# Biological network inference: a challenge to structured data-mining

Florence d'Alché-Buc<sup>1,2</sup>

<sup>1</sup>Visiting: INRIA, LRI umr CNRS 86 23, Université Paris-Sud, Orsay, France <sup>2</sup>Permanent address: IBISC EA 4526, Université d'Évry Val d'Essonne, Évry cedex, France



	$\sim$	$\sim$	0	<u>.</u>	10	١
(E	G	U.	2	υ	13	)

# Outline



- Introduction
- Gene regulatory networks
- Protein-protein interaction network

Gene regulatory network inference with operator-valued kernels

- Protein-protein interaction prediction with operator-valued kernels 3
- Conclusion

# Networks in molecular biology



• Circadian clock in mouse (Fig: Yan et al. Plos Comp. Biol. 2008)

# Gene regulatory networks

#### Circadian clock (CC)

- An example of a gene regulatory network involved in a cellular response to some input signal
- Involved in sustaining 24h-oscillations
- Feedback loop allows for control
- Cellular response to day-night alternance, meals etc...

## Gene regulatory networks

#### Definition of Transcriptional regulation

A gene a is said to regulate gene b if a codes for the protein A and A is a transcription factor of gene b. When A allows for the initiation of the transcription of b, gene a is said

to induce gene b. When A blocks the transcription, gene a is said to inhibit gene b.

## Gene regulatory networks

#### Definition of Transcriptional regulation

A gene a is said to regulate gene b if a codes for the protein A and A is a transcription factor of gene b. When A allows for the initiation of the transcription of b, gene a is said to induce gene b. When A blocks the transcription, gene a is said to

inhibit gene b.

#### A simple definition of a *gene regulatory network*

A gene (transcriptional) regulatory network is a dynamical system whose state variables are the mRNA's concentrations (possibly the proteins concentrations) and evolve through time.

・ 回 ト ・ ヨ ト ・ ヨ ト

## Reverse-engineering of gene regulatory networks

- Identify complex regulatory mechanisms at work in the cell
- Motivations: better understanding, predictive models for therapeutical targetting, biomarkers, personalized medecine, ...
- Main tasks related to data-mining
  - Parameter estimation in gene regulatory network models
  - Network inference

伺 ト イヨ ト イヨト

## **Biological network Inference**

#### Reverse-modeling of signalling and gene regulatory networks



# High throughput measurement techniques

#### Experimental techniques

- DNA chips
- Next Generation Sequencing methods
  - RNA Seq
  - Chip-seq

#### Data

- Gene expression level of tissue at a given time point after a long-term "run": steady state
- Time-course of gene expression: expression levels measured at a given time point for a given organism (one time-point: one organism)
- Perturbation data: Knock-out or knock-down of a gene and measurement of the gene expression level

# **Difficulties and limitations**

- Intrinsic noise, measurement on a cell population, measurement noise
- Limited size of data [ thousands of genes, tens of measurement]  $\neq$  BIG DATA
- Missing knowledge about the timing of regulation
- Nonlinearity of the behaviors
- Missing observations
- Other actors: role of chromatine, other kinds of regulations

A B A A B A

## Hopes

- Success of various clustering methods
- Reproducibility of data and behaviors
- Many other sources of knowledge/data: known functions of proteins, list of transcription factors, protein-protein interactions, metabolism (when relevant),...
- Multiple-view data
- Progress of experimental measurements and cost reduction

Gene regulatory network inference as a learning task

- Dimension reduction: clustering, biclustering
- Supervised link prediction : SIRENE (Mordelet et al. 2008)
- Model-free approaches to estimate the network structure: ARACNE (Margolin et al. 2006)
- Model-based and unsupervised approaches: Bayesian networks (Pe'er et al. 2001, Segal et al. 2003, ...), graphical Gaussian models (Strimmer et al. 2006)

A B b 4 B b

#### Network inference from time series

- Measurements of coupling: Kramer at al. 2009, mutual information : Zoppoli et al. 2010
- Differential equations: linear equations (Chen 1999), non-linear: S-systems (Voit et al. 2006)
- Autoregressive models

x<sub>t+1</sub> = h(x<sub>t</sub>) + e<sub>t</sub>, t > 0, x<sub>t</sub>: state vector, h ∈ H, e<sub>t</sub> : iid gaussian noise
 Dynamic Bayesian models

•  $x_{t+1}^i = h_i(Pa(i, t)) + \epsilon_t$ , t > 0, Pa(i,t): parent state variables at time t



Dynamic bayesian network without hidden variables

	0	$\sim$	0	n	1	3/	
( -	a	0	4	υ		J)	

## Network inference with autoregressive models

#### Sparse linear models

- Linear autoregressive models (Opgen-Rhein and Strimmer, 2007; Fujita et al. 2007, Shimamura et al. 2009)
- Granger causality (Shojaie and Michailidis, 2010 and 11)
- State-space models (Perrin et al. 2003, Rangel et al. 2004)
- Several order autoregressive models (Lozano 2009, Bolstad et al. 2011)
- Time-varying models (Lebre et al. 2011)

A B F A B F

# Network inference with nonlinear autoregressive models

#### Nonlinear nonparametric models

- Dynamic Bayesian Networks (Imoto et al. 2002, Husmeier et al. 2005, Li et al. 2007, Bansal et al. 2007)
- Gaussian processes for network inference (Aijo and Lahdesmaki 2009)

向下 イヨト イヨト

# Outline

#### Introduction

- Gene regulatory networks
- Protein-protein interaction network

2 Gene regulatory network inference with operator-valued kernels

#### Protein-protein interaction prediction with operator-valued kernels

#### 4 Conclusion

## Protein-protein interaction network



Protein-protein interaction network in yeast

(EGC 2013)

Network inference

# Experimental detection of protein-protein interactions

- in vivo large scale systems:
  - Y2H high false positive rate
- in vitro small scale methods: costly and laborious
  - protein-arrays
  - co-immunoprecipitations
  - FRET, NMR

( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( ) < ( )

#### Data

Combine indirect information on pairs of proteins with direct information on ppi data

- over-represented domains or motifs pairs
- structural information, primary sequences
- subcellular localization
- biological functions
- co-expression of genes
- conservation of pairs of sequences

## Difficulties

- Generally no structural information
- Very few labeled edges and no negative labels
- Relevant features ? Context ?
- Source of knowledge: some information used as input feature have been inferred from the outputs

#### Hopes

- Global improvement of datasets and databases
- Better encoding of structured knowledge
- Transfer learning, multi-task learning

# Existing approaches for link prediction (1): supervised classification

- Pairwise SVM [Ben-Hur and Noble 2005]
- Random forest, mixture of feature experts (Qi 2008), ensemble methods with original evolutionary features (De Vienne and AzŐ, 2012)
- Supervised Learning of a kernel or a similarity
  - With KCCA [Yamanishi et al. 2004], with metric learning [Yamanishi and Vert 2005]
  - With output kernel regression tree [Geurts et al. 2006,07], with output kernel gradient boosting [Geurts et al. 2007]
- Supervised classification linked to a node
  - Iocal classifiers [Bleakley et al. 2007]

< 回 > < 回 > < 回 >

Existing approaches (2): semi-supervised and transductive learning

- Kernel Matrix completion
  - Using EM [Tsuda et al. 2003] and [Kato et al. 2005]
- Transductive or semi-supervised learning
  - Link Propagation [Kashima et al. 2009]
  - Mixture of Wishart Matrices [Dit-Yeung 2009]
  - Training set expansion [Yip and Gerstein 2009]

# Outline

#### Introduction

- Gene regulatory networks
- Protein-protein interaction network

#### 2 Gene regulatory network inference with operator-valued kernels

#### Protein-protein interaction prediction with operator-valued kernels

#### 4 Conclusion

## Extend linear autoregressive models

#### Network inference with linear models

- Estimate B in  $\mathbf{x}_{t+1} = B\mathbf{x}_t + \epsilon_t$  with a **sparsity** constraint
- Threshold B to get an estimation of the true adjacency matrix A:

#### Operator-valued kernels based model

- Kernel-based models for vector prediction use operator-valued kernels (Micchelli and Pontil 2005, Caponnetto et al. 2008)
- Representer theorem for semi-supervised learning (\*\*) give:

•  $h(\mathbf{x}) = \sum_{\ell} K(\mathbf{x}_t, \mathbf{x}_\ell) \mathbf{c}_\ell$ 

## Extend linear autoregressive models

Network inference with operator-valued kernels\*

- $\mathbf{x}_{t+1} = h(\mathbf{x}_t) + \epsilon_t$
- $h(\mathbf{x}_t) = \sum_{i=1}^{N-1} K(\mathbf{x}_t, \mathbf{x}_\ell) \mathbf{c}_\ell$
- Use the following estimate:

$$\hat{A}_{ij} = sgn\left(\frac{1}{N+1}\sum_{t=1}^{N+1}\frac{\partial h(x_t)_i}{\partial (x_t)_j} - \theta\right)$$

• \*: collaboration with Nehemy Lim and George Michailidis

## Which matrix-valued kernel?

- Let us take  $k_{\gamma_1}$  as the (scalar) Gaussian kernel:
- $\forall (x, y) \in \mathbb{R} \times \mathbb{R}, k_{\gamma_1}(\mathbf{x}, \mathbf{y}) = exp(-\gamma_1 ||\mathbf{x} \mathbf{y}||^2).$
- and the matrix-valued kernel:  $K_{\gamma_2}(\mathbf{x}, \mathbf{y})_{ij} = exp(-\gamma_2(x^i y^j)^2)$
- Finally,  $K(\mathbf{x}, \mathbf{y}) = k_{\gamma_1}(\mathbf{x}, \mathbf{y}) B \circ K_{\gamma_2}(\mathbf{x}, \mathbf{y})$
- Important: Sparsity of B controls sparsity of the Jacobian

# Learning h

#### Learning algorithm

- B is estimated as well as c's
- Boosting algorithm (build  $H(\mathbf{x}_t) = \sum_m h_m(\mathbf{x}_t)$ )
  - Base model: a model h defined on a random subspace
  - Construction through  $\ell_2$ -boosting

A B F A B F

# Inference of a synthetic network in yeast [IRMA, Cantone, 2009]

	Switch-off		Switch-on		
	AUROC	AUPR	AUROC	AUPR	
OKVAR-Boost	0.807	0.807	1	1	
LASSO	0.500	0.253	0.583	0.474	
Äijö	0.875	0.848	0.838	0.836	

イロト イヨト イヨト

# Results on 10-size networks (DREAM3 challenge)

		Ecoli1		Eco		
	Size-10	AUROC	AUPR	AUROC	AUPR	AURC
	OKVAR + True B	0.932	0.712	0.814	0.754	0.85
	OKVAR-Boost (1 TS)	0.665	0.272	0.629	0.466	0.66
		$\pm$ 0.088	$\pm$ 0.081	$\pm$ 0.095	$\pm$ 0.065	$\pm$ 0.03
	OKVAR-Boost (4 TS)	0.853	0.583	0.749	0.536	0.68
	LASSO	0.500	0.119	0.547	0.531	0.528
	Team 236	0.621	0.197	0.650	0.378	0.64
	Team 190	0.573	0.152	0.515	0.181	0.63

Table: AUROC and AUPR for OKVAR-Boost ( $\lambda_1 = 1, \lambda_2 = 10$  selected by *Block-Stability*), LASSO, Team 236 and Team 190 (DREAM3 challenge) run on DREAM3 size-10 networks. OKVAR-Boost results using respectively one time series (OKVAR-Boost (1 TS)) (Average  $\pm$  Standard Deviations) and the four available time series (OKVAR-Boost (4 TS)) are from consensus networks. The numbers in **boldface** are the maximum values of each column. (\* Consensus thresholds for Yeast2 and Yeast3 are different due to 29/50

# Outline

#### Introduction

- Gene regulatory networks
- Protein-protein interaction network

2 Gene regulatory network inference with operator-valued kernels

#### Protein-protein interaction prediction with operator-valued kernels

### 4 Conclusion

## Protein-protein interaction network inference



- $\mathcal{V}$  : set of nodes (corresponding to proteins)
- An edge between nodes v and v' means a physical interaction between proteins v and v'

(EGC 2013)

Network inference

## Supervised link prediction

- Edges are known for the  $\ell$  first nodes ( $\mathcal{V}_{\ell}$ )
- **Goal**: learning a predictor  $f : \mathcal{V} \times \mathcal{V} \rightarrow \{0, 1\}$  from:
  - descriptions of proteins in  $\mathcal{V}_{\ell}$  (localization, sequence ...),
  - the adjacency matrix  $A_{\ell}$  of the training subgraph.



< B.

# Semi-supervised link prediction

- Let us use unlabeled data !
- Let  $\mathcal{V}_{\ell+u} = \{v_1, ..., v_{\ell+u}\} : \ell$  fully labeled nodes, u unlabeled nodes
- We assume that the description of  $v_{\ell+1}, ..., v_{\ell+u}$  is known
- Goal: learning a predictive model *f* : V × V → {0, 1} from descriptions of proteins in V<sub>ℓ+u</sub> and A<sub>ℓ</sub>, with ℓ << u.</li>



# Building a classifier f by learning similarity $\kappa_y$

•  $\kappa_y$ : similarity between two proteins as nodes in the known graph,

Similarity-based model:

$$f_{\theta}(\boldsymbol{v},\boldsymbol{v}') = sgn(\hat{\kappa_y}(\boldsymbol{v},\boldsymbol{v}') - \theta)$$

• Learning a proxy of  $\kappa_y$  and choosing  $\theta$  = learning the classifier  $f_{\theta}$ 

## Reminder: scalar-valued kernel

#### Definition

A symmetric function *k* from  $\mathcal{V} \times \mathcal{V}$  to *IR* is said to be a definite positive kernel if and only if: For any positive integer n, for any set of n objects  $(v_1, ..., v_n) \in \mathcal{V}$ , for any real  $c_1, ..., c_n$ ,

$$\sum_{i,j} c_i c_j k(v_i, v_j) \ge 0 \tag{1}$$

#### Theorem

For any positive definite kernel k on  $\mathcal{V} \times \mathcal{V}$ , there exists an Hilbert space  $\mathcal{F}$  with and a feature mapping  $\phi : \mathcal{V} \to \mathcal{F}$  such that for all (v, v'):  $k(v, v') = \langle \phi(v), \phi(v') \rangle_{\mathcal{F}}$ 

通 ト イ ヨ ト イ ヨ ト

## Now, let us introduce scalar-valued kernels (1)

#### Assumptions about the outputs

- Let us assume we only know for training data, the value of the  $\ell \times \ell$  Gram matrix  $K_y$  of an output kernel:  $(K_y)_{ij} = k_y(v_i, v_j)$ ,  $k_y : \mathcal{V} \times \mathcal{V} \to \mathbb{R}$
- For instance,  $K_y$  is a diffusion kernel matrix

# Output Gram matrix: diffusion kernel

Only 
$$\kappa_y(v_i, v_j) = \langle y(v_i), y(v_j) \rangle_{\mathcal{F}_y}$$
 for  $i, j = 1, \dots, \ell$  are known.

• Here we use the **diffusion kernel** [Kondor & Lafferty, 2002)] : The Gram matrix  $K_{Y_{\ell}}$  with  $K_{i,j} = \kappa_Y(v_i, v_j)$  is given by:

$$K_{Y_{\ell}} = \exp(-\beta L),$$

where the graph Laplacian *L* is defined by :

$$L=D_\ell-A_\ell,$$

with  $A_{\ell}$  the adjacency matrix and  $D_{\ell}$  the diagonal matrix of vertices degrees.

# Building the classifier f by learning an output kernel $k_y$

$$\forall (\mathbf{v}, \mathbf{v}') \in \mathcal{V} imes \mathcal{V}, k_{\mathbf{y}}(\mathbf{v}, \mathbf{v}') = < \mathbf{y}(\mathbf{v}), \mathbf{y}(\mathbf{v}') >_{\mathcal{F}_{\mathbf{y}}}$$

- Let us learn to predict *y* with a function  $h: \mathcal{V} \to \mathcal{F}_y$
- Then we will get:  $\hat{k_y}(v, v') = \langle h(v), h(v') \rangle_{\mathcal{F}_y}$

=> instead of learning a pairwise classifier, we learn a single variable function with output values in a Hilbert space

#### How to learn h?



э

イロト イポト イヨト イヨト

# Which family H of models to build ?

#### Output kernel tree (OK3)\* and extensions

• 
$$h_{tree}(v) = \sum_{m=1}^{M} \mathbf{1}_m(x(v)).\bar{\mathbf{y}}_m$$

• where M is the number of leaves in the tree,  $1_m(x(v)) = 1$  if x(v) falls into leaf m and 0 otherwise

• 
$$\bar{\mathbf{y}}_m = \frac{1}{N_m} \sum_{i=1}^n \mathbf{1}_m(x(v_i))$$

- \*: joint work with Pierre Geurts (Geurts et al. 2006, Geurts et al. 2007) and Louis Wehenkel (Geurts et al. 2007)
- Extension to boosting and random forests

# Which family H of models to build for semi-supervised learning ?

#### Operator-valued kernels based models\*\*

- Kernel-based models for vector prediction use operator-valued kernels (Micchelli and Pontil 2005, Caponnetto et al. 2008)
- Representer theorem for semi-supervised learning (\*\*) give:
   h(v) = ∑<sup>ℓ+u</sup><sub>i=1</sub> K<sub>x</sub>(v, v<sub>i</sub>)c<sub>i</sub>
- \*\*: joint work with Céline Brouard (PhD student) and Marie Szafranski (Brouard et al. 2011)

## Operator-valued kernel-based model 1

#### Choice of $K_x$

- Let us define an operator valued kernel  $K_x : \mathcal{V} \to \mathcal{L}(\mathcal{F}_y)$
- *L*(*F<sub>y</sub>*) is the set of bounded operators on *F<sub>y</sub>*
- A simple but efficient choice of  $K_x$  is  $K_x(v, v') = k_x(v, v')$ . Id
- With  $k_x(v, v') = \langle x(v), x(v') \rangle_{\mathcal{F}_x}, x : \mathcal{V} \to \mathcal{F}_x$

A B b A B b

#### Operator-valued kernel-based model 2

- $J(h) = \sum_{i=1}^{\ell} \|h(v_i) y_i\|_{\mathcal{F}_y}^2 + \lambda_1 \|h\|_{\mathcal{H}}^2 + \lambda_2 \sum_{i,j=1}^{\ell+u} k_{x,ij} \|h(v_i) h(v_j)\|_{\mathcal{F}_y}^2$
- Minimizing J(h) gives a closed-form solution for h

  as for k
  <sub>y</sub>

# Network inference on yeast ppinetwork, improvement brought by unlabeled data



- Results of Céline Brouard
- 10 random training/test sets

#### Results of network inference: AUROC on yeast ppi net



#### Results of Céline Brouard

(EGC 2013)

45 / 50

#### Results of network inference: AUPR on yeast ppi net



#### Results of Céline Brouard

(EGC 2013)

#### Network inference

46 / 50

# Outline

#### Introduction

- Gene regulatory networks
- Protein-protein interaction network

2 Gene regulatory network inference with operator-valued kernels

#### Protein-protein interaction prediction with operator-valued kernels

## 4 Conclusion

# Input and Output Kernel Regression (IOKR)

#### Take home message

Take car of the output space (choice of  $k_y$ ) as well as the input space (choice of  $K_x$ )

#### Issues

- Kernel design, kernel learning
- Multiple kernel learning / data integration
- Model selection, scaling
- Pre-image problem in general

## Remaining challenges in network inference

- Integration of various structured prior knowledge
- Scaling to a large number of genes
- Experimental design: active learning

## Acknowledgements

- IBISC CNRS Université d'Evry, Genopole, France
  - Céline Brouard (PhD student)
  - Marie Szafranski
- INRIA-Saclay LRI, CNRS, Orsay, France
  - Jérome Azé
  - Christine Froidevaux
- Université de Liège, GIGA, Belgium
  - Pierre Geurts
  - Louis Wehenkel
- Radbound University Nijmegen, The Netherlands
  - Adriana Birlutiu
  - Tom Heskes
- Hopital Necker, INSERM, Paris, France
  - Aleksander Edelman
  - Chiara Guerrera