y1$ ,  $X2 \rightarrow Y2$  is replaced by  $X1 \rightarrow Y2$ ,  $X2 \rightarrow Y1$ ) until the network is well randomized.

## **Generate Random Networks**

This yields networks with precisely the **same number of nodes** with *p* incoming and *q* outgoing nodes as the real network.

The corresponding randomized connectivity matrices, *Mrand*, have the same number of nonzero elements in each row and column as the corresponding row and column of the real connectivity matrix *M*:

$$\sum_{i} Mrand_{ij} = \sum_{i} M_{ij} \text{ and } \sum_{j} Mrand_{ij} = \sum_{j} M_{ij}$$

# **FFL dynamics**



In a **coherent** FFL: X **and** Y activate Z

#### Dynamics:

- input activates X
- X activates Y (delay)
- (X &&Y) activates Z

Delay between X and Y  $\rightarrow$  signal must persist longer than delay

#### (see lecture 12, slide 29)

- $\rightarrow$  reject transient signal, react only to **persistent** signals
- $\rightarrow$  enables fast shutdown

#### Helps with **decisions** based on **fluctuating signals.**

# Motif 2: Single-Input-Module



Set of operons controlled by a single transcription factor

- same sign
- no additional regulation
- control is usually autoregulatory (70% vs. 50% overall)

Example for this in *E. coli*: arginine biosynthetic operon *argCBH* plus other enzymes of arginine biosynthesis pathway.

Mainly found in genes that code for **parts** of a protein **complex** or metabolic **pathway** 

 $\rightarrow$  produces components in comparable amounts (stoichiometries).

# **SIM-Dynamics**



If different thresholds exist for each regulated operon:

- $\rightarrow$  first gene that is activated is the last one that is deactivated
  - $\rightarrow$  well defined temporal ordering (e.g. flagella synthesis) + stoichiometries

# Motif 3: Densely Overlapping Regulon



Dense layer between groups of TFs and operons

→ much denser than network average (≈ community)

Usually each operon is regulated by a different combination of TFs.

Sometimes: same set of TFs for group of operons  $\rightarrow$  "multiple input module"

# **Network with Motifs**



# Regulatory network rewiring in S. cerevisiae

**a**, Schematics and summary of properties for the endogenous and exogenous sub-networks.

b, Graphs of the static and condition-specific networks. Transcription factors and target genes are shown as nodes in the upper and lower sections of each graph respectively, and regulatory interactions are drawn as edges; they are coloured by the number of conditions in which they are active. Different conditions use distinct sections of the network.



**c**, Standard statistics (global topological measures and local network motifs) describing network structures. These vary between endogenous and exogenous conditions; those that are high compared with other conditions are shaded. (Note, the graph for the static state displays only sections that are active in at least one condition, but the table provides statistics for the entire network including inactive regions.)

Luscombe, Babu, ... Teichmann, Gerstein, Nature 431, 308 (2004)

# **Identification of Master regulatory genes**



A vertex *u* **dominates** another vertex *v* if there exists a directed arc (*u*,*v*).

<u>Idea</u>: find a **set of dominator nodes** of minimum size that controls all other vertices.

In the case of a GRN, a directed arc symbolizes that a transcription factor regulates a target gene.

In the figure, the MDS nodes  $\{A, B\}$  are the dominators of the network. Together, they regulate all other nodes of the network (*C*, *E*, *D*).

### **Identification of Master regulatory genes**



Core pluripotency network, Kim et al. Cell (2008)

The nodes of a MDS can be spread as isolates nodes over the entire graph. However, e.g. the set of core pluripotency factors is tightly connected (right).

#### Next idea: find a connected dominating set of minimum size (MCDS).

(Left) the respective set of MCDS nodes (*black and gray*). Here, node *C* is added in order to preserve the connection between the two dominators *A* and *B* to form an MCDS.

# **ILP for minimum dominating set**

<u>Aim</u>: we want to determine a set D of minimum cardinality such that for each  $v \in V$ , we have that  $v \in D$  or that there is a node  $u \in D$  and an arc  $(u,v) \in E$ .

Let  $\delta^{-}(v)$  be the set of incoming nodes of v such that  $(u,v) \in E$ ,  $x_u$  and  $x_v$  are binary variables associated with u and v.

We select a node *v* as dominator if its binary variable  $x_v$  has value 1, otherwise we do not select it.

$$\begin{array}{ll} \text{minimize} & \sum_{\nu \in V} x_{\nu} \\ \text{subject to} & x_{u} + \sum_{\nu \in \delta^{-}(u)} x_{\nu} \geq 1 \quad \forall u \in V \\ & x_{\nu} \in \{0, 1\} \qquad \quad \forall \nu \in V \end{array}$$

With the GLPK solver, the runtime was less than 1 min for all considered networks.

Nazarieh et al. BMC Syst Biol 10:88 (2016)

V I4 – 24

# ILP for minimum connected dominating set

A minimum connected dominating set (MCDS) for a directed graph G = (V,E) is a set of nodes  $D \subseteq V$  of minimum cardinality that is a dominating set and additionally has the property that the graph G[D] induced by D is **weakly connected**, i.e. such that in the underlying undirected graph there exists a path between any two nodes of D that only uses vertices in D.

This time we will use two binary valued variables  $y_v$  and  $x_e$ .  $y_v$  indicates whether node v is selected to belong to the MCDS.  $x_e$  for the edges then yields a tree that contains all selected vertices and no vertex that was not selected.

minimize  $\sum_{\nu \in V} y_{\nu}$ subject to  $\sum_{e \in E} x_e = \sum_{i \in V} y_i - 1$ 

This guarantees that the number of edges is one less than the number of vertices. This is necessary (but not sufficient) to form a (spanning) tree.

# ILP for minimum connected dominating set

minimize

 $\sum y_{\nu}$ 

subject to

$$\sum_{e \in E} x_e = \sum_{i \in V} y_i - 1$$
$$\sum_{e \in E(S)} x_e \le \sum_{i \in S \setminus \{j\}} y_i \quad \forall S \subset V, \forall j \in S$$

#### Second constraint

#### $\rightarrow$ selected edges imply a **tree**.

(Note that this defines an exponential number of constraints for all subgraphs of V!)

$$y_{u} + \sum_{v \in \delta^{-}(u)} y_{v} \ge 1 \qquad \forall u \in V$$
$$y_{v} \in \{0, 1\} \qquad \forall v \in V$$

 $x_e \in \{0, 1\} \qquad \qquad \forall e \in E$ 

Third constraint

 $\rightarrow$  node set forms **dominating set**.

For dense graphs, this yields a quick solution.

However, for sparse graphs, the running time may be considerable.

The second constraint was implemented by an iterative approach.

### **Example MDS**



(Left) this toy network includes 14 nodes and 14 edges.

(Right) The dark colored nodes {*J*, *B*, *C*, *H*, *L*} are the dominators of the network obtained by computing a MDS.

## **Example MCDS**



(Left) The nodes colored blue make up the **largest connected component** (LCC) of the underlying **undirected** graph.

(Right) MCDS nodes for this component are {J, D, B, C, G, H}.

### **Example MCDS**



(Left) The green colored nodes are elements of **the largest connected component** underlying the **directed** graph.

(Right) The two nodes  $\{B, C\}$  form the MCDS for this component.

# **MCDS of the strongly connected component**



Strongly connected component of directed path:

there is a directed path (all edges oriented into the same direction) between all pairs of vertices

Here, there are no directed paths from (G or F) to (E or D).

(Left) The nodes colored orange show the LSCC in the network.

(Right) The node A is the only element of the MCDS

# Studied networks: RegulonDB (E.coli)

This GRN contains 1807 genes, including 202 TFs and 4061 regulatory interactions. It forms a general network which controls all sorts of responses which are needed in different conditions.



# Periodic genes in cell cycle network of yeast

Take regulatory data from Yeast Promoter Atlas (YPA). It contains 5026 genes including I22 TFs.

From this set of regulatory interactions, we extracted a cell-cycle specific subnetwork of 302 genes that were differentially expressed along the cell cycle of yeast (MA study by Spellman et al. Mol Biol Cell (1998)).

# MCDS of cell cycle network of yeast

Tightly interwoven network of 17 TFs and target genes that organize the cell cycle of S. *cerevisiae*.

Shown on the circumference of the outer circle are 164 target genes that are differentially expressed during the cell cycle and are regulated by a TF in the MCDS (shown in the inner circle).

The inner circle consists of the 14 TFs from the heuristic MCDS and of 123 other target genes that are regulated by at least two of these TFs



## **Studied networks: PluriNetwork**

PluriNetWork was manually assembled as an interaction/regulation network describing the molecular mechanisms underlying pluripotency.

It contains 574 molecular interactions, stimulations and inhibitions, based on a collection of research data from 177 publications until June 2010, involving 274 mouse genes/proteins.



Som A, et al. (2010) PLoS ONE 5: e15165.

## **MCDS of mouse pluripotency network**

Connectivity among TFs in the heuristic MCDS of the largest strongly connected component of a GRN for mouse ESCs.

The red circle borders mark the 7 TFs belonging to the set of master regulatory genes identified experimentally.

The MCDS genes were functionally significantly more homogeneous than randomly selected gene pairs of the whole network (p = 6.41e-05, Kolmogorov-Smirnow test).



# **Overlap with most central nodes**



Percentage overlap of the genes of the MDS and MCDS with the list of top genes (same size as MCDS) according to 3 centrality measures. Shown is the percentage of genes in the MDS or MCDS that also belong to the list of top genes with respect to degree, betweenness and closeness centrality

MDS nodes tend to be central in the network (high closeness) and belong to the most connected notes (highest degree).

When considering only outdegree nodes in the directed network, most of the top nodes of the MCDS have the highest overlap with the top nodes of the degree centrality and the betweenness centrality ( $\rightarrow$  connector nodes).

### **Breast cancer network**

Analyze breast cancer data from TCGA  $\rightarrow$ ca. 1300 differentially expressed genes.

Hierarchical clustering of coexpression network yielded 10 segregated network **modules** that contain between 26 and 295 gene members.

Add regulatory info from databases Jaspar, Tred, MSigDB.

(b) - (d) are 3 modules.



### **Breast cancer network**

The MDS and MCDS sets of the nine modules contain 68 and 70 genes, respectively.

Intersect the proteins encoded by these genes with the targets of anti-cancer drugs.

20 of the 70 proteins in the MCDS are known drug targets (p = 0.03, hypergeometric test against the network with 1169 genes including 228 drug target genes).

Also, 16 out of the 68 proteins belonging to the MDS genes are binding targets of at least one anti-breast cancer drug.

# $|MDS| \leq |MCDS|$



Number of MCDS genes determined by the heuristic approach or by the ILP formulation and in the MDS.

Shown are the results for 9 modules of the breast cancer network

# **Summary**

#### Today:

- network co-expression modules are best identified by ICA
- Network **motifs**: FFLs, SIMs, DORs are overrepresented  $\rightarrow$  different functions, different temporal behavior
- MDS and MCDS identify candidate master regulatory genes
  who reliable are they when applied to noisy and incomplete data?

#### Next lecture VI5:

• Epigenetics, analysis of DNA methylation data