

Machine Learning in Systems Biology (MLSB 2007)

Florence d'Alché-Buc^{1*}, Louis Wehenkel²

¹IBISC CNRS FRE 3190, Université d'Evry-Val d'Essonne, Genopole, Tour Evry II, Evry, FRANCE

²GIGA-R and Dept. of EE&CS, University of Liège, B4000, Liège, BELGIUM

Email: Florence d'Alché-Buc - florence.dalche@ibisc.univ-evry.fr; Louis Wehenkel - L.Wehenkel@ulg.ac.be;

*Corresponding author

Abstract

This article introduces a supplement comprising extended versions of a selected subset of papers presented at the workshop MLSB2007, Machine Learning in Systems Biology, Evry, France, from September 24 to 25, 2007.

Introduction

Molecular biology and also all the biomedical sciences are undergoing a true revolution as a result of the emergence and growing impact of a series of new disciplines/tools sharing the “-omics” suffix in their name. These include in particular genomics, transcriptomics, proteomics and metabolomics devoted respectively to the examination of the entire systems of genes, transcripts, proteins and metabolites present in a given cell or tissue type.

The availability of these new, highly effective tools for biological exploration is dramatically changing the way one performs research in at least two respects. First of all, the amount of available experimental data is not at all a limiting factor any more; on the contrary, there is a plethora of it. The challenge has shifted towards identifying the relevant pieces of information given the question, and how to make sense out of it (a “data mining” issue). Secondly, rather than to focus on components in isolation, we can now try to understand how biological systems behave as the result of the integration and interaction between the

individual components that one can now monitor simultaneously (so called “systems biology”).

Taking advantage of this wealth of “genomic” information has become a *conditio sine qua non* for whoever ambitions to remain competitive in molecular biology and more generally in biomedical sciences. Machine learning naturally appears as one of the main drivers of progress in this context, where most of the targets of interest deal with complex structured objects: sequences, 2D and 3D structures, or interaction networks. At the same time bioinformatics and systems biology have already induced significant new developments of general interest in machine learning, for example in the context of learning with structured data, graph inference, semi-supervised learning, system identification, and novel combinations of optimization and learning algorithms.

The aim of the MLSB 2007 workshop on Machine Learning in Systems Biology, held at University of Evry, France, was to contribute to the cross-fertilization between the research in machine learning methods and their applications to complex biological and medical questions by bringing together method developers and experimentalists.

MLSB 2007, was a follow up of the PMSB 2006 workshop on Probabilistic Modeling and Machine Learning in Structural and Systems Biology, held in Tuusula, Finland, from June 17 to 18, 2006 (see also [4]). It has been followed by MLSB 2008, held in Brussels, Belgium, from September 13 to 14, 2008, and will be further followed by MLSB 2009, taking place in Bled, Slovenia, on September 5 to 6, 2009.

Summary of the supplement

Selected submissions were invited based on the papers presented in the workshop. This supplement contains a reviewed selection of six full papers that cover a large panel of high topics in Machine Learning devoted to Systems Biology.

Aastinen et al. [1] develop kernel methods for enzyme function prediction in the framework of structured output prediction methods, where the enzymatic reaction is the combinatorial target object for prediction. Ying et al. [7] address high throughput analysis of microarray data by using a variational Bayesian inference method for unsupervised clustering that allows latent process variables and model parameters to be dependent. The work of Omont et al. [6] analyzes genome-wide association studies results of Multiple Sclerosis with a new Bayesian model that integrates genotyping errors and genomic structure dependencies. Azé et al. [2] consider annotation of a protein with terms of the functional hierarchy that has been used to annotate *Bacillus subtilis* and learn a set of rules that predict classes in terms of elements of the functional hierarchy using two methods: first-order and multilabel attribute value decision-trees.

Kontos et al. [5] formulate the identification of putative NCR genes in the yeast *Saccharomyces cerevisiae* is formulated as a supervised two-class classification problem and use different classifiers and variable selection methods to predict whether genes are NCR-sensitive or not from a large number of variables related to the GATA motif in the upstream non-coding sequences of the genes.

Birmelé et al. [3] propose to cluster genes by co-regulation rather than by co-expression, present an inference algorithm for detecting co-regulated groups from gene expression data and then introduce a method to cluster genes given that inferred regulatory structure.

Competing interests

The authors declare that they have no competing interests.

Acknowledgements

We wish to thank particularly the local organizing committee, Farida Zehraoui, Pierre Geurts, Nicolas Brunel, Cyril Combe as well as the scientific program committee, Christophe Ambroise, Laurent Bréhelin, Nicolas Brunel, Vincent Frouin, Pierre Geurts, Mark Girolami, Samuel Kaski, Kathleen Marchal, Gunnar Raetsch, Juho Rousu, Céline Rouveirol, Yvan Saeys, Koji Tsuda, Jacques Van Helden, Jean-Philippe Vert, Farida Zehraoui, and Jean-Daniel Zucker.

We thank the University of Evry-Val d'Essonne, the EU FP6 Network of Excellence PASCAL (IST-2002-506778), Genopole, the University of Liège and its center of integrated geno-proteomics GIGA-R, the IAP-network of excellence BIOMAGNET funded by the Belgian Federal Science Policy Office for financial support, and finally, all contributors and participants for the successful workshop.

References

1. K. Astikainen, L. Holm, E. Pitkänen, S. Szedmak, and J. Rousu. Towards structured output prediction of enzyme function. 2008.
2. J. Azé, L. Gentils, C. Toffano-Nioche, V. Loux, J.-F. G. Gibrat, P. Bessières, C. Rouveirol, A. Poupon, and C. Froidevaux. Towards a semi-automatic functional annotation tool based on decision-tree techniques. 2008.
3. E. Birmelé, M. Elati, C. Rouveirol, and C. Ambroise. Identification of functional modules based on transcriptional regulation structure. 2008.

4. S. Kaski, J. Rousu, and E. Ukkonen. Probabilistic modeling and machine learning in structural and systems biology. *BMC Bioinformatics*, 8(Suppl 2):S1, 2007.
5. K. Kontos, P. Godard, B. André, J. van Helden, and G. Bontempi. Machine learning techniques to identify putative genes involved in nitrogen catabolite repression in the yeast *Saccharomyces cerevisiae*. 2008.
6. N. Omont, K. Forner, M. Lamarine, G. Martin, , G. Martin, F. Képès, and W. Jérôme. Gene-based bin analysis of genome-wide association studies. 2008.
7. Y. Ying and C. Campbell. A marginalized variational bayesian approach to the analysis of array data. 2008.