Inferring regulatory networks from expression data using tree-based methods

<u>Vân Anh HUYNH-THU,</u>

Alexandre Irrthum, Louis Wehenkel, and Pierre Geurts

MLSB 2010, Edinburgh October 2010

Department of Electrical Engineering and Computer Science, Systems and Modeling GIGA-Research, Bioinformatics and Modeling University of Liège, Belgium



Inferring regulatory networks is a challenging problem

unknown network

inferred network



Gene expression can be measured



Input : perturbation to the system (e.g. gene overexpression) Output : measure response to perturbation

Expression data is used to infer the network



A weight is learned for each edge



		Target gene				
		gene 1	gene 2	• • •	gene p	
Regulating gene	gene 1	-	0.05		0.56	
	gene 2	0.19	-		0.03	
	•••	•••	•••			
	gene p	0.11	0.42		-	

Tree-based methods

Network inference : GENIE3



Tree-based methods

Network inference : GENIE3



The inference problem decomposes into p sub-problems



Tree-based methods

Network inference : GENIE3



Tree-based ensemble methods are good candidates



Bagging Random Forests Extra-Trees Non-parametric

Can deal with interacting features

Work well with high-dimensional datasets

Scalable

The tree-based model is informative

The learned model can be used to find the most relevant inputs.



The variable importance is based on variance reduction At each tree node ${\cal N}$:

$$I(\mathcal{N}) = \#S\mathrm{Var}(S) - \#S_t\mathrm{Var}(S_t) - \#S_f\mathrm{Var}(S_f)$$

S: set of samples reaching node \mathcal{N}
 S_t (resp. S_f): subset of S for which the test is true (resp. false)
Var(.): variance of output variable in a subset

For a single tree : $w_i^t = \text{sum of } I$ at each node where variable i appears

For an ensemble of trees :

$$w_i = rac{1}{T}\sum_{t=1}^T w_i^t$$

Tree-based methods

Network inference : GENIE3



GENIE3 uses ensembles of trees to infer a network



A normalization is required

For an unpruned tree :

$$\sum_{i\neq j} w_{i\rightarrow j} \approx N \mathrm{Var}_j(S)$$

 $w_{i \rightarrow j}$: importance of gene i for the prediction of gene j N : number of experiments

 $\operatorname{Var}_{i}(S)$: variance of gene j in the learning sample from which the tree is built



Tree-based methods

Network inference : GENIE3



GENIE3 is best performer in DREAM4 challenge

DREAM4 *In silico Multifactorial network* challenge : inference of **synthetic** regulatory networks.

5 networks of 100 genes, 100 experiments per network.

David	Τ	Mean	Overall	Mean	Overall
капк	ream	AUPR	AUPR p-value	AUROC	AUROC p-value
1	GENIE3-Bagging	0.22	5.93e-54	0.76	1.93e-28
2	Team 549	0.14	7.45e-35	0.73	6.29e-23
• • • •	•••		• • •		•••

AUPR : Area Under Precision-Recall curve

AUROC : Area Under ROC curve

Quality of ranking decreases with in-degree of genes

In-degree = number of regulators



GENIE3 is able to predict a *directed* network

Predicted networks contain a significant number of asymmetric links.



GENIE3 can be used for directing an undirected network

Error on edge directionality :

i
ightarrow j present in gold standard j
ightarrow i not present in gold standard $w_{i
ightarrow j} < w_{j
ightarrow i}$

At 5% recall, mean error rate on edge directionality is 20%.

Results on *E. coli* are competitive to existing approaches

1471 genes, 907 experiments. Validation with RegulonDB.

Input genes = 172 known TFs



 CLR : Faith et al. (2007)
 ARACNE : Margolin et al. (2006)

 MRnet : Meyer et al. (2007)
 GGM : Schafer et al. (2005)

Results on E. coli are competitive to existing approaches

1471 genes, 907 experiments. Validation with RegulonDB.

Input genes = all genes



 CLR : Faith et al. (2007)
 ARACNE : Margolin et al. (2006)

 MRnet : Meyer et al. (2007)
 GGM : Schafer et al. (2005)

Conclusions

Good results for a non parametric approach

Scalable

Can be easily parallelized

Adaptable to other types of genomic data and interactions

Improvement on the way variable importance scores are normalized

Threshold on the ranking of interactions

Comparison with Bayesian networks

Software :

 $http://www.montefiore.ulg.ac.be/{\sim}huynh-thu/software.html$

Related paper :

Inferring regulatory networks from expression data using tree-based methods.

Vân Anh Huynh-Thu, Alexandre Irrthum, Louis Wehenkel, and Pierre Geurts PLoS ONE **5**(9) :e12776