

Inferring and using biological networks

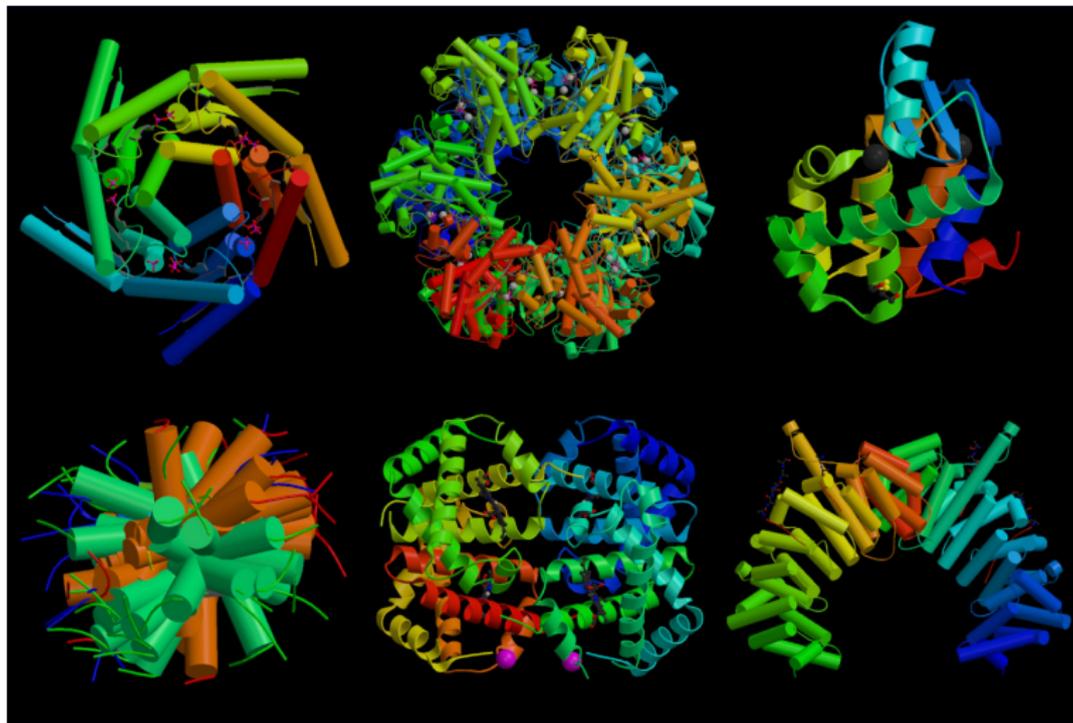
Jean-Philippe Vert

Jean-Philippe.Vert@mines-paristech.fr

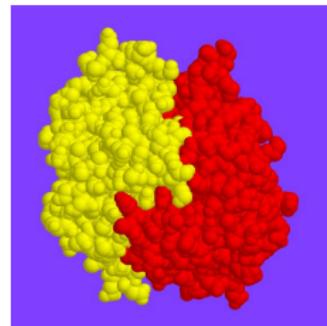
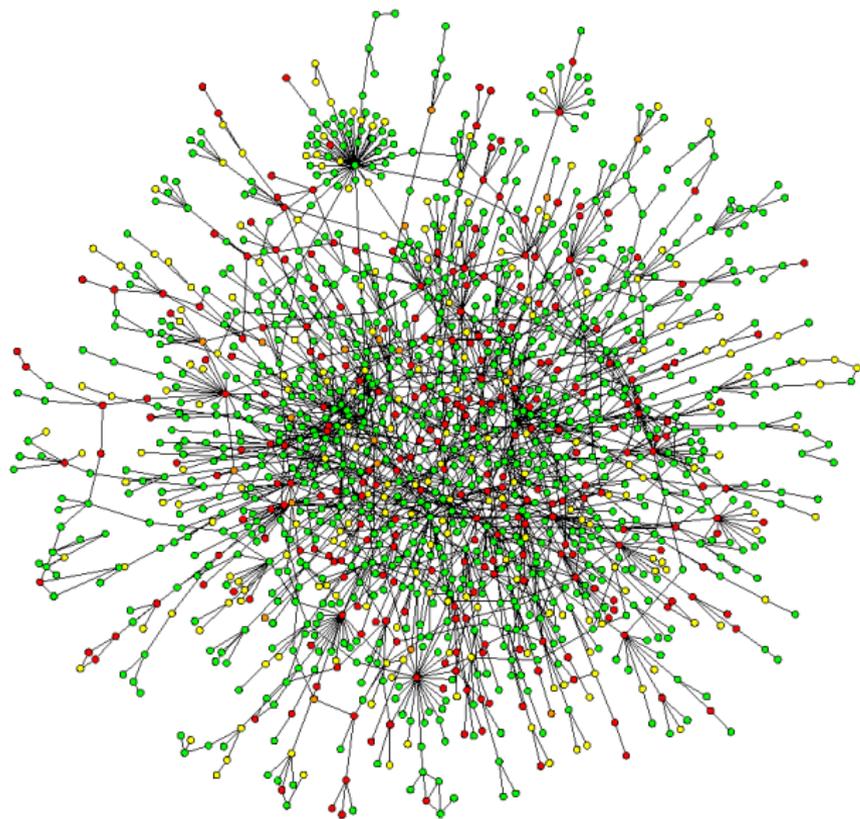
Mines ParisTech / Institut Curie / INSERM U900

University of Laval, Quebec, Canada, December 4, 2008.

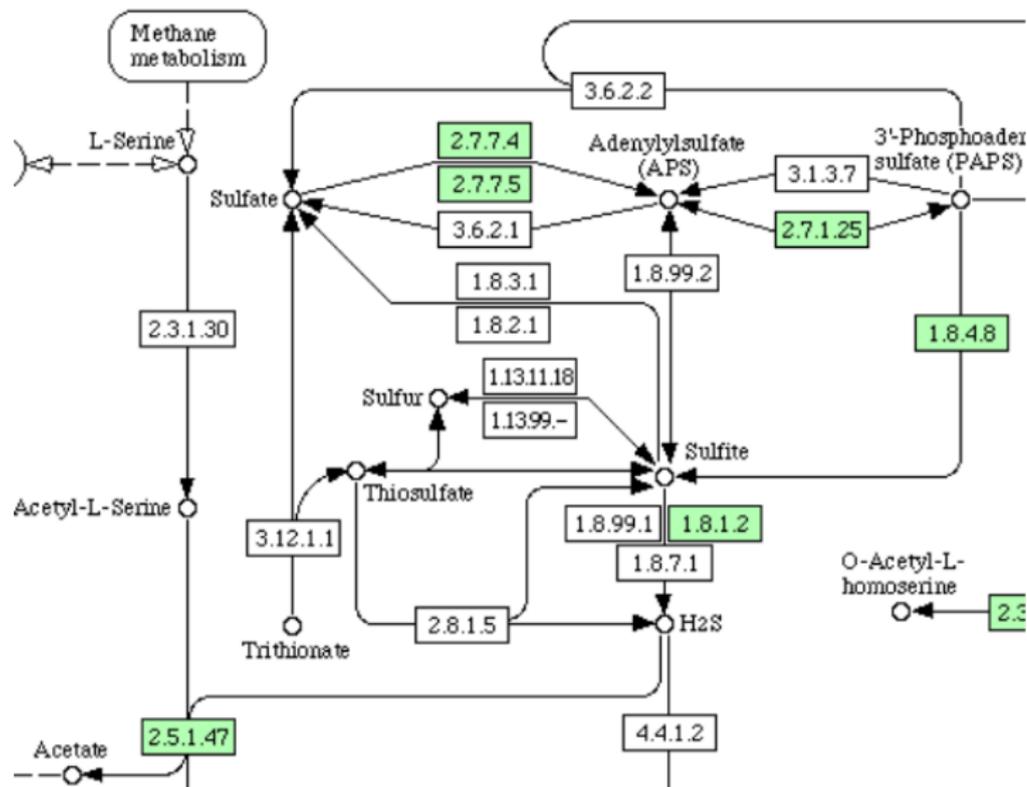
We have many genes and proteins..



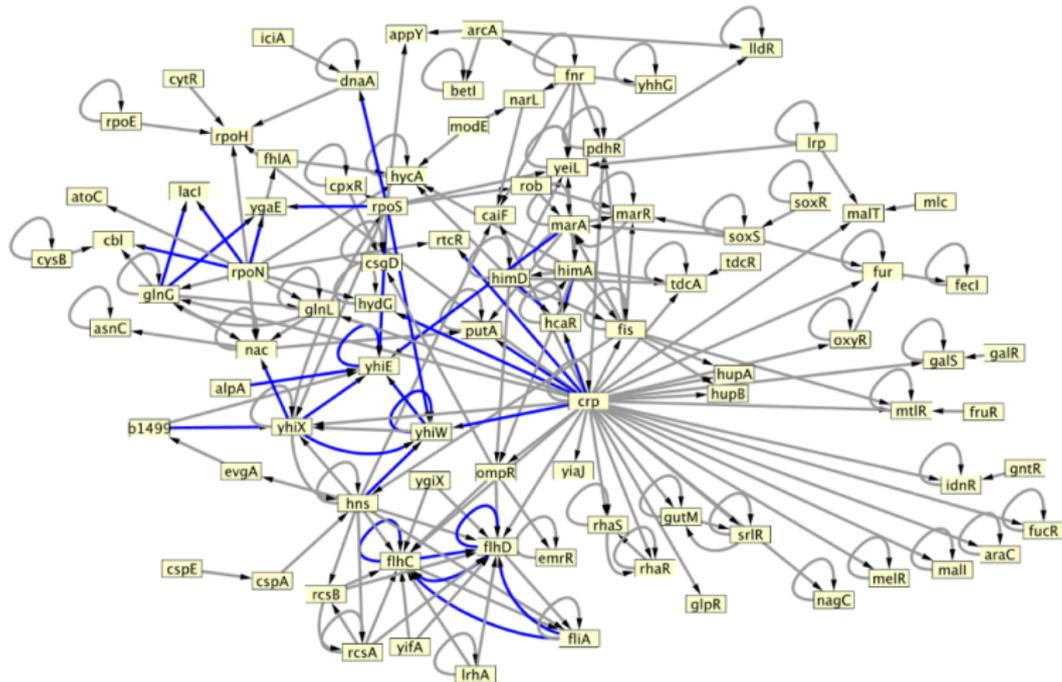
Network 1: protein-protein interaction



Network 2: metabolic network

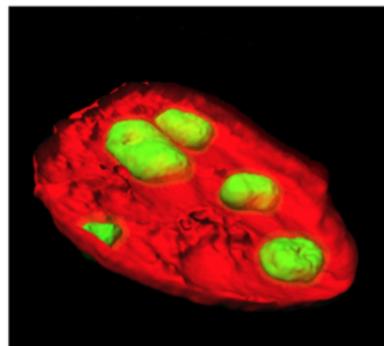
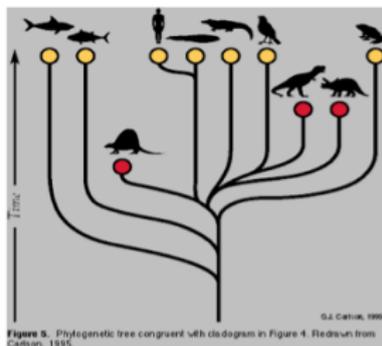
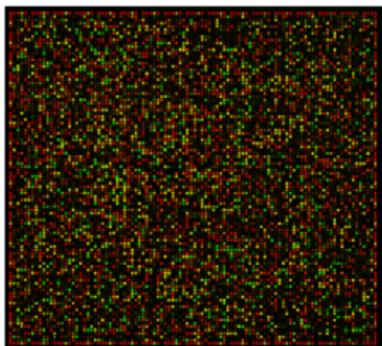


Network 3: gene regulatory network

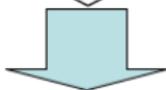
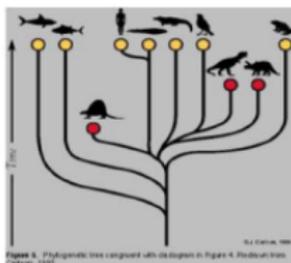
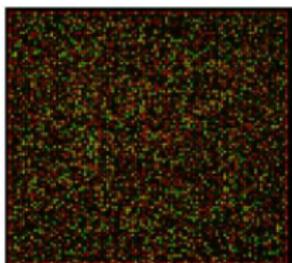


Biologists have collected a lot of data about proteins. e.g.,

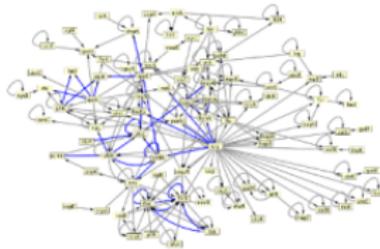
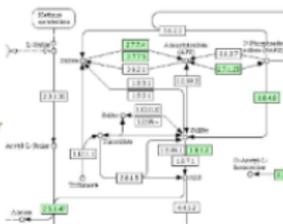
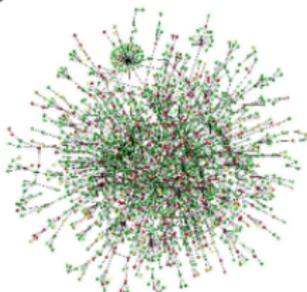
- Gene expression measurements
- Phylogenetic profiles
- Location of proteins/enzymes in the cell



Problem 1 : how to infer relationships between genes from biological data?



Inference



Problem 2 : how to use biological networks to help in the analysis of genomic data?

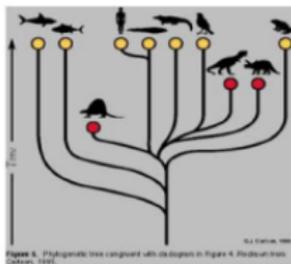
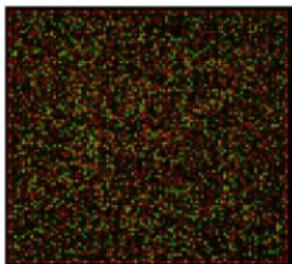
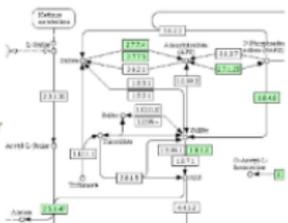
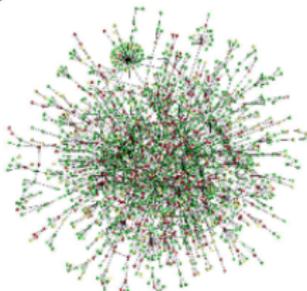


Figure 4. Phylogenetic tree compared with its diagram in Figure 4. *Protein Data Bank, 1999.*



Interpretation



- 1 How to infer relationships between genes from biological data?
- 2 How to use biological networks to help in the analysis of genomic data?
- 3 Conclusion

- 1 How to infer relationships between genes from biological data?
- 2 How to use biological networks to help in the analysis of genomic data?
- 3 Conclusion

“De novo” inference

- Given data about individual genes and proteins, ...
- ... Infer the edges between genes and proteins

Typical strategies

- Fit a **dynamical system** to time series (e.g., PDE, boolean networks, state-space models)
- Detect **statistical conditional independence or dependency** (Bayesian network, mutual information networks, co-expression)

“De novo” inference

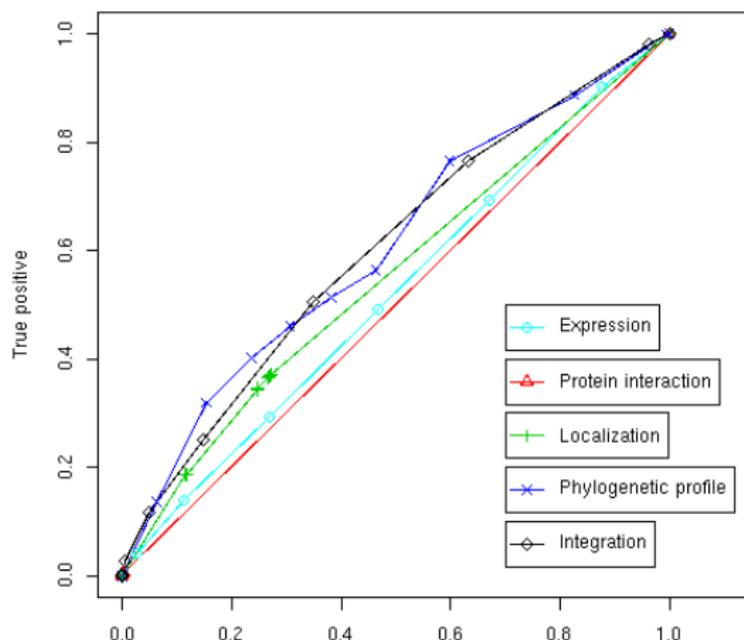
- Given data about individual genes and proteins, ...
- ... Infer the edges between genes and proteins

Typical strategies

- Fit a **dynamical system** to time series (e.g., PDE, boolean networks, state-space models)
- Detect **statistical conditional independence or dependency** (Bayesian network, mutual information networks, co-expression)

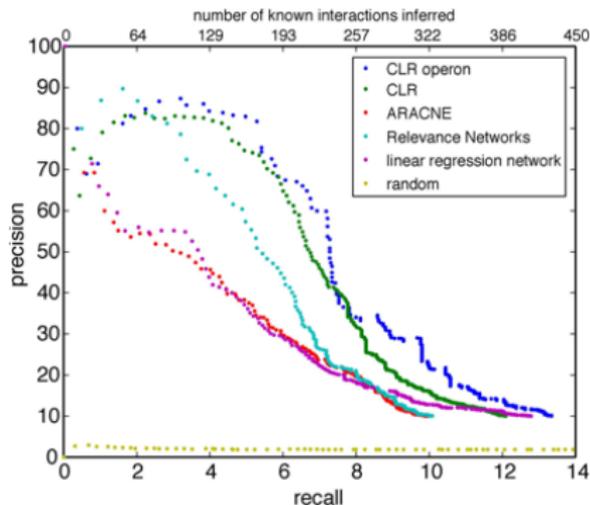
Evaluation on metabolic network reconstruction

- The known metabolic network of the yeast involves **769 proteins**.
- Predict edges from distances between a variety of genomic data (expression, localization, phylogenetic profiles, interactions).



Large-Scale Mapping and Validation of *Escherichia coli* Transcriptional Regulation from a Compendium of Expression Profiles

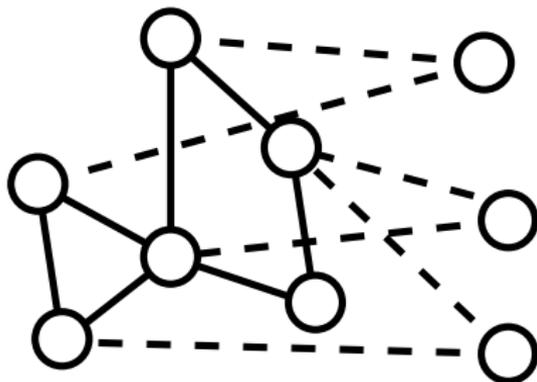
Jeremiah J. Faith¹, Boris Hayete¹, Joshua T. Thaden^{2,3}, Ilaria Mogno^{2,4}, Jamey Wierzbowski^{2,5}, Guillaume Cottarel^{2,5}, Simon Kasif^{1,2}, James J. Collins^{1,2}, Timothy S. Gardner^{1,2*}



Motivation

In actual applications,

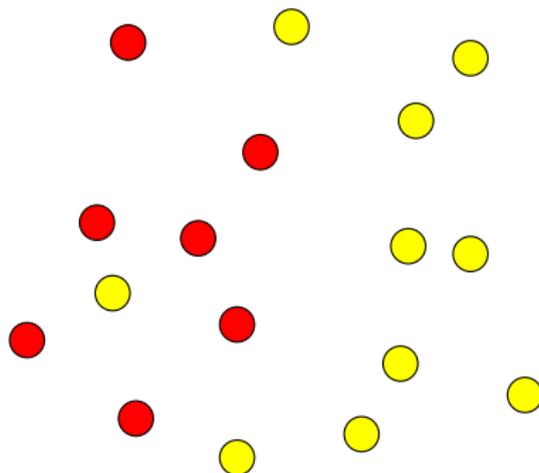
- we know in advance parts of the network to be inferred
- the problem is to add/remove nodes and edges using genomic data as side information



Supervised method

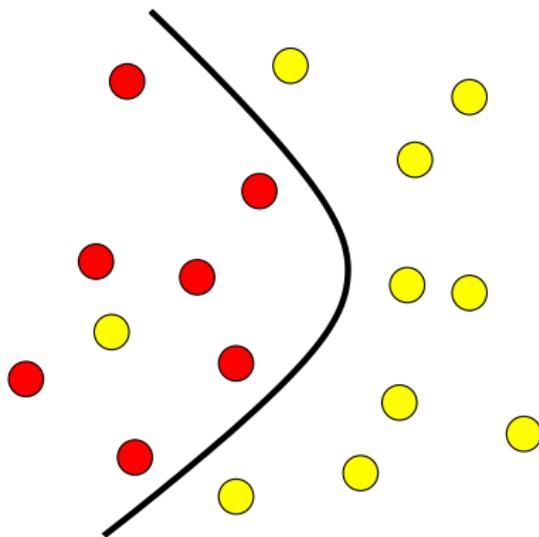
- Given genomic data **and** the currently known network...
- Infer **missing edges** between current nodes and additional nodes.

Interlude : Pattern recognition



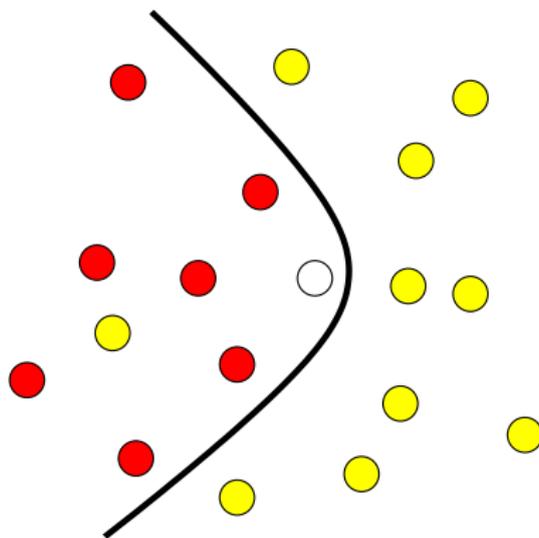
- Given a training set of patterns in two classes, learn to discriminate them
- Many algorithms (ANN, SVM, Decision trees, ...)

Interlude : Pattern recognition



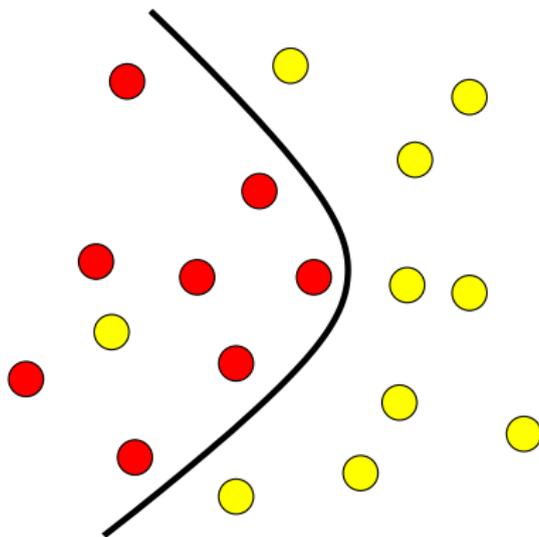
- Given a training set of patterns in two classes, learn to discriminate them
- Many algorithms (ANN, SVM, Decision trees, ...)

Interlude : Pattern recognition



- Given a training set of patterns in two classes, learn to discriminate them
- Many algorithms (ANN, SVM, Decision trees, ...)

Interlude : Pattern recognition



- Given a training set of patterns in two classes, learn to discriminate them
- Many algorithms (ANN, SVM, Decision trees, ...)

Pattern recognition and graph inference

Pattern recognition

Associate a binary label Y to each data X

Graph inference

Associate a binary label Y to each **pair** of data (X_1, X_2)

Two solutions

- Consider each pair (X_1, X_2) as a single data -> **learning over pairs**
- Reformulate the graph inference problem as a pattern recognition problem at the level of individual vertices -> **local models**

Pattern recognition and graph inference

Pattern recognition

Associate a binary label Y to each data X

Graph inference

Associate a binary label Y to each **pair** of data (X_1, X_2)

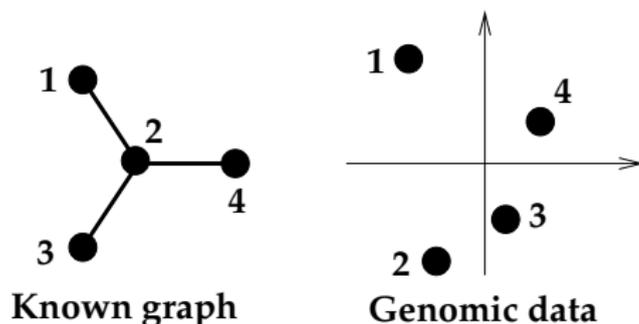
Two solutions

- Consider each pair (X_1, X_2) as a single data -> **learning over pairs**
- Reformulate the graph inference problem as a pattern recognition problem at the level of individual vertices -> **local models**

Pattern recognition for pairs

Formulation and basic issue

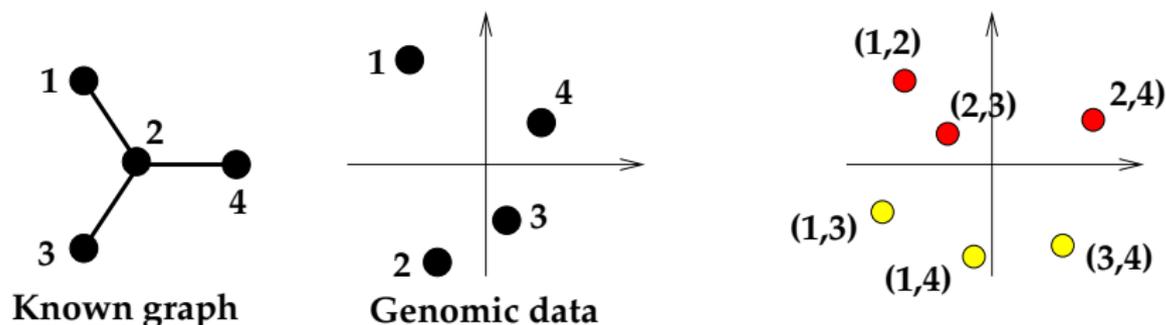
- A pair can be **connected (1)** or **not connected (-1)**
- From the known subgraph we can **extract examples** of connected and non-connected pairs
- However the genomic data characterize **individual** proteins; we need to work with **pairs** of proteins instead!



Pattern recognition for pairs

Formulation and basic issue

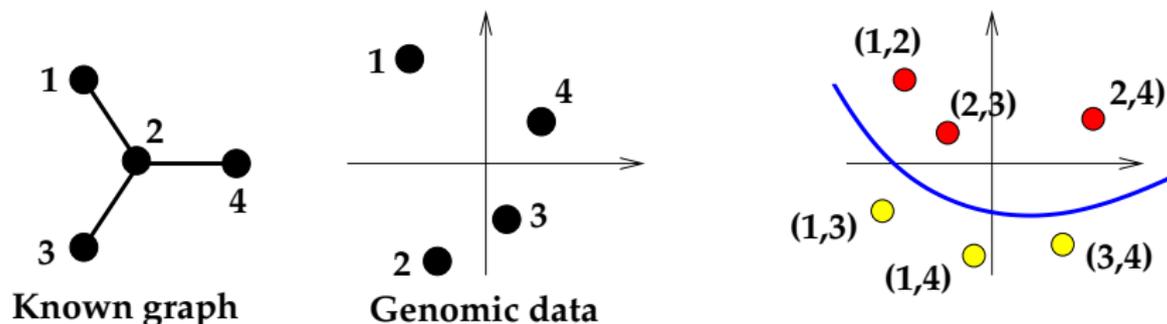
- A pair can be **connected** (1) or **not connected** (-1)
- From the known subgraph we can **extract examples** of connected and non-connected pairs
- However the genomic data characterize **individual** proteins; we need to work with **pairs** of proteins instead!



Pattern recognition for pairs

Formulation and basic issue

- A pair can be **connected** (1) or **not connected** (-1)
- From the known subgraph we can **extract examples** of connected and non-connected pairs
- However the genomic data characterize **individual** proteins; we need to work with **pairs** of proteins instead!



Concatenation

- A simple idea is to **concatenate** the vectors u and v describing two proteins to obtain a description of the pair:

$$\psi(u, v) = \begin{pmatrix} u \\ v \end{pmatrix}.$$

Symmetric tensor product (Ben-Hur and Noble, 2006)

$$K_{pair} [(A, B), (C, D)] = k(A, C)k(B, D) + k(A, D)k(B, C)$$

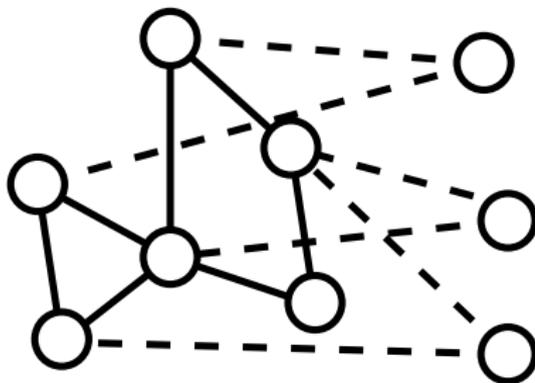
Intuition: a pair (A, B) is similar to a pair (C, D) if:

- A is similar to C **and** B is similar to D , **or**...
- A is similar to D **and** B is similar to C

Supervised inference with local models

The idea (Bleakley et al., 2007)

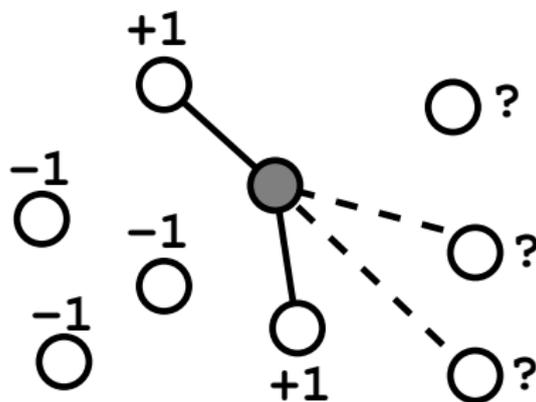
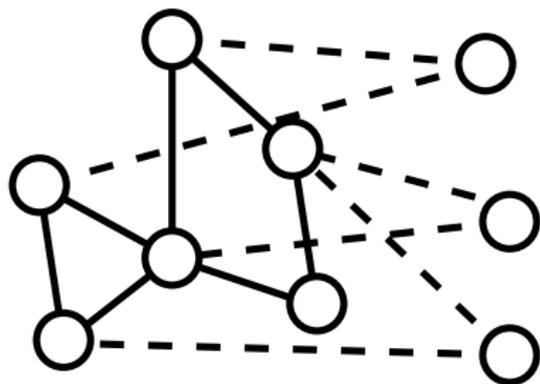
- Motivation: define **specific models** for **each target node** to discriminate between its neighbors and the others
- Treat each node independently from the other. Then **combine** predictions for ranking candidate edges.



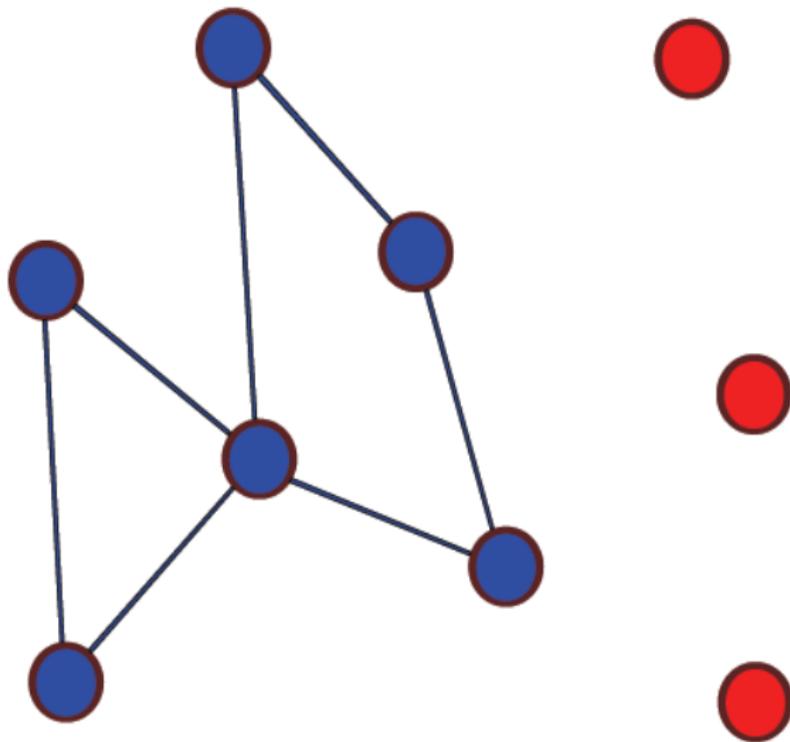
Supervised inference with local models

The idea (Bleakley et al., 2007)

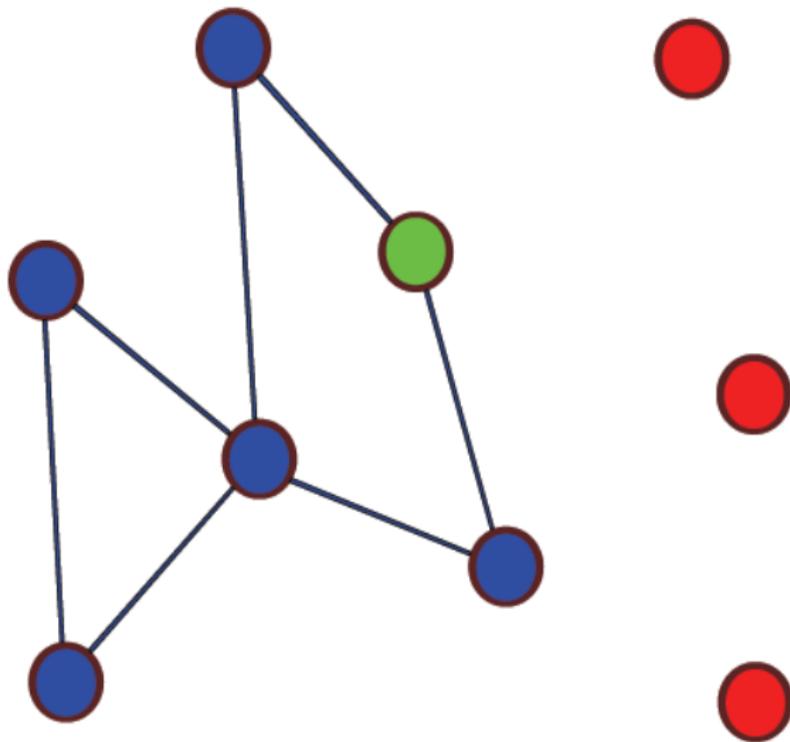
- Motivation: define **specific models** for **each target node** to discriminate between its neighbors and the others
- Treat each node independently from the other. Then **combine** predictions for ranking candidate edges.



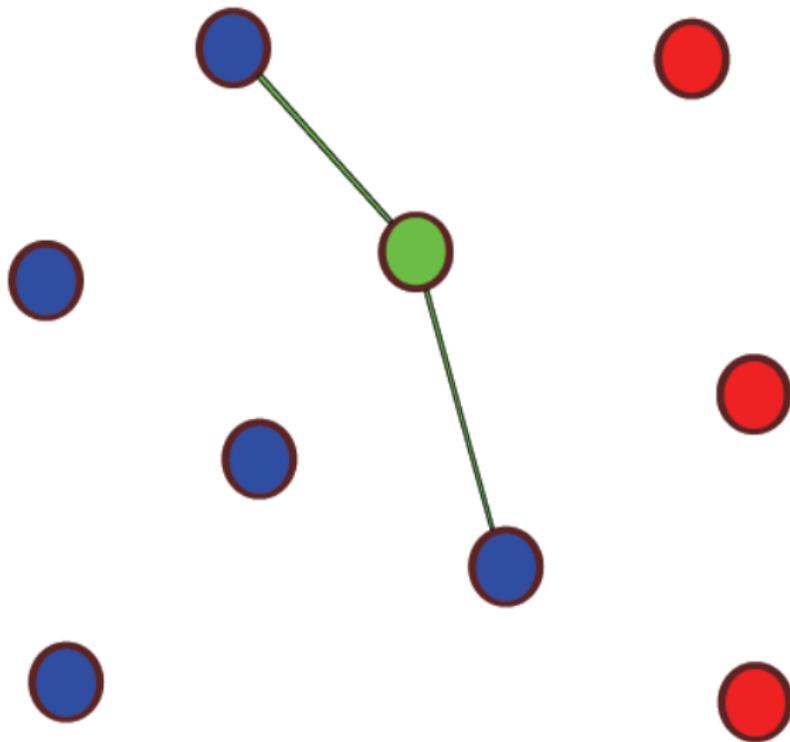
The LOCAL model



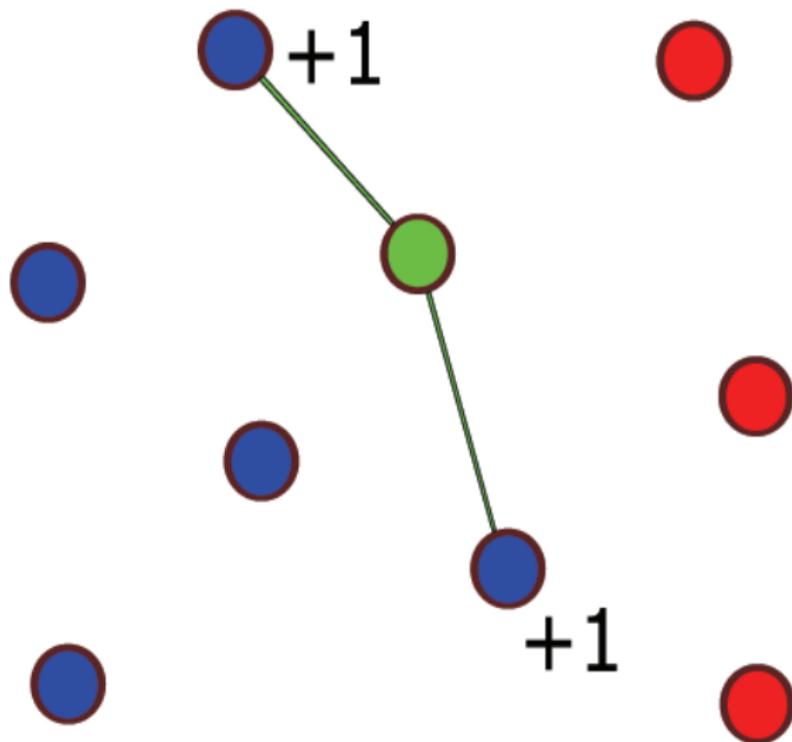
The LOCAL model



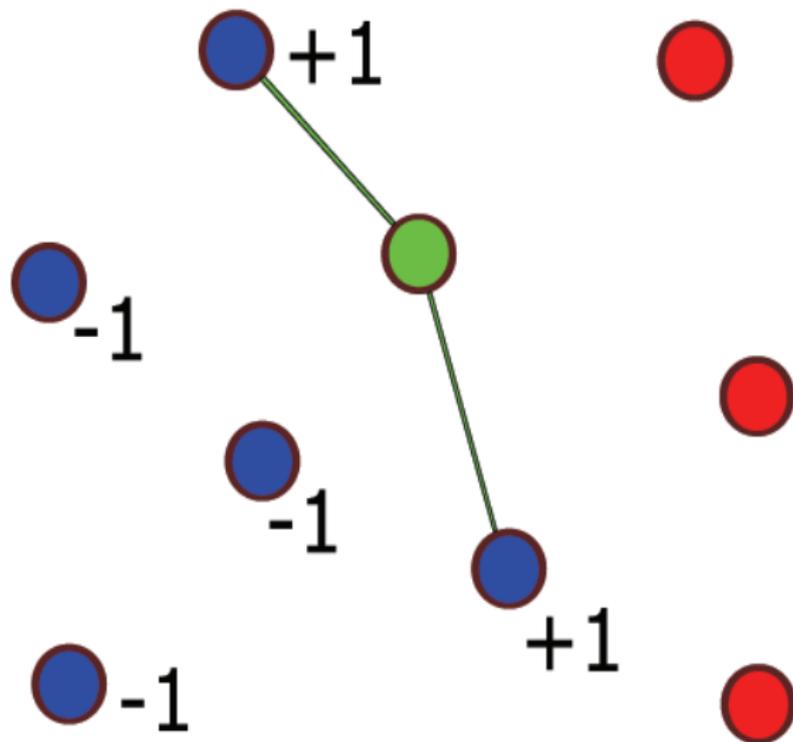
The LOCAL model



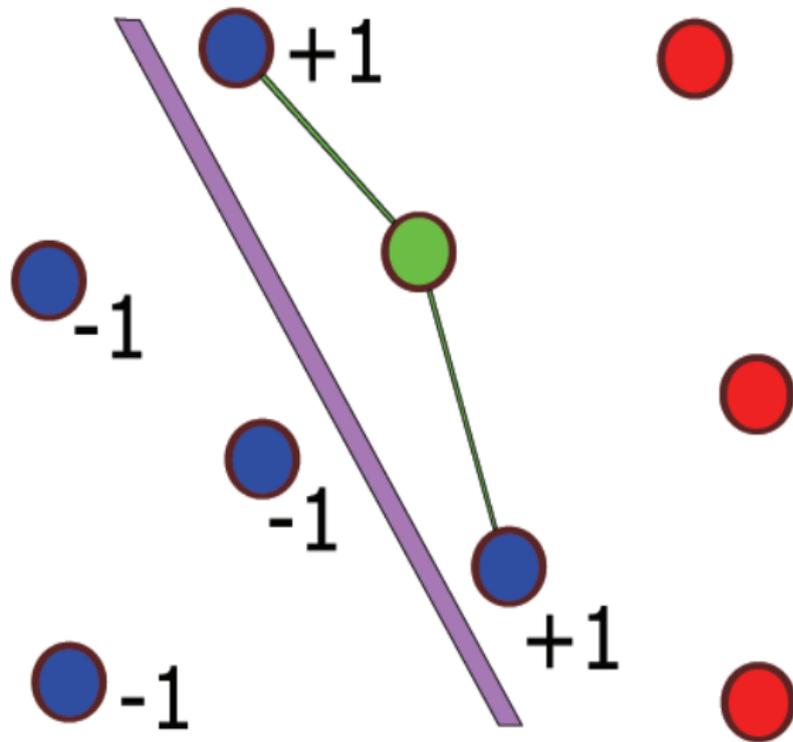
The LOCAL model



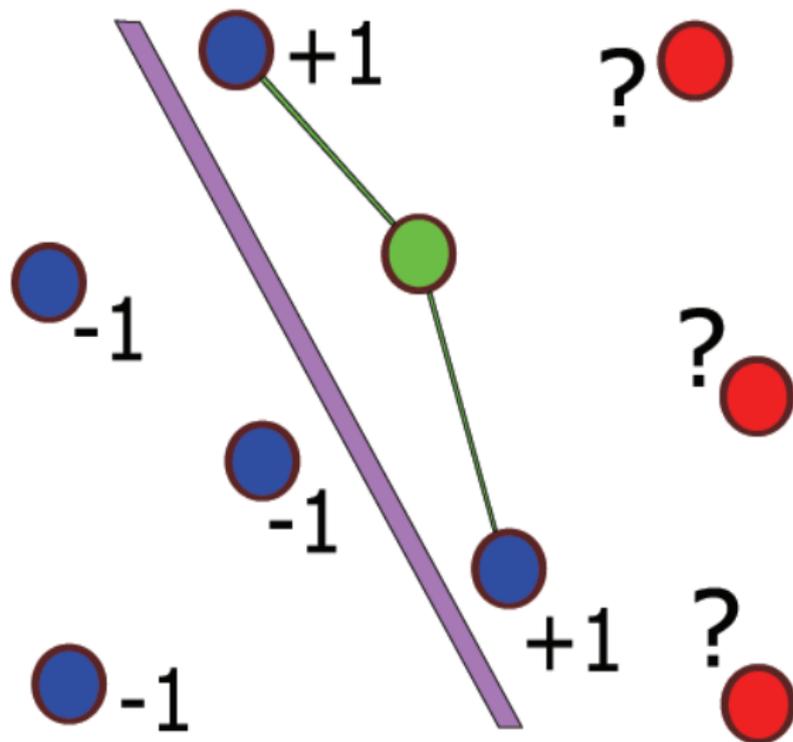
The LOCAL model



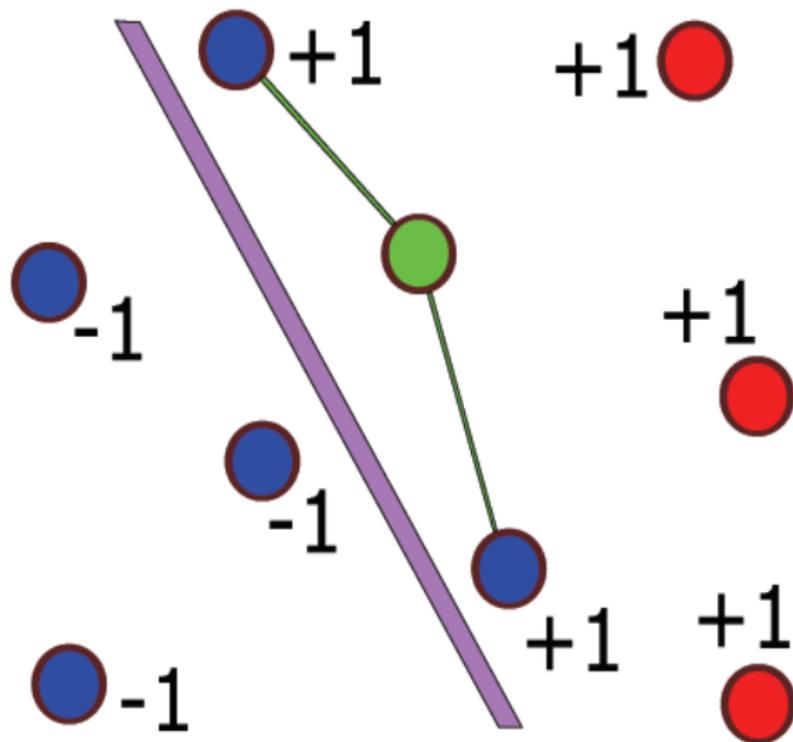
The LOCAL model



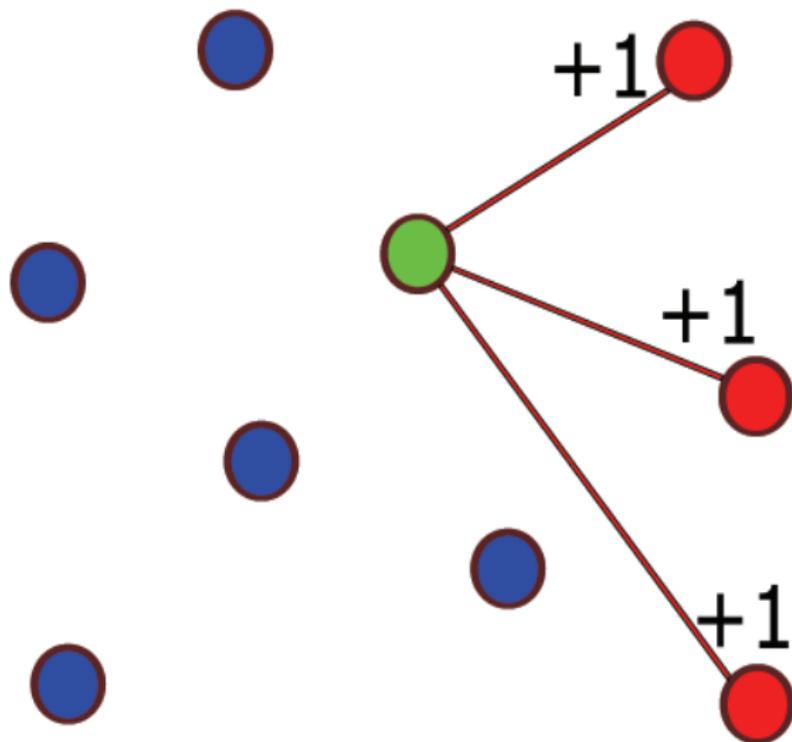
The LOCAL model



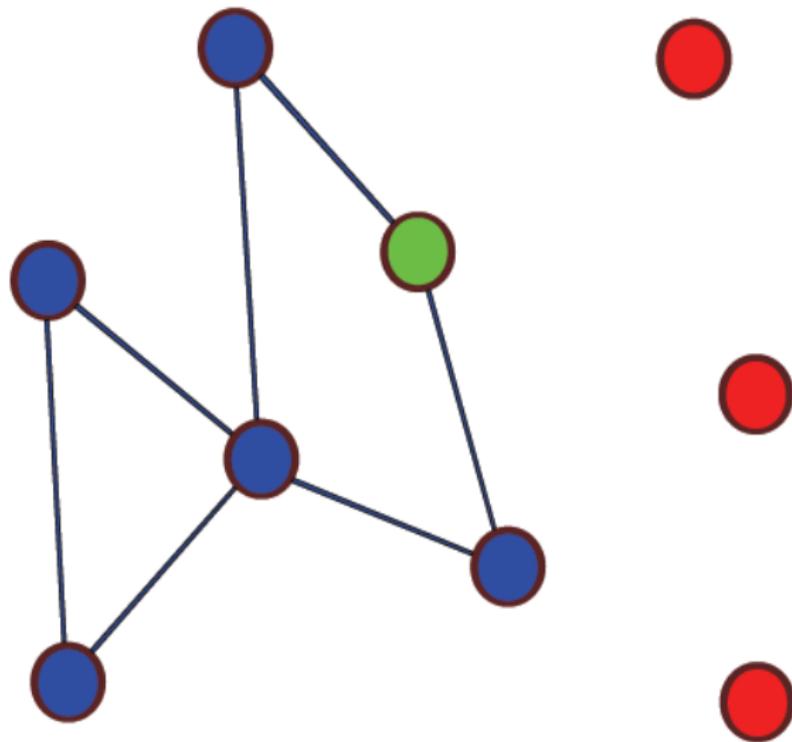
The LOCAL model



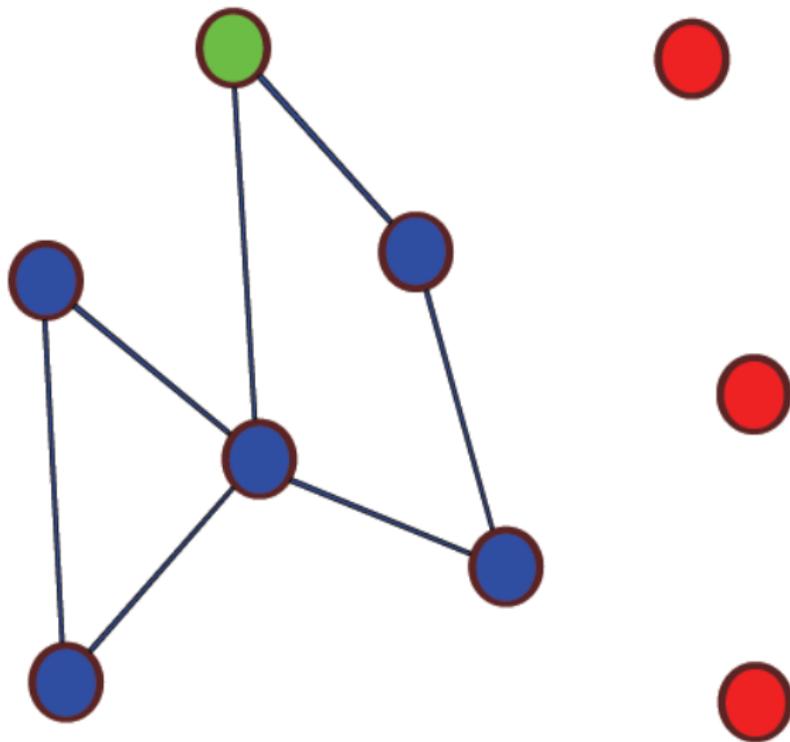
The LOCAL model



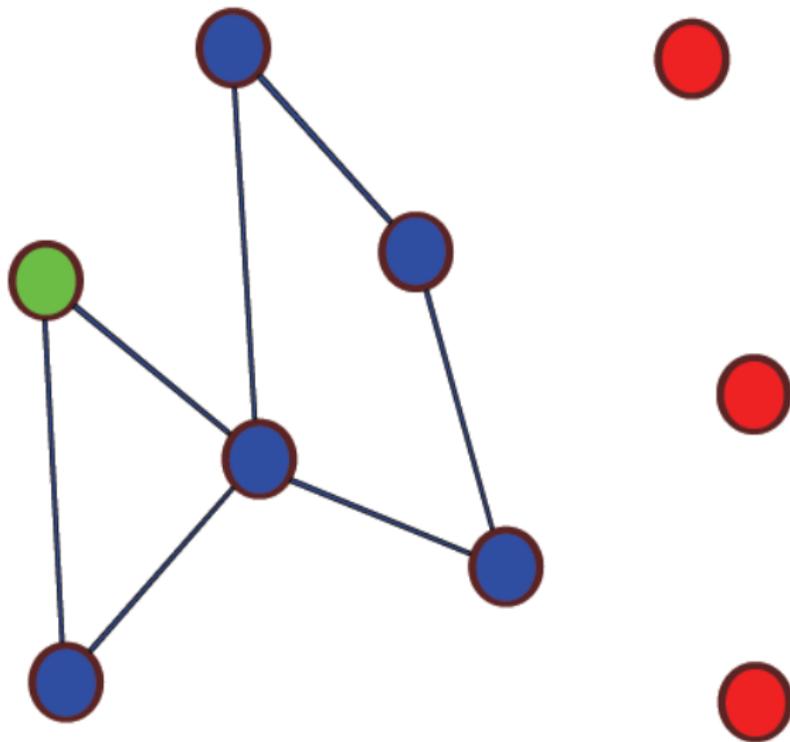
The LOCAL model



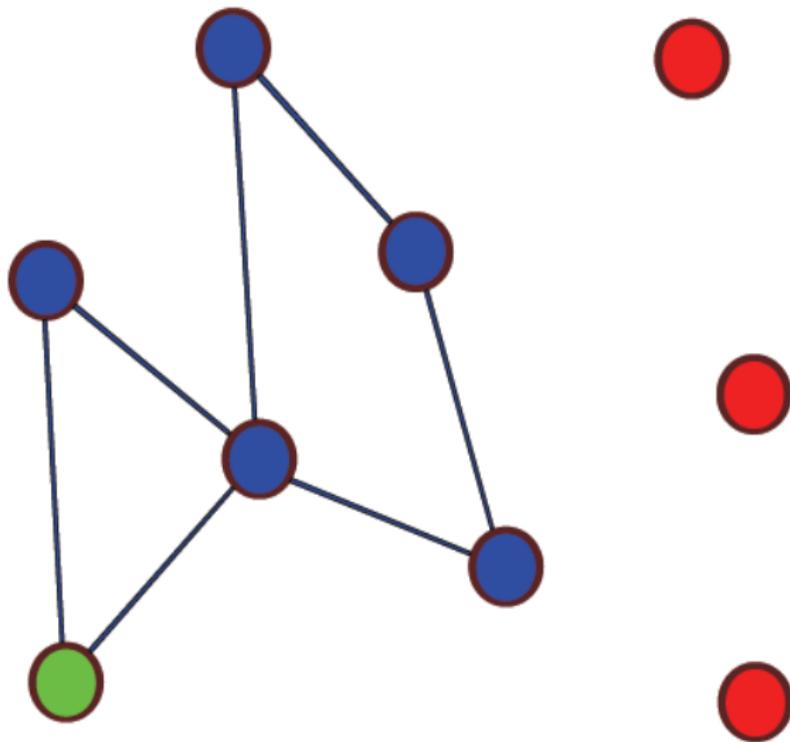
The LOCAL model



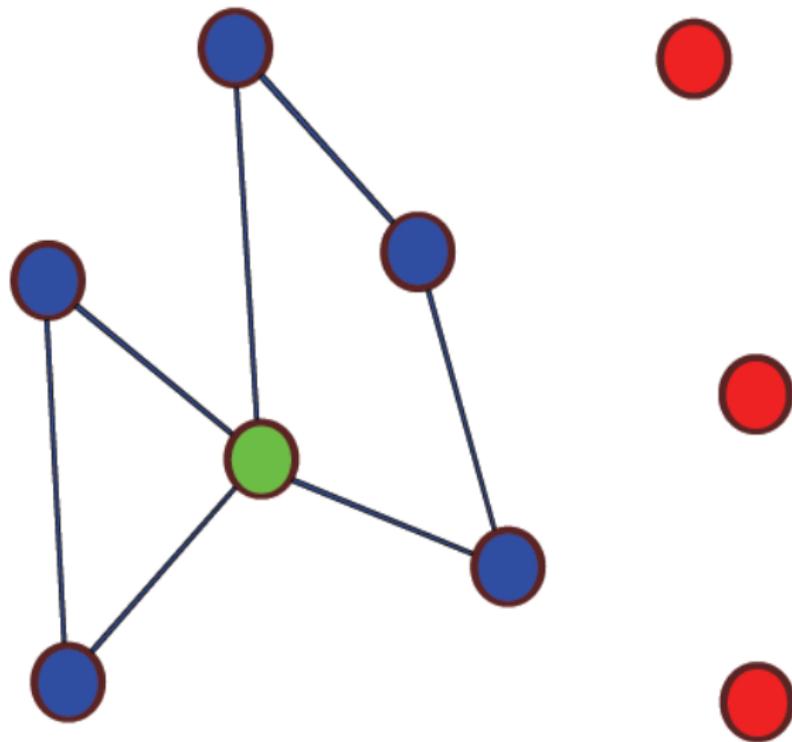
The LOCAL model



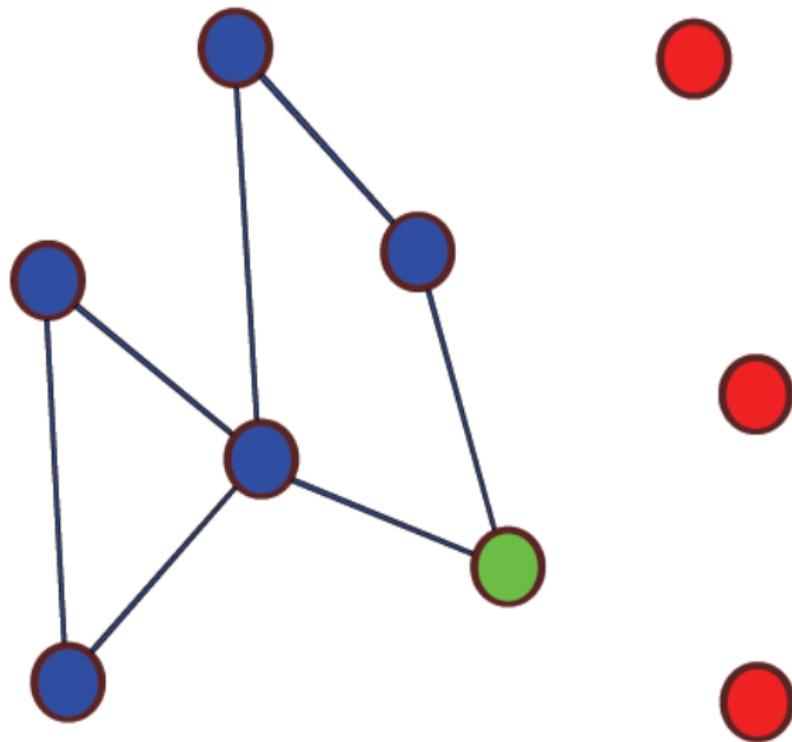
The LOCAL model



The LOCAL model



The LOCAL model



A few remarks about the local approach

- **Weak hypothesis:**
 - if A is connected to B,
 - if C is similar to B,
 - then A is likely to be connected to C.
- **Computationally:** much faster to train N local models with N training points each, than to train 1 model with N^2 training points.
- **Caveats:**
 - each local model may have very few training points
 - no sharing of information between different local models

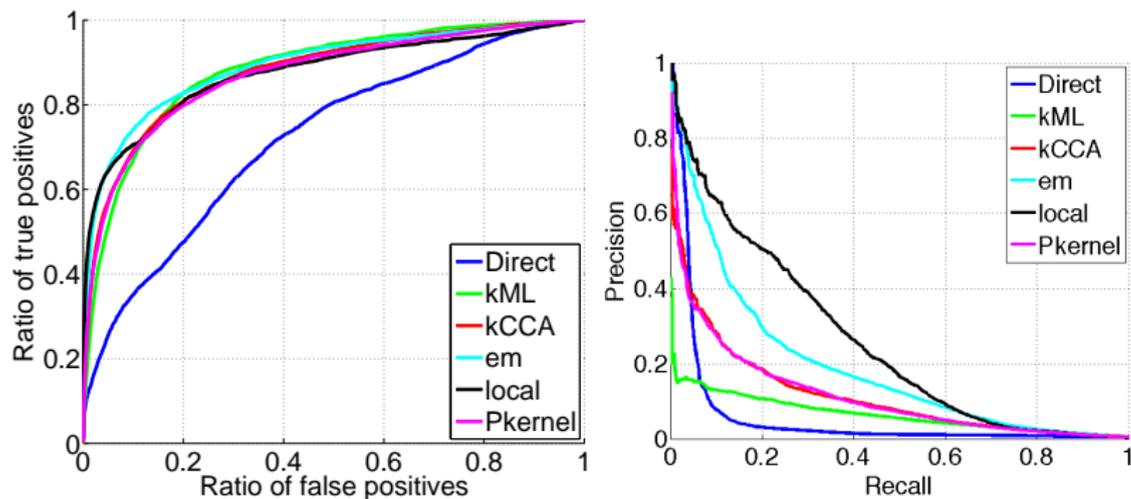
A few remarks about the local approach

- **Weak hypothesis:**
 - if A is connected to B,
 - if C is similar to B,
 - then A is likely to be connected to C.
- **Computationally:** much faster to train N local models with N training points each, than to train 1 model with N^2 training points.
- **Caveats:**
 - each local model may have very few training points
 - no sharing of information between different local models

A few remarks about the local approach

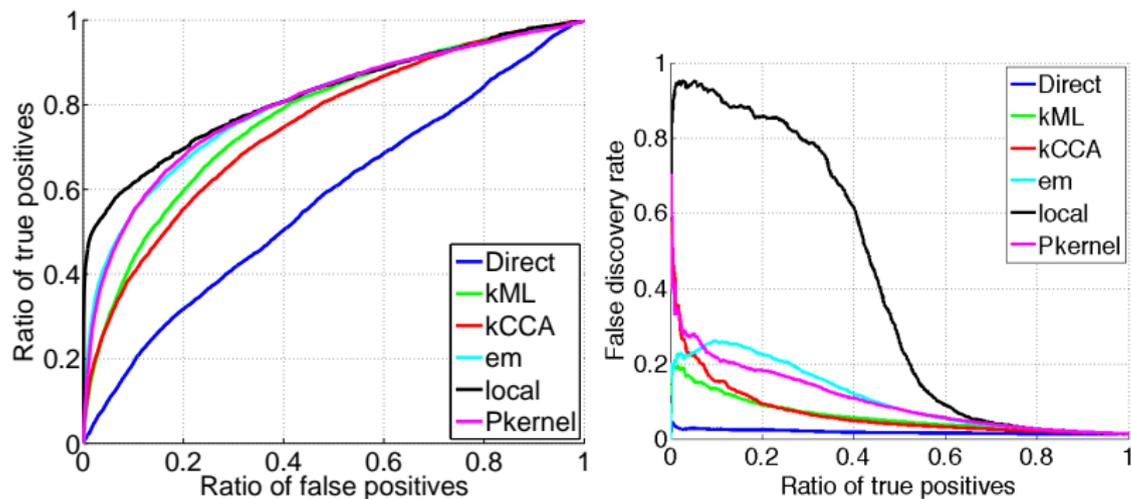
- **Weak hypothesis:**
 - if A is connected to B,
 - if C is similar to B,
 - then A is likely to be connected to C.
- **Computationally:** much faster to train N local models with N training points each, than to train 1 model with N^2 training points.
- **Caveats:**
 - each local model may have very few training points
 - no sharing of information between different local models

Results: protein-protein interaction (yeast)



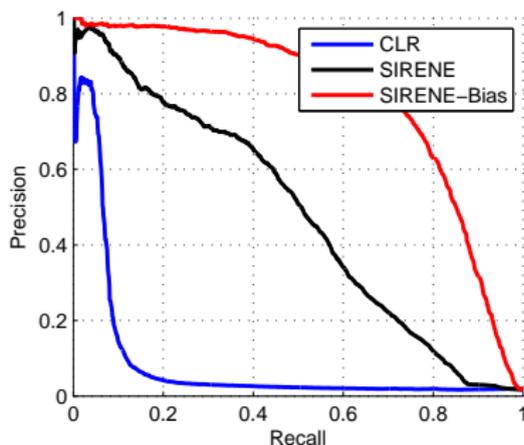
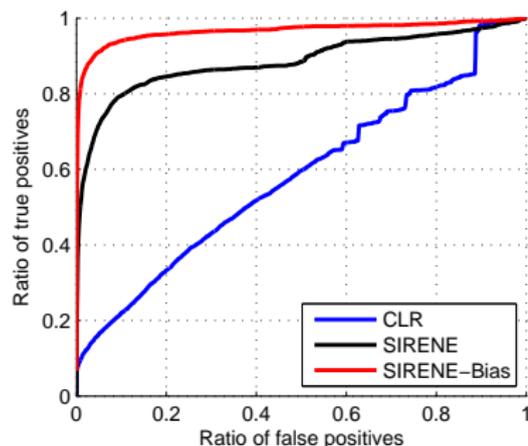
(from Bleakley et al., 2007)

Results: metabolic gene network (yeast)



(from Bleakley et al., 2007)

Results: regulatory network (E. coli)



Method	Recall at 60%	Recall at 80%
SIRENE	44.5%	17.6%
CLR	7.5%	5.5%
Relevance networks	4.7%	3.3%
ARACNe	1%	0%
Bayesian network	1%	0%

SIRENE = Supervised Inference of REgulatory Networks (Mordelet and V., 2008)

Prediction of missing enzyme genes in a bacterial metabolic network

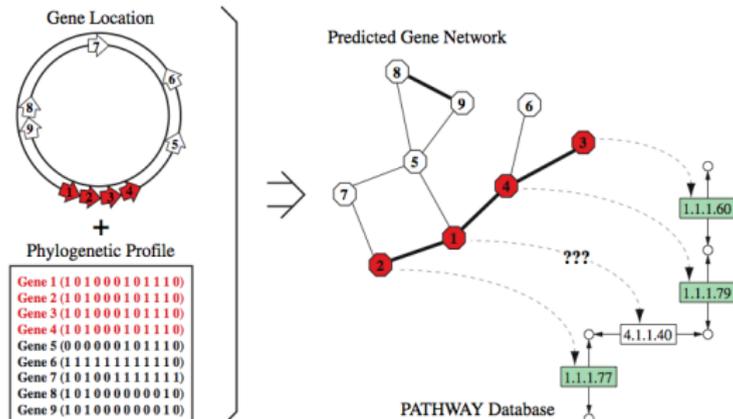
Reconstruction of the lysine-degradation pathway of *Pseudomonas aeruginosa*

Yoshihiro Yamanishi¹, Hisaaki Mihara², Motoharu Osaki², Hisashi Muramatsu³, Nobuyoshi Esaki², Tetsuya Sato¹, Yoshiyuki Hizukuri¹, Susumu Goto¹ and Minoru Kanehisa¹

¹ Bioinformatics Center, Institute for Chemical Research, Kyoto University, Japan

² Division of Environmental Chemistry, Institute for Chemical Research, Kyoto University, Japan

³ Department of Biology, Graduate School of Science, Osaka University, Japan



RESEARCH ARTICLE

Prediction of nitrogen metabolism-related genes in *Anabaena* by kernel-based network analysis

Shinobu Okamoto^{1*}, Yoshihiro Yamanishi¹, Shigeki Ehira², Shuichi Kawashima³, Koichiro Tonomura^{1**} and Minoru Kanehisa¹

¹ Bioinformatics Center, Institute for Chemical Research, Kyoto University, Uji, Japan

² Department of Biochemistry and Molecular Biology, Faculty of Science, Saitama University, Saitama, Japan

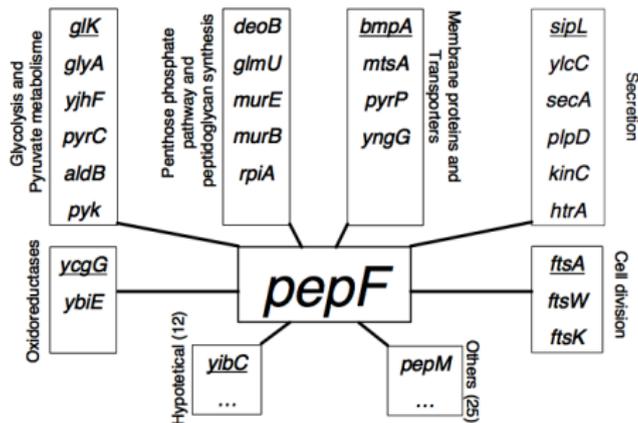
³ Human Genome Center, Institute of Medical Science, University of Tokyo, Meguro, Japan

Determination of the role of the bacterial peptidase PepF by statistical inference and further experimental validation

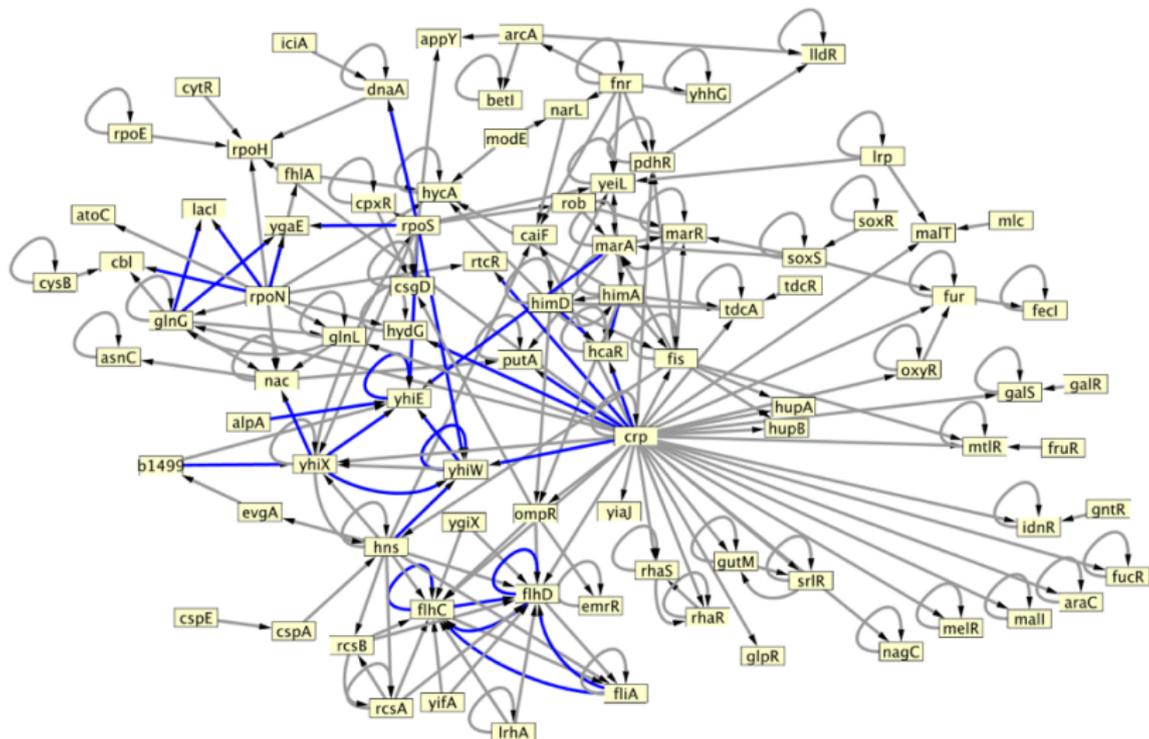
Liliana LOPEZ KLEINE^{1,2}, Alain TRUBUIL¹, Véronique MONNET²

¹Unité de Mathématiques et Informatiques Appliquées. INRA Jouy en Josas 78352, France.

²Unité de Biochimie Bactérienne. INRA Jouy en Josas 78352, France.



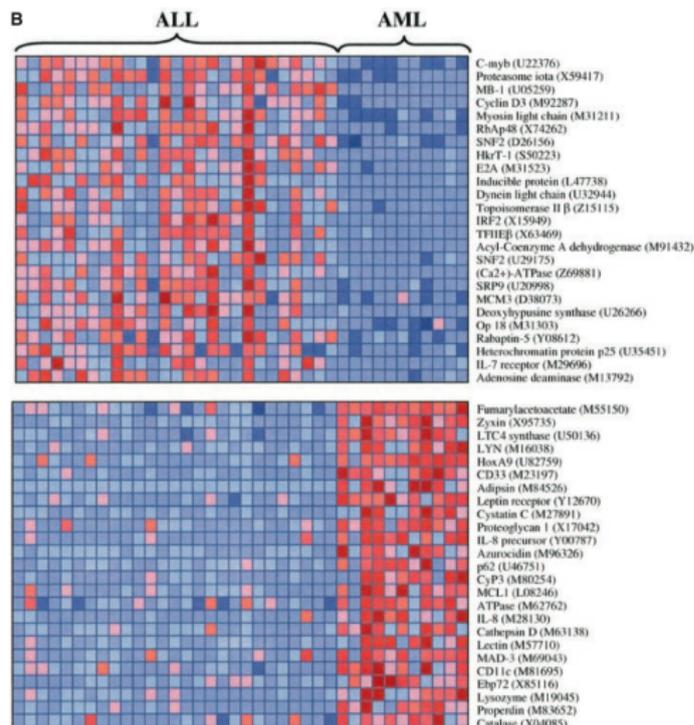
Application: predicted regulatory network (E. coli)



Prediction at 60% precision, restricted to transcription factors (from Mordelet and V., 2008).

- 1 How to infer relationships between genes from biological data?
- 2 How to use biological networks to help in the analysis of genomic data?
- 3 Conclusion

Tissue classification from microarray data



Goal

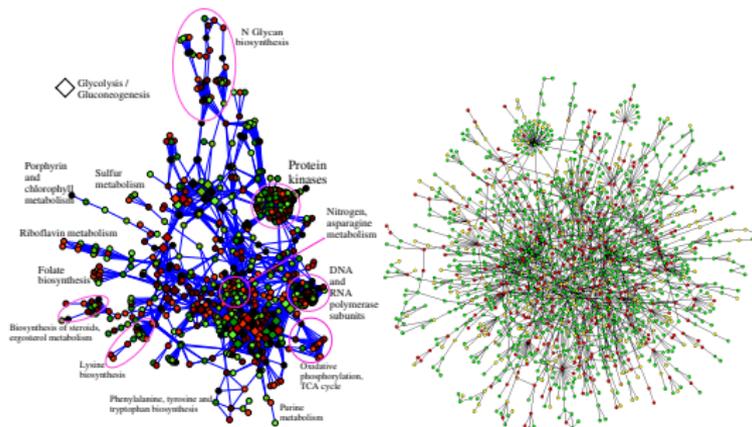
- Design a **classifier** to automatically assign a class to future samples from their expression profile
- **Interpret** biologically the differences between the classes

Issue

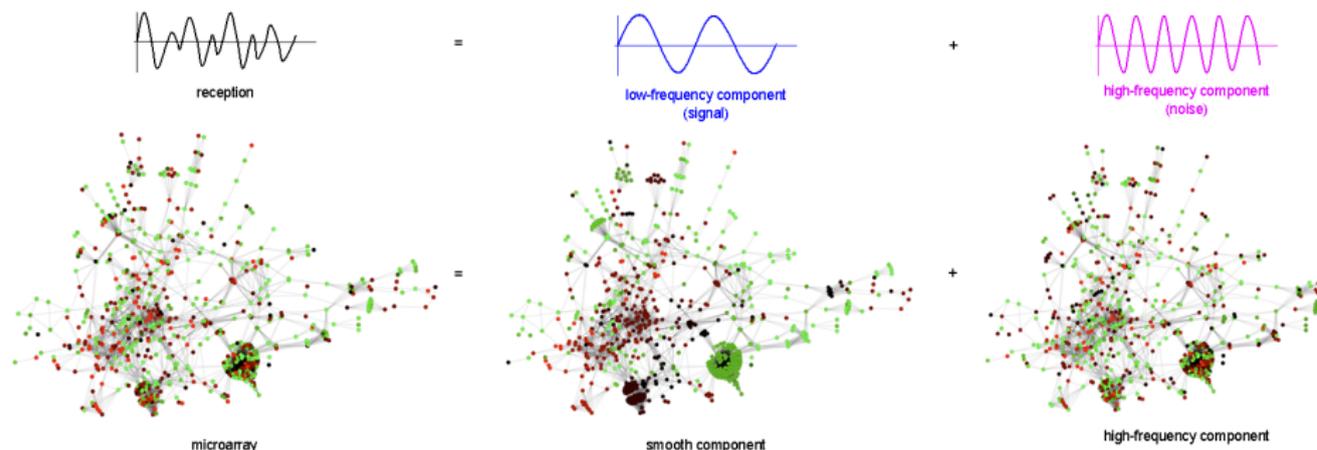
20K+ genes but only <100 tumours

Protein networks can help us

- Basic biological functions usually involve the **coordinated action of several proteins**:
 - Formation of **protein complexes**
 - Activation of metabolic, signalling or regulatory **pathways**
- Many pathways and protein-protein interactions are **already known**
- **Hypothesis**: the weights of the classifier should be “coherent” with respect to this **prior knowledge**



The idea



- 1 Use the gene network to extract the “important information” in gene expression profiles by **Fourier analysis** on the graph
- 2 Learn a linear classifier on the **smooth components**

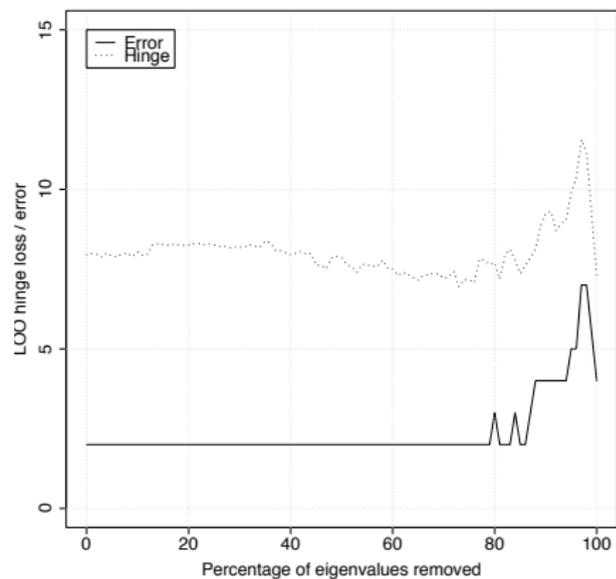
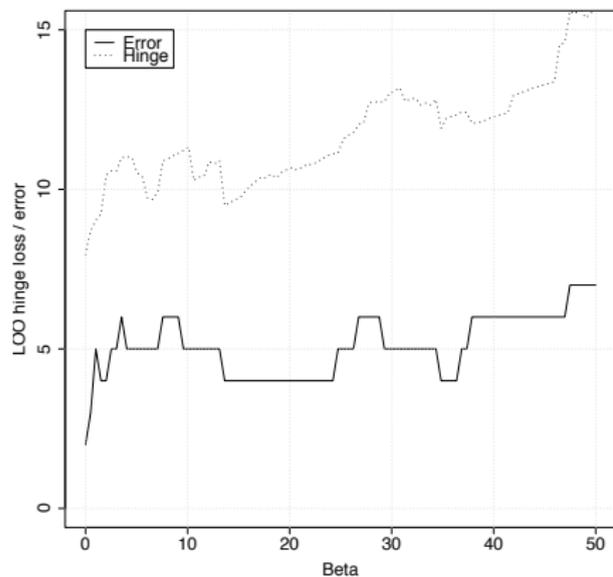
Expression

- Study the effect of low irradiation doses on the yeast
- 12 non irradiated vs 6 irradiated
- Which pathways are involved in the response at the transcriptomic level?

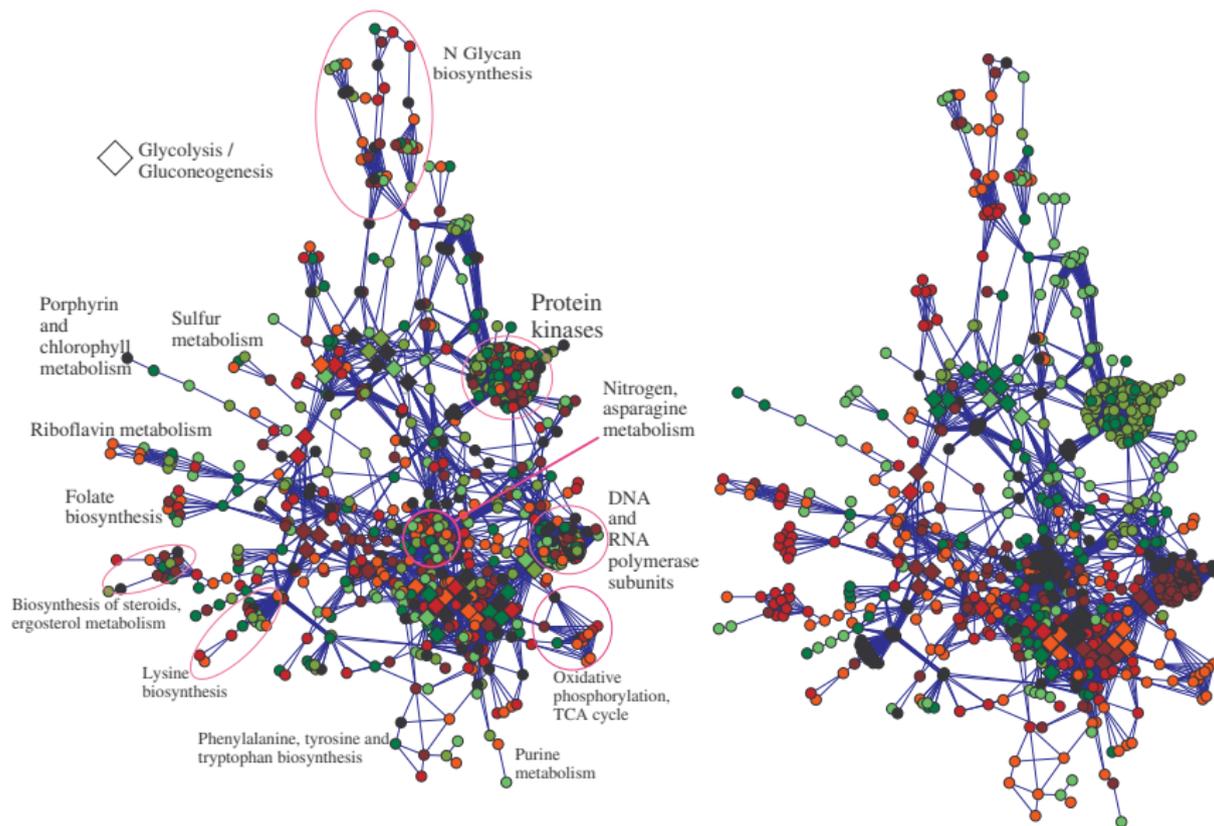
Graph

- KEGG database of metabolic pathways
- Two genes are connected if they code for enzymes that catalyze successive reactions in a pathway (**metabolic gene network**).
- 737 genes, 4694 vertices.

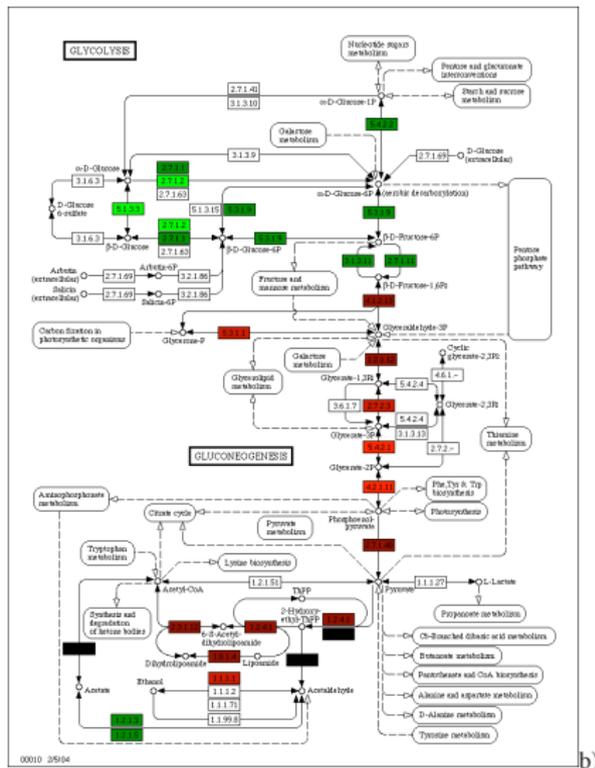
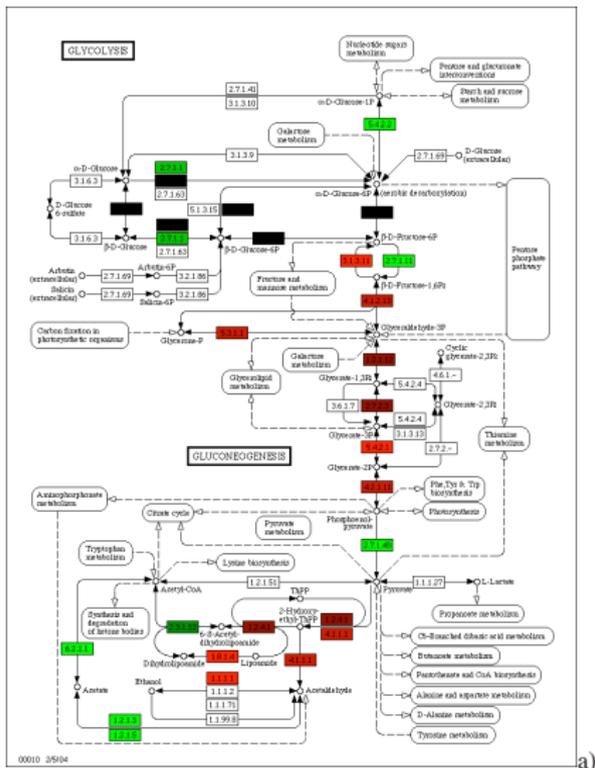
Classification performance



Classifier



Classifier



- 1 How to infer relationships between genes from biological data?
- 2 How to use biological networks to help in the analysis of genomic data?
- 3 Conclusion**

- A supervised machine learning formulation leads to promising results on the problem of inferring unknown relationships between genes and proteins.
- Conversely, biological networks can help fighting the curse of dimensionality
- All this is progression very quickly these days!

People I need to thank



- Graph inference : Yoshihiro Yamanishi, Minoru Kanehisa (Univ. Kyoto), Jian Qian, Bill Noble (Univ. Washington), Kevin Bleakley, Gerard Biau (Univ. Montpellier), Fantine Mordelet (ParisTech/Curie)
- Using graphs : Franck Rapaport, Emmanuel Barillot, Andrei Zinovyev (Institut Curie)

