

# Editorial: Computational Learning Models and Methods Driven by Omics for Precision Medicine

Lei Zhu<sup>1</sup>, Hongmin Cai<sup>1\*</sup>, Fa Zhang<sup>2</sup>, Quan Zou<sup>3</sup>, Yanjie Wei<sup>4</sup>, Huiru Zheng<sup>5</sup>

<sup>1</sup>South China University of Technology, China, <sup>2</sup>Chinese Academy of Sciences (CAS), China, <sup>3</sup>University of Electronic Science and Technology of China, China, <sup>4</sup>Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences (CAS), China, <sup>5</sup>Ulster University, United Kingdom

*Submitted to Journal:*  
Frontiers in Genetics

*Specialty Section:*  
Computational Genomics

*Article type:*  
Editorial Article

*Manuscript ID:*  
620976

*Received on:*  
24 Oct 2020

*Revised on:*  
27 Nov 2020

*Frontiers website link:*  
[www.frontiersin.org](http://www.frontiersin.org)

In Review

---

### *Conflict of interest statement*

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest

### *Author contribution statement*

The article was written by L Zhu. H Cai, F Zhang, Q Zou, Y Wei and H Zheng have provided guidance to the manuscript preparation, have also reviewed and edited the paper. All authors have approved the final version of the editorial.

### *Keywords*

omics, machine learning, drug, RNA, Disease, biomarker, Network analysis, deep learning

### *Contribution to the field*

This is an Editorial, summarizing 34 articles on the topic( Computational Learning Models and Methods Driven by Omics for Precision Medicine ). The article is organized through the following directions: sequencing alignment, correlation detection between omics data and biological traits, prediction of biological functionality, computational methods for cancer subtyping, finding of pathogenic causes, repositions and targeting, and computational methods specially designed for biological knowledge mining. I hope that readers will quickly understand the recent developments of omics data-driven calculation methods and models through this article.

## Editorial: Computational Learning Models and Methods Driven by Omics for Precision Medicine

1 **Lei Zhu<sup>1</sup>, Hongmin Cai<sup>1\*</sup>, Fa Zhang<sup>2</sup>, Quan Zou<sup>3</sup>, Yanjie Wei<sup>4</sup>, Huiru Zheng<sup>5</sup>**

2 <sup>1</sup> School of Computer Science and Engineering, South China University of Technology, Guangzhou,  
3 China

4 <sup>2</sup>Institute of Computing Technology, Chinese Academy of Sciences (CAS), Beijing, China

5 <sup>3</sup>Institute of Fundamental and Frontier Sciences, University of Electronic Science and Technology of  
6 China, Chengdu, China

7 <sup>4</sup>Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences (CAS), Shenzhen,  
8 China

9 <sup>5</sup>School of Computing, Engineering and Intelligent Systems, Faculty of Computing, Engineering and  
10 the Built Environment, Ulster University, Northern Ireland, United Kingdom

### 11 \* Correspondence:

12 Hongmin Cai  
13 hmcai@scut.edu.cn

14 **Keywords: omics, machine learning, drug, RNA, disease, biomarker, network analysis, deep**  
15 **learning**

### 16 Editorial on the Research Topic

#### 17 Computational Learning Models and Methods Driven by Omics for Precision Medicine

18 Due to the high experimental cost and the exponential decline in the cost of high-throughput  
19 sequencing, computational models and methods are preferred by scholars. The curse of  
20 dimensionality is the primary obstacle to dealing with the explosive growth of omics data. Machine  
21 learning methods are applied to reduce dimensionality and perform feature selection from massive  
22 data. Researchers meet the requirements of data sparsity by increasing the sparsity constraints of the  
23 computational models. The models combined with the deep learning method help to discover  
24 potential nonlinear associations. Improving data representation or adding embedding layers could  
25 provide better performance of the models. Computational methods for biomarker discovery, sample  
26 classification and disease process interpretation pave the way for precision medicine.

27 This topic includes 34 papers and a corrigendum. These papers introduce latest researches in the area  
28 of computational biology, catering for precision medicine and complex diseases. They include  
29 sequencing alignment, correlation detection between omics data and biological traits, prediction of  
30 biological functionality, computational methods for cancer subtyping, finding of pathogenic causes,  
31 repositions and targeting, and computational methods specially designed for biological knowledge  
32 mining.

33 **Sequence Alignment:** The raw sequencing data is unstructured short sequences. The structured data  
34 can be generated from downstream analysis through filtering, quality control, and assembly of these

35 unstructured data. Assembly reconciliation can generate high-quality assembly results. In Tang et al.,  
36 using the consensus blocks between contigs to construct adjacency graphs to avoid varying  
37 sequencing depth and sequencing errors, the authors propose a scoring function to rank the input  
38 assembly sets. They use an adjacency algebra model for accurate fusion, which performs well on  
39 *M.abscessus*, *B.fragilis*, *R.sphaeroides* and *V.cholerae*. Shi and Zhang apply the partition and recur  
40 platform to generate a high-level abstraction of the sequence alignments. The algorithm component  
41 library is verified by Apla language. The advantage of implementing the sequence assembly process  
42 through abstract components is that it can effectively improve stability and reduce the possibility of  
43 errors caused by manual selection.

44 **Establishing Omics - Disease Associations:** Four groups present research on RNA association  
45 prediction, including Long non-coding RNA(lncRNA)-protein interactions (LPI), LncRNA-Disease,  
46 microRNA(miRNA)-Disease, and Circular RNA(circRNA)-Disease. Peng et al. give us an overview  
47 of how to identify lncRNA-protein interactions(LPI), and they introduced 16 related repositories and  
48 methods. Among these network-based and deep learning-based methods for predicting LPI, the  
49 proposed SFPEL-LPI used assembly learning and achieved the best Area Under Curve(AUC)  
50 performance. Hu et al. combined the two methods of neural network and matrix factorization (MF) to  
51 predict lncRNA-disease associations. They achieved this combination by concatenating outputs and  
52 sharing inputs between the two methods. Both the MF and the neural network are trained  
53 simultaneously under the framework of TensorFlow. In Yu et al., prior information (lncRNA-miRNA  
54 and lncRNA-disease associations) and known miRNA-disease associations are integrated to construct  
55 a three-layer heterogeneous network of lncRNA, miRNA and disease. In this three-layer network,  
56 the edges between the layers are filled with prior information. Random walk is applied to predict  
57 miRNA-disease associations. The proposed methods are evaluated using cancer data. Their results  
58 show that most potential miRNAs can be confirmed by databases. In Lei et al., the circRNA similarity  
59 network and the disease similarity network are used as the input of the collaboration filtering  
60 recommendation system. Their experiments on predicting potential circRNA-disease associations  
61 indicate the effectiveness of the recommendation system algorithm.

62 Like RNA, microbes and pathogens are also the causes of diseases. In Li et al., a bipartite network is  
63 applied to avoid the omission of neighbor information for predicting Pathogen-Host associations.  
64 Among the top 20 pathogen-host pairs discovered, 16 pairs can be verified by biological experiments.  
65 In Ma et al., to explore the pathogenesis of complex diseases from the modular perspective, the  
66 similarity matrix is decomposed to generate microbe-disease co-modules by nonnegative matrix tri-  
67 factorization. Their method achieves nice performance in the enrichment index and the number of  
68 significantly enriched taxon sets. In Li et al., on the strength of a matrix containing microbes  
69 similarity, disease similarity and a bipartite graph network of the two interactions, the potential  
70 microbe-disease associations are calculated by Katz centrality. The prediction performance was  
71 evaluated by the leave-one-out cross validation and reached an AUC of 0.9098. Zhu et al. use a deep  
72 feedforward network to identify microbial markers and realize graph embedding by replacing the first  
73 two layers of the network with a sparse graph. Experiments show that this Graph Embedding Deep  
74 Feedforward Network has the best performance, comparing deep forest, random forest and Support  
75 Vector Machine(SVM).

76 **Prediction of biological functionality:** Identifying acetylation proteins is conducive to  
77 understanding the post-translational modification process. In Qiu et al., the authors first generate a k-  
78 nearest neighbors (KNN) score, and then use random forest to classify the acetylation proteins. The  
79 formation of KNN scores is based on domain annotation and subcellular localization. Five-fold cross-  
80 validation on the three data sets was performed, and finally, an average AUC of 0.8389 was obtained.

81 In Miao et al., the authors aim to identify which proteins are endoplasmic reticulum-resident proteins,  
82 and they achieved accuracy over 86%. Such work allows us to understand the functionality of  
83 proteins, which may be potential points of drug design. The promoter drives the flow of genetic  
84 information from DNA to RNA, and its sequence information determines the strength of the  
85 promoter. In Le et al., the promoter sequence is divided into 10-gram levels and is used to form a  
86 1000-dimensional vector. The vector is input into a deep neural networks model to classify the  
87 promoter strength. Compared with other latest methods in the same test set, this method improves  
88 1%-4% on all indicators.

89 **Computational approach for cancer subtyping:** Cancer subtyping is fundamental for precision  
90 therapy. Accurately identifying cancer subtypes enables us to understand cancer evolution. In Lu et  
91 al., Laplacian score and low-rank representation methods are integrated to obtain a low-rank  
92 expression of cancer gene expression data. This low rank matrix is hoping to preserve subtype  
93 information. By sorting the obtained matrix, the feature genes are heuristically selected to comprise  
94 of a gene subset for accurate cancer subtyping. The method is tested on five cancer dataset and is  
95 shown to achieve superior performance over k-means, non-negative matrix factorization (NMF) and  
96 several other baseline methods. Aouiche et al. obtained the cancer stages on copy number  
97 variation(CNV) data. The positive significance of distinct stage division is dependent on not only a  
98 high cure rate after cancer been detected, but also on critical markers, which are potential therapeutic  
99 targets. Li et al. identify differentially expressed genes(DEGs) in tumor by analyzing the residues of  
100 each gene via a regression model and found potential biomarkers of the individual sample from  
101 DEGs. Survival analysis is performed on samples collected from human and mouse cancer data, and  
102 is shown to be statistically differently.

103 **Quantitative understanding of pathogenic causes:** The goal of developing computational disease  
104 models is to find a therapeutic target. As the first step, computational tools are required to explain the  
105 cause of the disease. Regarding the identification of Schizophrenia (SZ), Xiang et al. construct a  
106 Brainnetome atlas based on resting-state functional magnetic resonance imaging. Brainnetome atlas  
107 is a weighted undirected graph constructed with brain regions as nodes and correlations as edges. The  
108 authors calculate the features from the atlas and, then use least absolute shrinkage and selection  
109 operator(lasso) learning to prune the features. The classification of SZ is achieved by using SVM with  
110 an accuracy of 93.10%. In Li et al., each single sample is classified by a pathway-based approach,  
111 into Ulcerative colitis (UC) and Crohn's disease (CD). Even though UC and CD have common  
112 clinical characteristics, they have different responses to drugs. According to the gene expression data  
113 of the sample, the author scores each pathway to form a pathway activation for single sample matrix,  
114 which is classified by a random forest classifier. In Zhang et al., the authors aim to select CNV  
115 markers to distinguish between three different states of mono-ADP-ribosylhydrolase 2 (MACROD2).  
116 The frequent deletions of MACROD2 locus may lead to chromosomal instability of human colorectal  
117 cancer. The authors firstly select 17 important single nucleotide polymorphism(SNP) site via mutual  
118 information, and then uses bootstrapping scheme to train multiple classifiers. The trained classifiers  
119 are finally ensembled to effectively distinguish three types of MACROD2. In Lei et al., the  
120 effectiveness of lipoprotein 2 on Subarachnoid hemorrhage (SAH) intervention is revealed from the  
121 perspective of the cell signaling pathway. The authors discover five biomarkers, three of which have  
122 been verified by previous experimental evidence. Finally, the early SAH prediction is performed  
123 based on the assembly learning of logistic regression, SVM and Naive-Bayes, achieving an accuracy  
124 of 79%. Zhang et al. clarify a pathway of polycistronic mRNA ORF73 involved in host apoptosis  
125 through protein p53, supplementing the pathogenic process of Kaposi sarcoma-associated herpes  
126 virus. This work is mainly done through protein-protein interactions (PPI) analysis, Gene Ontology  
127 and Kyoto Encyclopedia of Genes and Genomes pathway analyses. In Shao et al., 108 whole-non-

128 structural protein 5 sequences are analyzed in Zika virus, and 35 potential glycosylation and  
129 phosphorylation sites have been discussed. Mutations in amino acid sites are found to be correlated  
130 with their pathogenicity and transmission efficiency. The relatively stable nucleic acid sequence is  
131 shown to be helpful for detection and vaccine development.

132 A meta-analysis can combine multiple studies, and the two groups apply meta-analysis methods. In  
133 Fukutani et al., after the analysis of Human T-lymphotropic virus 1 (HTLV-1)-infected patients, the  
134 authors find that gene CD40LG and gene GBP2 can be used as two phenotypic classifications of  
135 HTLV-1 infection, with accuracy rates of 0.88 and 1. In Jin and Shi, a meta-analysis is performed to  
136 test SNP-environment interaction. Based on meta-regression (MR), the author proposes overlapping  
137 MR combined with the method of processing overlapping data. This method can reduce type I error  
138 and is more robust than MR in dealing with the nonlinear interaction effect.

139 Gao et al. screen 107 methylomic features in whole blood methylation samples and use Support  
140 Vector Regressor to predict age. What is interesting is that only gene CALB1 and gene KLF14 are  
141 both found in the male and female age prediction models.

142 **Drug repositions and targeting:** Four works focus on drug repositions. In Manibalan et al., the  
143 authors focus on the S100A8 protein, which has a strong interaction with the prevalence of polycystic  
144 ovary syndrome biomarkers. Therefore, they design a series of RNA aptamers targeting the S100A8,  
145 and select the one with minimal binding energy as the targeted drug. Wound Scratch experiments  
146 confirm that the synthesized 18-mer oligo has a significant inhibition effect on tumor cell migration.  
147 Wu et al. hope to level the differences in chemotherapy prognosis through cisplatin resistance  
148 analysis of oral squamous cell carcinoma. Through the analysis of differentially expressed genes, PPI  
149 network and miRNA-mRNA targeted regulatory network, they find that five hub genes and the miR-  
150 200 family members that regulate hub genes may be potential drug targets. In Yu et al., new targeted  
151 drugs for hepatocellular carcinoma (HCC) are found by the drug repositioning bioinformatics  
152 method. Finding HCC's kernel genes is the first step in work. The next step is to combine the  
153 relationship between the drug and gene expression in the Connectivity Map database to score the  
154 relationship between the drug and HCC. Among the top ten drugs screened by this method, eight  
155 drugs have been supported by publications. In Emdadi and Eslahchi, cell line similarity, drug  
156 similarity and half maximal inhibitory concentration are combined to predict the drug sensitivity of  
157 cells, and logistic matrix factorization is applied to obtain latent vectors. For the drug sensitivity  
158 prediction of the new cell line, the k-nearest neighbors of the cell line are estimated through the  
159 decision tree to obtain the latent vectors of the cell line. Finally, a threshold based on the probability  
160 of the latent vector is used to predict whether the cell line is sensitive to drugs. The genomics of drug  
161 sensitivity on haematopoietic cell lines in cancer was tested for model performance, with an accuracy  
162 of 0.721.

163 **Biology-oriented learning methods:** Traditional learning methods have achieved tremendous  
164 success and have provided solutions to even some difficult biological problems. In Wang et al.,  
165 Huber loss is applied to alleviate non-Gaussian noise contaminations. A sparsity penalty item is used  
166 to encourage the sparsity of representation of The Cancer Genome Atlas data, and a graph  
167 regularization is used to preserve the manifold structure. The clustering accuracy is improved by 5%  
168 compared with non-negative matrix factorization. Che et al. improve the traditional methods on the  
169 basis of Sparse Group Lasso (SGL) and proposed a weighted sparse group lasso (WSGL) by  
170 introducing prior constraint on the sparse term. Compared with lasso and SGL, the performance is  
171 significantly improved, indicating that prior biological knowledge carries on valuable message.  
172 Comparing the lasso and SGL methods, WSGL can screen less genes, and the ratio of candidate

173 genes is higher using *Arabidopsis* flowering time data. In Lemaçon et al., a visualization method is  
174 proposed based on a scoring system for rating susceptibility loci. In general, this is a visualization  
175 method for searching for the best potential variants through aggregating prediction approaches. In  
176 Guo, Kullback-Leibler divergence is used to measure the distance between two SNPs, and these  
177 distances are used as k-means clustering. Then, statistical testing methods are applied to find epistatic  
178 interactions, and the time cost of this method is about one-tenth that of Bayesian inference-based  
179 method. Zheng et al. use sparse subspace clustering to perform single-cell clustering. This method  
180 assumes that the feature vector of a sample can be expressed as a linear combination of other samples  
181 in the same subspace. In the test of ten single-cell datasets, this method maintains the leading position  
182 in normalized mutual information and adjusted rand index.

183 These teams work together to continuously improve model accuracy. Most articles related to  
184 computational methods are tailored from early established models for biology knowledge learning.

### 185 **Author Contributions**

186 The article was written by L Zhu. H Cai, F Zhang, Q Zou, Y Wei and H Zheng have provided  
187 guidance to the manuscript preparation, have also reviewed and edited the paper. All authors have  
188 approved the final version of the editorial.

### 189 **Acknowledgments**

190 We thank all the authors who contributed to this topic.