

WOODHEAD PUBLISHING SERIES IN BIOMEDICINE



DATA MINING FOR BIOINFORMATICS APPLICATIONS

ZENGYOU HE



Woodhead Publishing Series in
Biomedicine: Number 76

Data Mining for Bioinformatics Applications

Zengyou He



AMSTERDAM • BOSTON • CAMBRIDGE • HEIDELBERG
LONDON • NEW YORK • OXFORD • PARIS • SAN DIEGO
SAN FRANCISCO • SINGAPORE • SYDNEY • TOKYO

Woodhead Publishing is an imprint of Elsevier



Woodhead Publishing Limited is an imprint of Elsevier
80 High Street, Sawston, Cambridge, CB22 3HJ, UK
225 Wyman Street, Waltham, MA 02451, USA
Langford Lane, Kidlington, OX5 1GB, UK

© 2015 Elsevier Ltd. All rights reserved.

No part of this publication may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording, or any information storage and retrieval system, without permission in writing from the publisher. Details on how to seek permission, further information about the Publisher's permissions policies and our arrangements with organizations such as the Copyright Clearance Center and the Copyright Licensing Agency, can be found at our website: www.elsevier.com/permissions.

This book and the individual contributions contained in it are protected under copyright by the Publisher (other than as may be noted herein).

Notices

Knowledge and best practice in this field are constantly changing. As new research and experience broaden our understanding, changes in research methods, professional practices, or medical treatment may become necessary.

Practitioners and researchers must always rely on their own experience and knowledge in evaluating and using any information, methods, compounds, or experiments described herein. In using such information or methods they should be mindful of their own safety and the safety of others, including parties for whom they have a professional responsibility.

To the fullest extent of the law, neither the Publisher nor the authors, contributors, or editors, assume any liability for any injury and/or damage to persons or property as a matter of products liability, negligence or otherwise, or from any use or operation of any methods, products, instructions, or ideas contained in the material herein.

ISBN: 978-0-08-100100-4 (print)

ISBN: 978-0-08-100107-3 (online)

British Library Cataloguing-in-Publication Data

A catalogue record for this book is available from the British Library

Library of Congress Control Number: 2015934370

For information on all Woodhead Publishing publications
visit our website at <http://store.elsevier.com/>

Printed in the United States of America



**Working together
to grow libraries in
developing countries**

www.elsevier.com • www.bookaid.org

Data Mining for Bioinformatics Applications

Related titles

From plant genomics to plant biotechnology
(ISBN: 978-1-907568-29-9)

An introduction to biotechnology
(ISBN: 978-1-907568-28-2)

MATLABs in bioscience and biotechnology
(ISBN: 978-1-907568-04-6)

About the author

Zengyou He is an associate professor in the School of Software at Dalian University of Technology, China.

He received his BS, MS, and PhD degrees in computer science from the Harbin Institute of Technology, China, in 2000, 2002, and 2006, respectively. Prior to becoming an associate professor, he was a research associate in the Department of Electronic and Computer Engineering at the Hong Kong University of Science and Technology (2007–2010).

His research interests include computational proteomics and biological data mining. He has published more than 30 papers on leading journals in the field of bioinformatics, including *Bioinformatics*, *BMC Bioinformatics*, *Briefings in Bioinformatics*, *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, and *Journal of Computational Biology*.

Dedication

Introduction

Data mining methods have been widely used for solving real bioinformatics problems. However, the data mining process is not trivial. It consists of many steps: problem definition, data collection, data preprocessing, modeling, and validation. For each step, different techniques may be applied. Due to the complexity of data mining process and data mining methods, people cannot easily use data mining tools to solve their bioinformatics problems.

In this book, I will use an example-based method to illustrate how to apply data mining techniques to solving real bioinformatics problems. More precisely, I will use six bioinformatics problems that have been investigated in my recent research as examples. For each example, I will describe the entire data mining process, ranging from data preprocessing to modeling and result validation. In addition, I will describe how to use *different* data mining methods to solve the *same* bioinformatics problem in some examples.

In this problem-driven book, I will cover the most commonly used data mining methods, such as frequent pattern mining, discriminative pattern mining, classification, and clustering to show how to select one feasible data mining method to solve a real bioinformatics problem at hand.

Audience

This book will have obvious appeal for a broad audience of computer scientists who are interested in designing new data mining algorithms and biologists who are trying to solve bioinformatics problems using existing data mining tools. To achieve this objective, this book is organized with the following distinct features.

- Providing an example-based description on the whole data mining process for bioinformatics applications. This is distinct from method-based description, in which the chapters are organized according to different data mining techniques. Such an example-based organization is beneficial as it may help the readers to understand how to solve a real problem at hand by choosing proper data mining methods.
- Covering most popular data mining techniques throughout the book. Currently, there are many data mining methods in the literature. This book covers most of them and shows their applications in practical bioinformatics problems.
- Giving detailed illustrations and examples of how to use different data mining techniques to solve the same bioinformatics problem. Due to the complex nature of bioinformatics problems, the same problem can be solved using different data mining techniques. Different

solutions vary from underlying assumptions to algorithmic details. Such kinds of examples will not only enable the reader to understand the target problem more deeply, but also provide hints on how to apply data mining methods in his or her future bioinformatics research. Using frontier bioinformatics problems as examples in each chapter. All the examples discussed in this book will be frontier bioinformatics problems that are under investigation by the author and other researchers. Students who are interested in developing new and better algorithms for these problems may use this book as a starting point.

Acknowledgments

Many articles and books have been referenced in the writing of this book, and citations have been given for these works. To enhance the readability, I have tried to minimize literature references in the text. The sources and some additional literature that may be interesting to the readers are included in the reference section at the end of each chapter. I hope that this organization will provide the authors of the source literature with the appropriate acknowledgments. Moreover, this work was partially supported by the Natural Science Foundation of China under the Grant No. 61003176, the Fundamental Research Funds for the Central Universities of China under the Grant No. DUT14QY07.

This book could not have been written without the help of many people. First of all, most contents in this book are based on the research articles coauthored by my students and myself. Therefore, I would like to first thank my former and current students: Ting Huang, Haipeng Gong, Yang Zhang, Ben Teng, Xiaoqing Liu, Jun Wu, Can Zhao, and Feiyang Gu. In particular, I am grateful to Ben Teng and Xiaoqing Liu, who have carefully read the book and given very detailed comments and suggestions.

In addition, I would like to thank my former colleagues at the Hong Kong University of Science and Technology: Weichuan Yu, Can Yang, Chao Yang, and Xiang Wan, who have assisted me with the study of bioinformatics problems that have contributed to this book.

Furthermore, I am especially thankful to my PhD supervisor, Professor Xiaofei Xu, for his continuous efforts on “transforming” my research style from algorithm-driven research to application-driven research. Without such a transition, I would not have been able to write and finish this book.

Last, but definitely not least, I am indebted to my parents, my wife, and my son for their continuous support and patience.

*Zengyou He
Dalian, China*

Contents

List of figures	vii
List of tables	xi
About the author	xiii
Dedication	xv
1 An overview of data mining	1
1.1 What's data mining?	1
1.2 Data mining process models	1
1.3 Data collection	1
1.4 Data preprocessing	2
1.5 Data modeling	3
1.6 Model assessment	8
1.7 Model deployment	9
1.8 Summary	9
References	10
2 Introduction to bioinformatics	11
2.1 A primer to molecular biology	11
2.2 What is bioinformatics?	11
2.3 Data mining issues in bioinformatics	12
2.4 Challenges in biological data mining	17
2.5 Summary	17
References	17
3 Phosphorylation motif discovery	19
3.1 Background and problem description	19
3.2 The nature of the problem	20
3.3 Data collection	20
3.4 Data preprocessing	21
3.5 Modeling: A discriminative pattern mining perspective	21
3.6 Validation: Permutation p -value calculation	24
3.7 Discussion and future perspective	26
References	27
4 Phosphorylation site prediction	29
4.1 Background and problem description	29
4.2 Data collection and data preprocessing	29

4.3	Modeling: Different learning schemes	32
4.4	Validation: Cross-validation and independent test	35
4.5	Discussion and future perspective	35
	References	36
5	Protein inference in shotgun proteomics	39
5.1	Introduction to proteomics	39
5.2	Protein identification in proteomics	40
5.3	Protein inference: Problem formulation	40
5.4	Data collection	41
5.5	Modeling with different data mining techniques	41
5.6	Validation: Target-decoy versus decoy-free	44
5.7	Discussion and future perspective	48
	References	48
6	PPI network inference from AP-MS data	51
6.1	Introduction to protein–protein interactions	51
6.2	AP-MS data generation	51
6.3	Data collection and preprocessing	52
6.4	Modeling with different data mining techniques	52
6.5	Validation	56
6.6	Discussion and future perspective	57
	References	58
7	Protein complex identification from AP-MS data	61
7.1	An introduction to protein complex identification	61
7.2	Data collection and data preprocessing	61
7.3	Modeling: A graph clustering framework	61
7.4	Validation	67
7.5	Discussion and future perspective	68
	References	68
8	Biomarker discovery	69
8.1	An introduction to biomarker discovery	69
8.2	Data preprocessing	69
8.3	Modeling	70
8.4	Validation	75
8.5	Case study	76
8.6	Discussion and future perspective	77
	References	78
	Conclusions	79
	Index	81

List of figures

Figure 1.1	Typical phases involved in a data mining process model.	2
Figure 2.1	An example of the alignment of five biological sequences. Here “-” denotes the gap inserted between different residues.	13
Figure 3.1	Overview of the Motif-All algorithm. In the first phase, it finds frequent motifs from P to reduce the number of candidate motifs. In the second phase, it performs the significance testing procedure to report all statistically significant motifs to the user.	22
Figure 3.2	Overview of the C-Motif algorithm. The algorithm generates and tests candidate phosphorylation motifs in a breath-first manner, where the support and the statistical significance values are evaluated simultaneously.	23
Figure 3.3	The calculation of conditional significance in C-Motif. In the figure, $\text{Sig}(m, P(m_i), N(m_i))$ denotes the new significance value of m on its i th submotif induced data sets.	23
Figure 4.1	An illustration on the training data construction methods for non-kinase-specific phosphorylation site prediction. Here the shadowed part denotes the set of phosphorylated proteins and the unshadowed area represents the set of unphosphorylated proteins.	30
Figure 4.2	An illustration on the training data construction methods for kinase-specific phosphorylation site prediction. The proteins are divided into three parts: (I) the set of proteins that are phosphorylated by the target kinase, (II) the set of proteins that are phosphorylated by the other kinases, and (III) the set of unphosphorylated proteins.	31
Figure 4.3	An illustration on the basic idea of the active learning procedure for phosphorylation site prediction. (a) The SVM classifier (solid line) generated from the original training data. (b) The new SVM classifier (dashed line) built from the enlarged training data. The enlarged training data are composed of the initial training data and a new labeled sample.	33
Figure 4.4	An overview of the PHOSFER method. The training data are constructed with peptides from both soybean and other organisms, in which different training peptides have different weights. The classifier (e.g., random forest) is built on the training data set to predict the phosphorylation status of remaining S/T/Y residues in the soybean organism.	34
Figure 5.1	The protein identification process. In shotgun proteomics, the protein identification procedure has two main steps: peptide identification and protein inference.	40
Figure 5.2	An overview of the BagReg method. It is composed of three major steps: feature extraction, prediction model construction, and prediction result combination. In feature extraction, the BagReg method generates five features that are highly correlated with the presence probabilities of proteins. In prediction model construction, five classification models are built and applied to predict the presence probability of proteins, respectively. In	

	prediction result combination, the presence probabilities from different classification models are combined to obtain a consensus probability.	41
Figure 5.3	The feature extraction process. Five features are extracted from the original input data for each protein: the number of matched peptides (MP), the number of unique peptides (UP), the number of matched spectra (MS), the maximal score of matched peptides (MSP), and the average score of matched peptides (AMP).	42
Figure 5.4	A single learning process. Each separate learning process accomplishes a typical supervised learning procedure. The model construction phase involves constructing the training set and learning the classification model. And the prediction phase is to predict the presence probabilities of all candidate proteins with the classifier obtained in the previous phase.	43
Figure 5.5	The basic idea of ProteinLasso. ProteinLasso formulates the protein inference problem as a minimization problem, where y_i is the peptide probability, D_i represents the vector of peptide detectabilities for the i th peptide, x_j denotes the unknown protein probability of the j th protein, and λ is a user-specified parameter. This optimization problem is the well-known Lasso regression problem in statistics and data mining.	44
Figure 5.6	The target-decoy strategy for evaluating protein inference results. The MS/MS spectra are searched against the target-decoy database, and the identified proteins are sorted according to their scores or probabilities. The false discovery rate at a threshold can be estimated as the ratio of the number of decoy matches to that of target matches.	45
Figure 5.7	An overview of the decoy-free FDR estimation algorithm.	46
Figure 5.8	The correct and incorrect procedure for assessing the performance of protein inference algorithms. In model selection, we cannot use any ground truth information that should only be visible in the model assessment stage. Otherwise, we may overestimate the actual performance of inference algorithms.	47
Figure 6.1	A typical AP-MS workflow for constructing PPI network. A typical AP-MS study performs a set of experiments on bait proteins of interest, with the goal of identifying their interaction partners. In each experiment, a bait protein is first tagged and expressed in the cell. Then, the bait protein and their potential interaction partners (prey proteins) are affinity purified using AP. The resulting proteins (both bait and prey proteins) are digested into peptides and passed to tandem mass spectrometer for analysis. Peptides are identified from the MS/MS spectra with peptide identification algorithms and proteins are inferred from identified peptides with protein inference algorithms. In addition, the label-free quantification method such as spectral counting is typically used to estimate the protein abundance in each experiment. Such pull-down bait-prey data from all AP-MS runs are used to filter contaminants and construct the PPI network.	52
Figure 6.2	A sample AP-MS data set with six purifications.	54
Figure 6.3	The PPI network constructed from the sample data. Here DC is used as the correlation measure and the score threshold is 0.5, that is, a protein pair is considered to be a true interaction if the DC score is above 0.5. In the figure, the width of the edge that connects two proteins is proportional to the corresponding DC score.	55

Figure 6.4	An illustration of database-free method for validating the interaction prediction results. Under the null hypothesis that each bait protein captures a prey protein is a random event, some simulated data sets are generated such that they are comparable to the original one. Then, an empirical p -value representing the probability that an original interaction score for a protein pair would occur in the random data sets by chance can be calculated. Finally, the false discovery rate is calculated according to these p -values.	58
Figure 7.1	An example bait–prey graph. In this figure, each B_i ($i = 1, 2, 3, 4$) denotes a bait protein and each P_i ($i = 1, 2, 3, 4, 5, 6$) represents a prey protein. The score that measures interaction strength between a bait–prey pair is provided as well.	63
Figure 7.2	Three maximal bicliques are identified. Among these three bicliques, $C1$ and $C2$ are reliable and only $C1$ is finally reported as a protein–complex core.	63
Figure 7.3	The final protein complex by including both the protein complex core $C1$ and an attachment $B3$.	64
Figure 8.1	A typical data analysis pipeline for biomarker discovery from mass spectrometry data. In this workflow, there are three preprocessing steps: feature extraction, feature alignment, and feature transformation. After preprocessing the raw data, feature selection techniques are employed to identify a subset of features as the biomarker.	70
Figure 8.2	An illustration of feature transformation based on protein–protein interaction (PPI) information. The PPI information is used to find groups of correlated features in terms of proteins. These identified feature groups are transformed into a set of new features for biomarker identification.	71
Figure 8.3	Filter methods for feature selection. In the filter method, the goodness of a feature subset is evaluated using only the intrinsic properties of the data.	71
Figure 8.4	Wrapper methods for feature selection. In the wrapper method, the feature subset selection is based on the performance of a classification algorithm.	72
Figure 8.5	Embedded methods for feature selection. In the embedded method, the selection of feature subset is integrated with the construction of the classifier.	72
Figure 8.6	An illustration of the FBTC method.	73
Figure 8.7	Stability evaluation by random sampling the original training data. Suppose multiple random subsets of the original training data set are generated. For each random subset of samples, the feature selection method is used to identify a feature subset. If the candidate biomarker (feature subset) identified from the original data has good stability, then it should occur frequently in the set of feature subsets obtained from the randomly selected data sets.	77

List of tables

Table 1.1	An example data set with eight samples and five features	3
Table 3.1	A sample data set used for phosphorylation motif discovery	20
Table 3.2	The original foreground data P and the background data N before the permutation	25
Table 3.3	One example randomized data set after permutation	25
Table 6.1	Some AP-MS data sets available online	53
Table 6.2	The transformed type-value table of the sample data in Figure 6.2	54
Table 6.3	Some high-quality protein–protein interaction database available online	57
Table 7.1	The binary purification matrix	64
Table 7.2	The adjacency matrix derived from Table 7.1	65
Table 8.1	Confusion matrix defines four possible scenarios when classifying samples in the context of biomarker discovery	74
Table 8.2	Performance metrics for evaluating classifiers	74

An overview of data mining

1

1.1 What's data mining?

Data mining lies at the intersection of computer science, optimization, and statistics, and often appears in other disciplines. Generally, data mining is the process of searching for knowledge in data from different perspectives. Here knowledge can refer to any kinds of summarized or unknown information that are hidden underlying the raw data. For instance, it can be a set of discriminative rules generated from the data collected on some patients of a certain disease and healthy people. These rules can be used for predicting the disease status of new patients.

In general, data mining tasks can be classified into two categories: *descriptive* and *predictive*. Descriptive mining tasks characterize a target data set in concise, informative, discriminative forms. Predictive mining tasks conduct the induction and inference on the current data to make future predictions.

1.2 Data mining process models

Data mining is an iterative process that consists of many steps. There are already some generic reference models on the data mining process, such as the Cross Industry Standard Process for Data Mining (CRISP-DM) process model. From a data-centric perspective, these models are structured as sequences of steps to transform the raw data into information or knowledge that is practically useful. As shown in Figure 1.1, a data mining process model typically involves the following phases: data collection, data preprocessing, data modeling, model assessment, and model deployment.

1.3 Data collection

The first step in the data mining process is to collect the relevant data according to the analysis goal in the applications. Generally, all the data that are helpful to achieve the objective in the analysis should be included. The key point here is how to define and understand the rather subjective term of “relevant data.” Its correct interpretation highly depends on our understanding of the target problem and application background. Although this point will be further illustrated in subsequent chapters, we offer some general remarks here:

- In some cases, people definitely know that some kinds of data are highly relevant to the data mining task at hand. However, the acquisition of such data is very difficult or even impossible due to the device deficiency or cost. For instance, to accurately identify peptides in mass-spectrometry-based shotgun proteomics, it is necessary to generate at least one mass

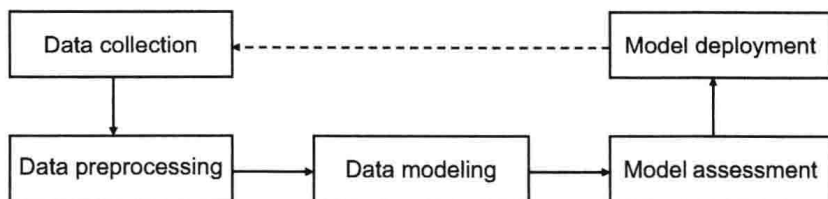


Figure 1.1 Typical phases involved in a data mining process model.

spectrum for each peptide in the sample. However, due to the limitation of current mass spectrometers, it is not always possible to obtain mass spectra data that can cover all peptides present in the sample.

- On the other hand, the inclusion of new relevant data may significantly change the models and methods in the consequent steps of the data mining process. Furthermore, it is necessary to check thoroughly if the use of more relevant data will boost the performance of data mining procedures.

1.4 Data preprocessing

The objective of data preprocessing is twofold: (1) The real-world data are usually low quality; hence preprocessing is used to improve the quality of data, and consequently, the quality of data mining results. (2) In the data modeling step, some specific modeling algorithms cannot operate on the raw data, which should be transformed into some predefined data formats.

There are several general-purpose data preprocessing methods: *data cleaning*, *data integration*, *data reduction*, and *data transformation*.

Data cleaning: Real-world data are usually noisy, inconsistent, and incomplete. Data cleaning procedures aim at removing the noise, correcting inconsistencies, and filling in missing values in the data.

Data integration: In the data collection phase, data sets from different sources are relevant to the analysis problem. Data integration merges data from different sources into an integrated data set for subsequent data mining analysis. The main objective of data integration is to reduce and avoid redundancies and inconsistencies in the resulting data set.

Data reduction: The purpose of data reduction is to generate a new yet smaller representation of the original data set. Generally, the reduced data should contain approximately the same information of the original data that is of primary importance to the analysis target. The most commonly used data reduction technique includes dimension reduction (vertically, reduce the number of features) and sampling (horizontally, reduce the number of samples).

Data transformation: Different data mining algorithms may require different forms of data. Data transformation techniques consolidate the original data into forms appropriate for subsequent mining tasks. For instance, data normalization will transform

the feature values into a predefined range such as [0.0, 1.0]. Data discretization will replace a continuous feature with a discrete one by dividing numeric values into intervals.

1.5 Data modeling

Before discussing the data modeling algorithms, it would be best to explain some terminology. Typically, the data preprocessing step would transform the raw data into a tabular form, in which the columns represent features/variables and rows correspond to samples/instances. For instance, Table 1.1 is a sample data set that has eight samples and five features (class is a special feature for which we are aiming to predict its feature value for a new given sample). The first four features are called *predictive features* and the class feature is the *target feature*. Here the predictive features can be symptoms of some disease, where the value of 1 indicates the existence of a symptom and 0 indicates otherwise. Similarly, the class feature value is 1 if the corresponding person (sample) has the disease.

1.5.1 Pattern mining

Pattern discovery is a core data mining problem, which generates a set of interesting patterns that characterize the data sets in concise and informative forms. Initially, the studies on pattern discovery were dominated by the frequent pattern discovery paradigm, where only frequent patterns were explored. Currently, the issue of frequent pattern discovery has been thoroughly investigated, rendering its limitations well understood. Many alternative pattern discovery formulations are emerging and investigated in the literature. For example, many research efforts impose statistical significance tests over candidate patterns to control the risk of false discoveries.

Table 1.1 An example data set with eight samples and five features

	Feature 1	Feature 2	Feature 3	Feature 4	Class
1	0	0	1	1	0
2	0	1	1	1	0
3	1	0	0	1	0
4	1	0	0	0	0
5	1	1	1	1	1
6	1	1	1	1	1
7	0	0	1	1	1
8	1	1	0	0	1

Here *class* is a special feature representing the category to which each sample belongs.