

M.N. Murty
Rashmi Raghava

Support Vector Machines and Perceptrons

Learning, Optimization,
Classification, and
Application to Social
Networks



Springer

SpringerBriefs in Computer Science

Series editors

- Stan Zdonik, Brown University, Providence, Rhode Island, USA
Shashi Shekhar, University of Minnesota, Minneapolis, Minnesota, USA
Jonathan Katz, University of Maryland, College Park, Maryland, USA
Xindong Wu, University of Vermont, Burlington, Vermont, USA
Lakhmi C. Jain, University of South Australia, Adelaide, South Australia, Australia
David Padua, University of Illinois Urbana-Champaign, Urbana, Illinois, USA
Xuemin (Sherman) Shen, University of Waterloo, Waterloo, Ontario, Canada
Borko Furht, Florida Atlantic University, Boca Raton, Florida, USA
V.S. Subrahmanian, University of Maryland, College Park, Maryland, USA
Martial Hebert, Carnegie Mellon University, Pittsburgh, Pennsylvania, USA
Katsushi Ikeuchi, University of Tokyo, Tokyo, Japan
Bruno Siciliano, Università di Napoli Federico II, Napoli, Italy
Sushil Jajodia, George Mason University, Fairfax, Virginia, USA
Newton Lee, Newton Lee Laboratories, LLC, Tujunga, California, USA

More information about this series at <http://www.springer.com/series/10028>

M.N. Murty · Rashmi Raghava

Support Vector Machines and Perceptrons

Learning, Optimization, Classification,
and Application to Social Networks



Springer

M.N. Murty
Department of Computer Science
and Automation
Indian Institute of Science
Bangalore, Karnataka
India

Rashmi Raghava
IBM India Private Limited
Bangalore, Karnataka
India

ISSN 2191-5768
SpringerBriefs in Computer Science
ISBN 978-3-319-41062-3
DOI 10.1007/978-3-319-41063-0

ISSN 2191-5776 (electronic)
ISBN 978-3-319-41063-0 (eBook)

Library of Congress Control Number: 2016943387

© The Author(s) 2016

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made.

Printed on acid-free paper

This Springer imprint is published by Springer Nature
The registered company is Springer International Publishing AG Switzerland

Preface

Overview

Support Vector Machines (SVMs) have been widely used in *Classification, Clustering and Regression*. In this book, we deal primarily with classification. Classifiers can be either *linear* or *nonlinear*. The linear classifiers typically are learnt based on a *linear discriminant function* that separates the feature space into two half-spaces, where one half-space corresponds to one of the two classes and the other half-space corresponds to the remaining class. So, these half-space classifiers are ideally suited to solve *binary classification* or two-class classification problems. There are a variety of schemes to build multiclass classifiers based on combinations of several binary classifiers.

Linear discriminant functions are characterized by a *weight vector* and a *threshold* weight that is a scalar. These two are learnt from the training data. Once these entities are obtained we can use them to classify patterns into any one of the two classes. It is possible to extend the notion of linear discriminant functions (LDFs) to deal with even nonlinearly separable data with the help of a suitable mapping of the data points from the low-dimensional *input* space to a possibly higher dimensional *feature space*.

Perceptron is an early classifier that successfully dealt with linearly separable classes. Perceptron could be viewed as the simplest form of *artificial neural network*. An excellent theory to characterize parallel and distributed computing was put forth by *Minsky and Papert* in the form of a book on perceptrons. They use logic, geometry, and group theory to provide a computational framework for perceptrons. This can be used to show that any computable function can be characterized as a linear discriminant function possibly in a high-dimensional space based on *minterms* corresponding to the input Boolean variables. However, for some types of problems one needs to use all the minterms which correspond to using an exponential number of minterms that could be realized from the primitive variables.

SVMs have revolutionized the research in the areas of *machine learning* and *pattern recognition*, specifically classification, so much that for a period of more

than two decades they are used as state-of-the-art classifiers. Two distinct properties of SVMs are:

1. The problem of learning the LDF corresponding to SVM is posed as a convex optimization problem. This is based on the intuition that the hyperplane separating the two classes is learnt so that it corresponds to maximizing the *margin* or some kind of separation between the two classes. So, they are also called as *maximum-margin classifiers*.
2. Another important notion associated with SVMs is the *kernel trick* which permits us to perform all the computations in the low-dimensional input space rather than in a higher dimensional feature space.

These two ideas become so popular that the first one lead to the increase of interest in the area of *convex optimization*, whereas the second idea was exploited to deal with a variety of other classifiers and clustering algorithms using an appropriate kernel/similarity function.

The current popularity of SVMs can be attributed to excellent and popular software packages like *LIBSVM*. Even though SVMs can be used in nonlinear classification scenarios based on the kernel trick, the linear SVMs are more popular in the real-world applications that are high-dimensional. Further learning the parameters could be time-consuming. There is a renewal of energy, in the recent times, to examine other linear classifiers like perceptrons. Keeping this in mind, we have dealt with both perceptron and SVM classifiers in this book.

Audience

This book is intended for senior undergraduate and graduate students and researchers working in *machine learning, data mining, and pattern recognition*. Even though SVMs and perceptrons are popular, people find it difficult to understand the underlying theory. We present material in this book so that it is accessible to a wide variety of readers with some basic exposure to undergraduate level mathematics. The presentation is intentionally made simpler to make the reader feel comfortable.

Organization

This book is organized as follows:

1. *Literature and Background*: Chapter 1 presents literature and state-of-the-art techniques in SVM-based classification. Further, we also discuss relevant background required for pattern classification. We define some of the important terms that are used in the rest of the book. Some of the concepts are explained with the help of easy to understand examples.

2. *Linear Discriminant Function:* In Chap. 2 we introduce the notion of *linear discriminant function* that forms the basis for the linear classifiers described in the text. The role of *weight vector* W and the *threshold* b are explained in describing linear classifiers. We also describe other linear classifiers including the *minimal distance classifier* and the *Naïve Bayes* classifier. It also explains how nonlinear discriminant functions could be viewed as linear discriminant functions in higher dimensional spaces.
3. *Perceptron:* In Chap. 3 we describe perceptron and how it can be used for classification. We deal with *perceptron learning algorithm* and explain how it can be used to learn Boolean functions. We provide a simple proof to show how the algorithm converges. We explain the notion of *order of a perceptron* that has bearing on the computational complexity. We illustrate it on two different classification datasets.
4. *Linear SVM:* In this Chap. 4, we start with the similarity between SVM and perceptron as both of them are used for *linear classification*. We discuss the difference between them in terms of the form of computation of w , the optimization problem underlying each, and the *kernel trick*. We introduce the linear SVM which possibly is the most popular classifier in machine learning. We introduce the notion of *maximum margin* and the geometric and semantic interpretation of the same. We explain how a binary classifier could be used in building a multiclass classifier. We provide experimental results on two datasets.
5. *Kernel Based SVM:* In Chap. 5, we discuss the notion of *kernel* or similarity function. We discuss how the optimization problem changes when the classes are not linearly separable or when there are some data points on the margin. We explain in simple terms the *kernel trick* and explain how it is used in classification. We illustrate using two practical datasets.
6. *Application to Social Networks:* In Chap. 6 we consider *social networks*. Specifically, issues related to representation of social networks using *graphs*; these graphs are in turn represented as matrices or lists. We consider the problem of community detection in social networks and *link prediction*. We examine several existing schemes for link prediction including the one based on SVM classifier. We illustrate its working based on some network datasets.
7. *Conclusion:* We conclude in Chap. 7 and also present potential future directions.

Bangalore, India

M.N. Murty
Rashmi Raghava

Contents

1	Introduction	1
1.1	Terminology	1
1.1.1	What Is a Pattern?	1
1.1.2	Why Pattern Representation?	2
1.1.3	What Is Pattern Representation?	2
1.1.4	How to Represent Patterns?	2
1.1.5	Why Represent Patterns as Vectors?	2
1.1.6	Notation	3
1.2	Proximity Function	3
1.2.1	Distance Function	3
1.2.2	Similarity Function	4
1.2.3	Relation Between Dot Product and Cosine Similarity	5
1.3	Classification	6
1.3.1	Class	6
1.3.2	Representation of a Class	6
1.3.3	Choice of $G(X)$	7
1.4	Classifiers	7
1.4.1	Nearest Neighbor Classifier (NNC)	7
1.4.2	K-Nearest Neighbor Classifier (KNNC)	7
1.4.3	Minimum-Distance Classifier (MDC)	8
1.4.4	Minimum Mahalanobis Distance Classifier	9
1.4.5	Decision Tree Classifier: (DTC)	10
1.4.6	Classification Based on a Linear Discriminant Function	12
1.4.7	Nonlinear Discriminant Function	12
1.4.8	Naïve Bayes Classifier: (NBC)	13
1.5	Summary	14
	References	14

2 Linear Discriminant Function	15
2.1 Introduction	15
2.1.1 Associated Terms	15
2.2 Linear Classifier	17
2.3 Linear Discriminant Function	19
2.3.1 Decision Boundary	19
2.3.2 Negative Half Space	19
2.3.3 Positive Half Space	19
2.3.4 Linear Separability	20
2.3.5 Linear Classification Based on a Linear Discriminant Function	20
2.4 Example Linear Classifiers	23
2.4.1 Minimum-Distance Classifier (MDC)	23
2.4.2 Naïve Bayes Classifier (NBC)	23
2.4.3 Nonlinear Discriminant Function	24
References	25
3 Perceptron	27
3.1 Introduction	27
3.2 Perceptron Learning Algorithm	28
3.2.1 Learning Boolean Functions	28
3.2.2 W Is Not Unique	30
3.2.3 Why Should the Learning Algorithm Work?	30
3.2.4 Convergence of the Algorithm	31
3.3 Perceptron Optimization	32
3.3.1 Incremental Rule	33
3.3.2 Nonlinearly Separable Case	33
3.4 Classification Based on Perceptrons	34
3.4.1 Order of the Perceptron	35
3.4.2 Permutation Invariance	37
3.4.3 Incremental Computation	37
3.5 Experimental Results	38
3.6 Summary	39
References	40
4 Linear Support Vector Machines	41
4.1 Introduction	41
4.1.1 Similarity with Perceptron	41
4.1.2 Differences Between Perceptron and SVM	42
4.1.3 Important Properties of SVM	42
4.2 Linear SVM	43
4.2.1 Linear Separability	43
4.2.2 Margin	44
4.2.3 Maximum Margin	46
4.2.4 An Example	47

4.3	Dual Problem	49
4.3.1	An Example	50
4.4	Multiclass Problems.	51
4.5	Experimental Results	52
4.5.1	Results on Multiclass Classification	52
4.6	Summary	54
	References	56
5	Kernel-Based SVM	57
5.1	Introduction	57
5.1.1	What Happens if the Data Is Not Linearly Separable?	57
5.1.2	Error in Classification	58
5.2	Soft Margin Formulation	59
5.2.1	The Solution.	59
5.2.2	Computing b	60
5.2.3	Difference Between the Soft and Hard Margin Formulations	60
5.3	Similarity Between SVM and Perceptron	60
5.4	Nonlinear Decision Boundary	62
5.4.1	Why Transformed Space?	63
5.4.2	Kernel Trick	63
5.4.3	An Example	64
5.4.4	Example Kernel Functions	64
5.5	Success of SVM	64
5.6	Experimental Results	65
5.6.1	Iris Versicolour and Iris Virginica	65
5.6.2	Handwritten Digit Classification	66
5.6.3	Multiclass Classification with Varying Values of the Parameter C	66
5.7	Summary	67
	References	67
6	Application to Social Networks.	69
6.1	Introduction	69
6.1.1	What Is a Network?	69
6.1.2	How Do We Represent It?	69
6.2	What Is a Social Network?.	72
6.2.1	Citation Networks	73
6.2.2	Coauthor Networks	73
6.2.3	Customer Networks.	73
6.2.4	Homogeneous and Heterogeneous Networks.	73
6.3	Important Properties of Social Networks.	74
6.4	Characterization of Communities	75
6.4.1	What Is a Community?	75
6.4.2	Clustering Coefficient of a Subgraph	76