Chapter 60 Relative Relations in Biomedical Data Classification

Marcin Czajkowski

Bialystok University of Technology, Poland

ABSTRACT

Advances in data science continue to improve the precision of biomedical research, and machine learning solutions are increasingly enabling the integration and exploration of molecular data. Recently, there is a strong need for "white box," a comprehensive machine learning model that may actually reveal and evaluate patterns with diagnostic or prognostic value in omics data. In this article, the authors focus on algorithms for biomedical analysis in the field of explainable artificial intelligence. In particular, they present computational methods that address the concept of relative expression analysis (RXA). The classification algorithms that apply this idea access the interactions among genes/molecules to study their relative expression (i.e., the ordering among the expression values, rather than their absolute expression values). One then searches for characteristic perturbations in this ordering from one phenotype to another. They cover the concept of RXA, challenges of biomedical data analysis, and the innovations that the use of relative relationship-based algorithms brings.

INTRODUCTION

Advances in data science continue to improve the precision of biomedical research, and machine learning solutions are increasingly enabling the integration and exploration of molecular data (Huang, Chaudhary & Garmire, 2017). To enable a better understanding of cancer and enhance advances in personalized medicine such data need to be converted to knowledge. An interdisciplinary subfield of computer science called data mining (Han, Kamber & Pei, 2012) aims to reveal important and insightful information hidden in data. It requires appropriate tools and algorithms to effectively identify correlations and patterns within the data. However, the overwhelming majority of systems focus almost exclusively on the prediction accuracy of core data mining tasks like classification or regression. Far less effort has gone into the crucial task of extracting meaningful patterns or molecular signatures of biological processes.

DOI: 10.4018/979-8-3693-3026-5.ch060

Recently, there is a strong need for "white box", comprehensive machine learning models which may actually reveal and evaluate patterns that have diagnostic or prognostic value in biomedical data (McDermott & Wang, 2013). In this chapter, we focus on algorithms for biomedical analysis in the field of eXplainable Artificial Intelligence (XAI) (Angelov, Soares, et. al., 2021). In particular, we present computational methods that address the concept of Relative Expression Analysis (RXA) (Eddy, Sung, et. al., 2010). The algorithms that are based on this idea access the interactions among genes/molecules to study their relative expression, i.e., the ordering among the expression values, rather than their absolute expression values. One then searches for characteristic perturbations in this ordering from one phenotype to another. The simplest form of such an interaction is the ordering of expression among two genes, in which case one seeks to identify typical reversals' pairs of genes (ordering is usually present in one phenotype and rarely present in the other). Such pairs of genes can be viewed as 'biological switches" which can be directly related to regulatory "motifs" or other properties of transcriptional networks. The classification algorithms based on RXA are often data-driven and due to the comparison between feature relative expression levels within the same sample, the predictor is robust to inter- and intra-platforms variabilities as well as complex analytical and data processing methods like normalization and standardization procedures.

The purpose of this chapter is to illustrate the concept of RXA and the innovations that the use of relative relationship-based algorithms brings. We will also cover the issues and challenges of biomedical data analysis.

BACKGROUND

Data mining is an umbrella term covering a broad range of tools and techniques for extracting hidden knowledge from large quantities of data. Biomedical data can be very challenging due to the enormous dimensionality, biological and experimental noise as well as other perturbations. In the literature, we will find that nearly all standard, off-the-shelf techniques were initially designed for other purposes than omics data (Bacardit, Widera, et. al. 2014), such as neural networks, random forests, SVMs, and linear discriminant analysis. When applied for omics data, the prediction models usually involve nonlinear functions of hundreds or thousands of features, many parameters, and are therefore constrain the process of uncovering new biological understanding that, after all, is the ultimate goal of data-driven biology. Deep learning approaches have also been getting attention (Min, Lee & Yoon, 2016) as they can better recognize complex features through representation learning with multiple layers and can facilitate the integrative analysis by effectively addressing the challenges discussed above. However, we know very little about how such results are derived internally. Such lack of knowledge discovery itself in those 'black box' systems impedes biological understanding and are obstacles to mature applications.

In contrast to data mining systems, statistical methods for analyzing high-dimensional biomolecular data generated with high-throughput technologies permeate the literature in computational biology. Those analyses have uncovered a great deal of information about biological processes (Zhao, Shi, et. al., 2015) such as important mutations and lists of "marker genes" associated with common diseases and key interactions in transcriptional regulation. Statistical methods can enhance our understanding by detecting the presence of disease (e.g., "tumor" vs "normal"), discriminating among cancer subtypes (e.g., "GIST" vs "LMS" or "BRCA1 mutation" vs "no BRCA1 mutation") and predicting clinical outcomes (e.g., "poor prognosis" vs "good prognosis"). The statistical analysis is often based on a relatively small

10 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/chapter/relative-relations-biomedical-data-

classification/342578

Related Content

Relevance of Mesh Dimension Optimization, Geometry Simplification and Discretization Accuracy in the Study of Mechanical Behaviour of Bare Metal Stents

Mariacristina Gagliardi (2013). *Methods, Models, and Computation for Medical Informatics (pp. 1-15).* www.irma-international.org/chapter/relevance-mesh-dimension-optimization-geometry/73068

An Overview of Biological Data Mining

Seetharaman Balaji (2017). Library and Information Services for Bioinformatics Education and Research (pp. 130-154).

www.irma-international.org/chapter/an-overview-of-biological-data-mining/176140

Language Focus for Genetics and Molecular Biology Students

Brett Andrew Lidbury (2013). *Bioinformatics: Concepts, Methodologies, Tools, and Applications (pp. 1474-1493).*

www.irma-international.org/chapter/language-focus-genetics-molecular-biology/76129

Matrix Genetics and Culture

Sergey Petoukhovand Matthew He (2010). *Symmetrical Analysis Techniques for Genetic Systems and Bioinformatics: Advanced Patterns and Applications (pp. 248-263).* www.irma-international.org/chapter/matrix-genetics-culture/37905

Biological and Medical Big Data Mining

George Tzanis (2014). *International Journal of Knowledge Discovery in Bioinformatics (pp. 42-56).* www.irma-international.org/article/biological-and-medical-big-data-mining/105100