

Shibiao Wan, Man-Wai Mak

# Machine Learning for Protein Subcellular Localization Prediction

—

DE GRUYTER

### **Authors**

Dr. Shibiao Wan  
The Hong Kong Polytechnic University  
Department of Electronic and Information Engineering  
Hung Hom, Kowloon  
Hong Kong SAR  
shibiao.wan@connect.polyu.hk

Dr. Man-Wai Mak  
The Hong Kong Polytechnic University  
Department of Electronic and Information Engineering  
Hung Hom, Kowloon  
Hong Kong SAR  
enmwamak@polyu.edu.hk

ISBN 978-1-5015-1048-9  
e-ISBN (PDF) 978-1-5015-0150-0  
e-ISBN (EPUB) 978-1-5015-0152-4  
Set-ISBN (print/ebook) 978-1-5015-0151-7

### **Library of Congress Cataloging-in-Publication Data**

A CIP catalog record for this book has been applied for at the Library of Congress.

### **Bibliographic information published by the Deutsche Nationalbibliothek**

The Deutsche Nationalbibliothek lists this publication in the Deutsche Nationalbibliografie; detailed bibliographic data are available on the Internet at <http://www.dnb.de>.

© 2015 Walter de Gruyter Inc., Boston/Berlin  
Typesetting: PTP-Berlin, Protago T<sub>E</sub>X-Production GmbH  
Printing and binding: CPI books GmbH, Leck  
Cover image: THOMAS DEERINCK, NCMIR/Science Photo Library/gettyimages  
© Printed on acid-free paper  
Printed in Germany

[www.degruyter.com](http://www.degruyter.com)

## Preface

Proteins, which are essential macromolecules for organisms, need to be located in appropriate physiological contexts within a cell to exhibit tremendous diversity of biological functions. Aberrant protein subcellular localization may lead to a broad range of diseases. Knowing where a protein resides within a cell can give insights into drug target discovery and drug design. This book explores machine-learning approaches to the automatic prediction of protein subcellular localization. The approaches exploit the gene ontology database to extract relevant information. With the ever increasing numbers of new protein sequences in the postgenomic era, machine-learning approaches have become an indispensable tool for assisting the laborious and time-consuming wet-lab experiments and for accurate, fast, and large-scale predictions in proteomics research.

Recent years have witnessed an incredibly fast development of molecular biology and computer science, which makes it possible to utilize computational methods to determine the subcellular locations of proteins. It is of paramount significance for wet-lab biologists, bioinformaticians, and computational biologists to be informed of the up-to-date development in this field. Compared to traditional books on protein subcellular localization, this book has the following advantages:

1. This book elaborately presents the latest state-of-the-art machine-learning approaches for protein subcellular localization prediction.
2. This book comprehensively covers many aspects of protein subcellular localization, from single- to multi-label prediction, from prediction of *Homo sapiens* proteins, *Viridiplantae* proteins, *Eukaryota* proteins to prediction of *Virus* proteins.
3. This book systematically introduces three machine-learning approaches to improving predictors' performance, including classification refinement, deeper feature extraction and dimensionality reduction.
4. This book not only proposes several advanced and accurate single- and multi-label predictors but also introduces their easy-to-use online web-servers.

This book is organized into four related parts:

1. Part I – Chapters 1, 2, and 3 – introduces the significance of computationally predicting protein subcellular localization, provides an overview of state-of-the-art approaches, and details the legitimacy of using gene ontology (GO) information for predicting subcellular localization of proteins.
2. Part II – Chapters 4, 5, 6, and 7 – proposes several state-of-the-art predictors for single- and multi-location protein subcellular localization. In Chapter 4, two predictors, namely GOASVM and FusionSVM, both based on GO information, are proposed for single-location protein subcellular localization. Subsequently, multi-location protein subcellular localization is described in Chapter 5. In this chapter, several multi-label predictors, including mGOASVM, AD-SVM, and mPLR-Loc,

which were developed based on different classifiers, are introduced for accurate prediction of subcellular localization of both single- and multi-location proteins. Next, Chapter 6 presents the predictors, namely SS-Loc and HybridGO-Loc, which exploit the deep information embedded in the hierarchical structure of the GO Database. These predictors incorporate the information of semantic similarity over GO terms. For large-scale protein subcellular localization, Chapter 7 introduces ensemble random projection to construct two dimension-reduced multi-label predictors, namely RP-SVM and R3P-Loc. In addition, two compact databases (ProSeq and ProSeq-GO) are proposed to replace the conventional databases (Swiss-Prot and GOA) for fast and efficient feature extraction.

3. Part III – Chapters 8, 9, and 10 – presents the experimental setup and results for all of the proposed predictors and further discusses the properties of the proposed predictors. Chapter 8 details the specific experimental setup, including datasets construction and performance metrics. Extensive experimental results and analyses for all the proposed predictors are detailed in Chapter 9. Further discussions are provided in Chapter 10.
4. Part IV – Chapter 11 – gives a conclusion and possible future directions for further research in this field.

It is confidently believed that this book will provide bioinformaticians and computational biologists with the latest state-of-the-art machine-learning approaches for protein subcellular localization prediction and will enlighten them with a systematic scheme to improve predictors' performance. For wet-lab biologists, this book offers accurate and fast subcellular-localization predictors and easy-to-use online web-servers.

**Acknowledgement:** This book is an outgrowth of four years of research on the topics of bioinformatics and machine learning. First, the authors would like to express their sincere gratitude and appreciation to Prof. Sun-Yuan Kung from Princeton University, whose insightful comments and invaluable suggestions have facilitated the research.

The authors are also indebted to Prof. Yue Wang from Virginia Tech (VT), USA, and Dr. Zhen Zhang and Dr. Bai Zhang from Johns Hopkins University (JHU), USA. Our gratitude also goes to all of the CBIL members of VT and collaborators at JHU. Deep thanks should also go to Prof. Hong Yan from City University of Hong Kong, Hong Kong SAR, and Dr. Haiying Wang from the University of Ulster, UK. Their critical and constructive suggestions were imperative for the accomplishment of the book.

We are also grateful to senior editorial director Mr. Alexander Greene, project editor Ms. Julia Lauterbach, and project editor Ms. Lara Wysong of the De Gruyter publisher, who have provided professional assistance throughout the project.

Both authors are with the Center for Multimedia Signal Processing, Department of Electronic and Information Engineering, The Hong Kong Polytechnic University, Hong Kong SAR, China. The authors are grateful to the university and the department for their generous and consistent support.

We are pleased to acknowledge that the work presented in this book was in part supported by The Hong Kong Polytechnic University (Grant No. G-YJ86, G-YL78, and G-YN18) and the Research Grant Council of Hong Kong SAR (Grant No. PolyU5264/09E and PolyU 152117/14E).

The authors would also like to thank many collaborators and colleagues, including Wei Wang, Jian Guo, and others.

Particularly, Shibiao Wan would like to give special thanks to his partner Jieqiong Wang for her unreserved love and support. Last but not the least, the authors wish to give their deepest gratitude to their families. Without their generous support and full understanding, this book would not have been so smoothly completed.

**Dr. Shibiao Wan** is currently a Postdoctoral Fellow of the Department of Electronic and Information Engineering at the Hong Kong Polytechnic University. He obtained his BEng degree in telecommunication engineering from Wuhan University, China in 2010, and his PhD degree in bioinformatics from the Hong Kong Polytechnic University in 2014. He was a visiting scholar in the Virginia Tech and the Johns Hopkins School of Medicine from Spring 2013 to Summer 2013. His current research interests include bioinformatics, computational biology, and machine learning. He has published a number of technical articles on top bioinformatics journals such as *BMC Bioinformatics*, *PLoS ONE*, *Journal of Theoretical Biology*, etc, and key international conferences on signal processing, bioinformatics, and machine learning such as ICASSP, BIBM, MLSP, etc. He serves as a reviewer for a number of journals, such as *IEEE Trans. on Nanobioscience*, *AMC*, *JAM*, *IJBI*, and *IJMLC*.

**Dr. Man-Wai Mak** is an Associate Professor of the Department of Electronic and Information Engineering at the Hong Kong Polytechnic University. He has authored more than 150 technical articles in speaker recognition, machine learning, and bioinformatics. Dr. Mak is also a coauthor of the postgraduate textbook *Biometric Authentication: A Machine Learning Approach* (Prentice Hall, 2005). He served as a member of the IEEE Machine Learning for Signal Processing Technical Committee from 2005–2007. He has been serving as an associate editor of *IEEE/ACM Trans. on Audio, Speech and Language Processing*, *Journal of Signal Processing Systems*, and *Advances in Artificial Neural Systems*. He has been a technical committee member of a number of international conferences, such as Interspeech, ICSLP, and IEEE Workshop on MLSP.



# Contents

Preface — v

List of Abbreviations — xv

## 1 Introduction — 1

- 1.1 Proteins and their subcellular locations — 1
- 1.2 Why computationally predict protein subcellular localization? — 3
  - 1.2.1 Significance of the subcellular localization of proteins — 3
  - 1.2.2 Conventional wet-lab techniques — 3
  - 1.2.3 Computational prediction of protein subcellular localization — 4
- 1.3 Organization of this book — 5

## 2 Overview of subcellular localization prediction — 7

- 2.1 Sequence-based methods — 7
  - 2.1.1 Composition-based methods — 7
  - 2.1.2 Sorting signal-based methods — 14
  - 2.1.3 Homology-based methods — 17
- 2.2 Knowledge-based methods — 21
  - 2.2.1 GO-term extraction — 22
  - 2.2.2 GO-vector construction — 24
- 2.3 Limitations of existing methods — 26
  - 2.3.1 Limitations of sequence-based methods — 26
  - 2.3.2 Limitations of knowledge-based methods — 27

## 3 Legitimacy of using gene ontology information — 28

- 3.1 Direct table lookup? — 29
  - 3.1.1 Table lookup procedure for single-label prediction — 30
  - 3.1.2 Table-lookup procedure for multi-label prediction — 32
  - 3.1.3 Problems of table lookup — 33
- 3.2 Using only cellular component GO terms? — 34
- 3.3 Equivalent to homologous transfer? — 34
- 3.4 More reasons for using GO information — 35

## 4 Single-location protein subcellular localization — 37

- 4.1 Extracting GO from the Gene Ontology Annotation Database — 37
  - 4.1.1 Gene Ontology Annotation Database — 37
  - 4.1.2 Retrieval of GO terms — 40
  - 4.1.3 Construction of GO vectors — 41
  - 4.1.4 Multiclass SVM classification — 43

4.2	FusionSVM: Fusion of gene ontology and homology-based features — 45
4.2.1	InterProGOSVM: Extracting GO from InterProScan — 45
4.2.2	PairProSVM: A homology-based method — 47
4.2.3	Fusion of InterProGOSVM and PairProSVM — 48
4.3	Summary — 49
<b>5</b>	<b>From single- to multi-location — 51</b>
5.1	Significance of multi-location proteins — 51
5.2	Multi-label classification — 51
5.2.1	Algorithm-adaptation methods — 52
5.2.2	Problem transformation methods — 52
5.2.3	Multi-label classification in bioinformatics — 54
5.3	mGOASVM: A predictor for both single- and multi-location proteins — 54
5.3.1	Feature extraction — 55
5.3.2	Multi-label multiclass SVM classification — 56
5.4	AD-SVM: An adaptive decision multi-label predictor — 58
5.4.1	Multi-label SVM scoring — 59
5.4.2	Adaptive decision for SVM (AD-SVM) — 59
5.4.3	Analysis of AD-SVM — 60
5.5	mPLR-Loc: A multi-label predictor based on penalized logistic regression — 63
5.5.1	Single-label penalized logistic regression — 64
5.5.2	Multi-label penalized logistic regression — 65
5.5.3	Adaptive decision for LR (mPLR-Loc) — 65
5.6	Summary — 66
<b>6</b>	<b>Mining deeper on GO for protein subcellular localization — 67</b>
6.1	Related work — 67
6.2	SS-Loc: Using semantic similarity over GO — 69
6.2.1	Semantic similarity measures — 70
6.2.2	SS vector construction — 71
6.3	HybridGO-Loc: Hybridizing GO frequency and semantic similarity features — 72
6.3.1	Hybridization of two GO features — 73
6.3.2	Multi-label multiclass SVM classification — 73
6.4	Summary — 75
<b>7</b>	<b>Ensemble random projection for large-scale predictions — 77</b>
7.1	Random projection — 77



7.2	RP-SVM: A multi-label classifier with ensemble random projection — 79
7.2.1	Ensemble multi-label classifier — 80
7.2.2	Multi-label classification — 80
7.3	R3P-Loc: A compact predictor based on ridge regression and ensemble random projection — 82
7.3.1	Limitation of using current databases — 82
7.3.2	Creating compact databases — 84
7.3.3	Single-label ridge regression — 85
7.3.4	Multi-label ridge regression — 86
7.4	Summary — 88
<b>8</b>	<b>Experimental setup — 89</b>
8.1	Prediction of single-label proteins — 89
8.1.1	Datasets construction — 89
8.1.2	Performance metrics — 93
8.2	Prediction of multi-label proteins — 94
8.2.1	Dataset construction — 94
8.2.2	Datasets analysis — 97
8.2.3	Performance metrics — 100
8.3	Statistical evaluation methods — 103
8.4	Summary — 104
<b>9</b>	<b>Results and analysis — 105</b>
9.1	Performance of GOASVM — 105
9.1.1	Comparing GO vector construction methods — 105
9.1.2	Performance of successive-search strategy — 106
9.1.3	Comparing with methods based on other features — 108
9.1.4	Comparing with state-of-the-art GO methods — 109
9.1.5	GOASVM using old GOA databases — 110
9.2	Performance of FusionSVM — 111
9.2.1	Comparing GO vector construction and normalization methods — 111
9.2.2	Performance of PairProSVM — 112
9.2.3	Performance of FusionSVM — 113
9.2.4	Effect of the fusion weights on the performance of FusionSVM — 114
9.3	Performance of mGOASVM — 115
9.3.1	Kernel selection and optimization — 115
9.3.2	Term-frequency for mGOASVM — 115
9.3.3	Multi-label properties for mGOASVM — 117
9.3.4	Further analysis of mGOASVM — 118
9.3.5	Comparing prediction results of novel proteins — 120
9.4	Performance of AD-SVM — 121

9.5	Performance of mPLR-Loc —	122
9.5.1	Effect of adaptive decisions on mPLR-Loc —	122
9.5.2	Effect of regularization on mPLR-Loc —	124
9.6	Performance of HybridGO-Loc —	125
9.6.1	Comparing different features —	125
9.7	Performance of RP-SVM —	127
9.7.1	Performance of ensemble random projection —	127
9.7.2	Comparison with other dimension-reduction methods —	127
9.7.3	Performance of single random-projection —	130
9.7.4	Effect of dimensions and ensemble size —	130
9.8	Performance of R3P-Loc —	130
9.8.1	Performance on the compact databases —	130
9.8.2	Effect of dimensions and ensemble size —	134
9.8.3	Performance of ensemble random projection —	137
9.9	Comprehensive comparison of proposed predictors —	137
9.9.1	Comparison of benchmark datasets —	137
9.9.2	Comparison of novel datasets —	140
9.10	Summary —	143
<b>10</b>	<b>Properties of the proposed predictors —</b>	<b>145</b>
10.1	Noise data in the GOA Database —	145
10.2	Analysis of single-label predictors —	146
10.2.1	GOASVM vs FusionSVM —	146
10.2.2	Can GOASVM be combined with PairProSVM? —	147
10.3	Advantages of mGOASVM —	148
10.3.1	GO-vector construction —	148
10.3.2	GO subspace selection —	148
10.3.3	Capability of handling multi-label problems —	149
10.4	Analysis for HybridGO-Loc —	149
10.4.1	Semantic similarity measures —	149
10.4.2	GO-frequency features vs SS features —	150
10.4.3	Bias analysis —	150
10.5	Analysis for RP-SVM —	151
10.5.1	Legitimacy of using RP —	151
10.5.2	Ensemble random projection for robust performance —	152
10.6	Comparing the proposed multi-label predictors —	153
10.7	Summary —	154
<b>11</b>	<b>Conclusions and future directions —</b>	<b>157</b>
11.1	Conclusions —	157
11.2	Future directions —	158

**A Webservers for protein subcellular localization — 161**

- A.1 GOASVM webserver — 161
- A.2 mGOASVM webserver — 162
- A.3 HybridGO-Loc webserver — 163
- A.4 mPLR-Loc webserver — 164

**B Support vector machines — 167**

- B.1 Binary SVM classification — 167
- B.2 One-vs-rest SVM classification — 170

**C Proof of no bias in LOOCV — 173**

**D Derivatives for penalized logistic regression — 174**

**Bibliography — 177**

**Index — 191**



## List of Abbreviations

AA	amino-acid compositions
AC	accession number
ACC	overall accuracy
ACR	acrosome
AD	adaptive decision
BLAST	basic local alignment search tool
BR	binary relevance
c-region	C-terminal flanking region
CEL	cell wall
CEN	centrosome
CHL	chloroplast
CM	cell membrane
cTP	chloroplast transit peptide
CYA	cyanelle
CYK	cytoskeleton
CYT	cytoplasm
DLS	distinct label set
EBI	European Bioinformatics Institute
ECC	ensembles of classifier chains
END	endosome
ER	endoplasmic reticulum
EU16	the 16-class eukaryotic dataset
EXP	inferred from experiment
EXT	extracellular
F1	F1-score
GapAA	gapped amino-acid pair compositions
GO	gene ontology
GOA	gene ontology annotation database
GOL	Golgi apparatus
h-region	central hydrophobic region
HCYT	host cytoplasm
HER	host endoplasmic reticulum
HL	Hamming loss
HMMs	hidden Markov models
HNUC	host nucleus
HUM12	the 12-class human dataset
HYD	hydrogenosome
IDA	inferred from direct assay
IEA	inferred from electronic annotation

IMP	inferred from mutant phenotype
IPI	inferred from physical interaction
ISF	inverse sequence-frequency
ISS	inferred from structural and sequence similarity
LC	label cardinality
LCA	lowest common ancestors
LD	label density
LOOCV	leave-one-out cross validation
LP	label powerset
LR	logistic regression
LYS	lysosome
MCC	Mathew's correlation coefficient
MEL	melanosome
MIC	microsome
MIT	mitochondrion
mTP	mitochondrial targeting peptide
n-region	N-terminal flanking region
NE16	the 16-class novel eukaryotic dataset
NNs	neural networks
NUC	nucleus
OAA	overall actual accuracy
OE11	the 11-class old eukaryotic dataset
OET-KNN	optimized evidence-theoretic K-nearest neighbors
OLA	overall locative accuracy
OLS	ordinary least squares
OMCC	overall Mathew's correlation coefficient
PairAA	amino-acid pair compositions
PDLS	proportion of distinct label set
PER	peroxisome
PLA	plastid
PM	plasma membrane
PseAA	pseudo amino-acid compositions
PSFM	position-specific frequency matrix
PSI-BLAST	position-specific iterative BLAST
PSSM	position-specific scoring matrix
RP	random projection
RR	ridge regression
RS	relevance similarity
SEC	secreted
SP	signal peptide
SPI	spindle pole body
SS	semantic similarity

SVMs	support vector machines
SYN	synapse
TF	term-frequency
TF-ISF	term-frequency–inverse sequence-frequency
TLN	total locative number
VAC	vacuole
VC	viral capsid
WAMCC	weighted average Mathew’s correlation coefficient