LEARNING TO CLASSIFY TEXT USING SUPPORT VECTOR MACHINES

THE KLUWER INTERNATIONAL SERIES IN ENGINEERING AND COMPUTER SCIENCE

LEARNING TO CLASSIFY TEXT USING SUPPORT VECTOR MACHINES

by

Thorsten Joachims

Cornell University, U.S.A.

Dissertation, Universität Dortmund Fachbereich Informatik February 2001



SPRINGER SCIENCE+BUSINESS MEDIA, LLC

Library of Congress Cataloging-in-Publication Data

Joachims, Thorsten.

Learning to classify text using support vector machines: methods, theory, and algorithms / by Thorsten Joachims.

p. cm. – (Kluwer international series in engineering and computer science; SECS 668) Originally presented as the author's thesis (Universität Dortmund) under the title: "The maximum-margin approach to learning text classifiers—methods, theory, and algorithms." Includes bibliographical references and index.

ISBN 978-1-4613-5298-3 ISBN 978-1-4615-0907-3 (eBook) DOI 10.1007/978-1-4615-0907-3

1. Text processing (Computer science) 2. Machine learning. I. Title. II. Series.

QA76.9.T48 J63 2002 005—dc21

2002022127

Copyright © 2002 by Springer Science+Business Media New York Originally published by Kluwer Academic Publishers in 2002

Chighnany published by Kluwer Academic Fublishers in 20

Softcover reprint of the hardcover 1st edition 2002

All rights reserved. No part of this work may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, microfilming, recording, or otherwise, without the written permission from the Publisher, with the exception of any material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work.

Contents

Fo	rewo Pro	ord of. Tom	Mitchell and Prof. Katharina Morik	xi
Pr	eface	2	5	xiii
Ac	kno	wledgm	ients	xv
No	otatic	on		xvii
1.	INT	INTRODUCTION		
	1	Chall	lenges	2
	2	Goal	S	3
	3	Over	view and Structure of the Argument	4
		3.1	Theory	4
		3.2	Methods	5
		3.3	Algorithms	6
	4	Sum	mary	6
2.	TEXT CLASSIFICATION			7
	1	Lean	ning Task	7
		1.1	Binary Setting	8
		1.2	Multi-Class Setting	9
		1.3	Multi-Label Setting	10
	2	Repr	esenting Text	12
		2.1	Word Level	13
		2.2	Sub-Word Level	15
		2.3	Multi-Word Level	15
		2.4	Semantic Level	16
	3	Featu	are Selection	16
		3.1	Feature Subset Selection	17

vi LEARNING TO CLASSIFY TEXT USING SUPPORT VECTOR MACHINES

		3.2	Feature Construction	19
	4	Term	Weighting	20
	5	Conve	entional Learning Methods	22
		5.1	Naive Bayes Classifier	22
		5.2	Rocchio Algorithm	24
		5.3	k-Nearest Neighbors	25
		5.4	Decision Tree Classifier	25
		5.5	Other Methods	26
	6	Perfor	mance Measures	27
		6.1	Error Rate and Asymmetric Cost	28
		6.2	Precision and Recall	29
		6.3	Precision/Recall Breakeven Point and F_{β} -Measure	30
		6.4	Micro- and Macro-Averaging	30
	7	Exper	imental Setup	31
		7.1	Test Collections	31
		7.2	Design Choices	32
3.	SUF	PORT	VECTOR MACHINES	35
	1	Linear Hard-Margin SVMs		
	2	Soft-Margin SVMs		39
	3	Non-Linear SVMs		41
	4	Asymmetric Misclassification Cost		43
	5	Other Maximum-Margin Methods		
	6	Further Work and Further Information		

Part Theory

4.	A STATISTICAL LEARNING MODEL OF TEXT CLASSIFICATION			
	FOR	SVMS		45
	1	Propert	ties of Text-Classification Tasks	46
		1.1	High-Dimensional Feature Space	46
		1.2	Sparse Document Vectors	47
		1.3	Heterogeneous Use of Terms	47
		1.4	High Level of Redundancy	48
		1.5	Frequency Distribution of Words and Zipf's Law	49
	2	A Disc	riminative Model of Text Classification	51
		2.1	Step 1: Bounding the Expected Error Based on the Margin	51

Contents

		2.2	Step 2: Homogeneous TCat-Concepts as a Model of Text-Classification Tasks	53	
		2.3	Step 3: Learnability of TCat-Concepts	59	
	3	Compa	ring the Theoretical Model with Experimental Results	64	
	4	Sensiti	vity Analysis: Difficult and Easy Learning Tasks	66	
		4.1	Influence of Occurrence Frequency	66	
		4.2	Discriminative Power of Term Sets	68	
		4.3	Level of Redundancy	68	
	5	Noisy 7	TCat-Concepts	69	
	6	Limitat	tions of the Model and Open Questions	72	
	7	Related	1 Work	72	
	8	Summa	ary and Conclusions	74	
5.	EFFICIENT PERFORMANCE ESTIMATORS FOR SVMS				
	1	Generi	c Performance Estimators	76	
		1.1	Training Error	76	
		1.2	Hold-Out Testing	77	
		1.3	Bootstrap and Jackknife	78	
		1.4	Cross-Validation and Leave-One-Out	79	
	2	$\xi \alpha$ -Esti	imators	81	
		2.1	Error Rate	82	
		2.2	Recall, Precision, and F_1	89	
	3	Fast Le	ave-One-Out Estimation	93	
	4	Experin	ments	94	
		4.1	How Large are Bias and Variance of the $\xi \alpha$ -Estimators?	95	
		4.2	What is the Influence of the Training Set Size?	99	
		4.3	How Large is the Efficiency Improvement for Exact		
			Leave-One-Out?	101	
	5	Summa	ary and Conclusions	102	

Part Methods

6.	INDUCTIVE TEXT CLASSIFICATION			103
	1	1 Learning Task		
	2	Auto	matic Model and Parameter Selection	105
		2.1	Leave-One-Out Estimator of the PRBEP	106
		2.2	$\xi \alpha$ -Estimator of the PRBEP	106
		2.3	Model-Selection Algorithm	108

vii

	3	Experi	ments	108
		3.1	Word Weighting, Stemming and Stopword Removal	108
		3.2	Trading Off Training Error vs. Complexity	111
		3.3	Non-Linear Classification Rules	113
		3.4	Comparison with Conventional Methods	113
	4	Relate	d Work	116
	5	Summ	ary and Conclusions	117
7.	TRA	NSDU	CTIVE TEXT CLASSIFICATION	119
	1	Learni	ng Task	120
	2	Transc	luctive Support Vector Machines	121
	3	What]	Makes TSVMs well Suited for Text Classification?	123
		3.1	An Intuitive Example	123
		3.2	Transductive Learning of TCat-Concepts	125
	4	Experi	iments	127
	5	Constr	aints on the Transductive Hyperplane	130
	6	Relatio	on to Other Approaches Using Unlabeled Data	133
		6.1	Probabilistic Approaches using EM	133
		6.2	Co-Training	134
		6.3	Other Work on Transduction	139
	7	Summ	ary and Conclusions	139

Part Algorithms

8.	TRA	INING	INDUCTIVE SUPPORT VECTOR MACHINES	141
	1	Problem	m and Approach	142
	2	Genera	l Decomposition Algorithm	143
	3	Selecti	ng a Good Working Set	145
		3.1	Convergence	145
		3.2	How to Compute the Working Set	146
	4	Shrink	ing: Reducing the Number of Variables	146
	5	Efficie	nt Implementation	148
		5.1	Termination Criteria	148
		5.2	Computing the Gradient and the Termination Criteria Efficiently	149
		5.3	What are the Computational Resources Needed in each Iteration?	150
		5.4	Caching Kernel Evaluations	151

Contents

		5.5	How to Solve the QP on the Working Set	152
	6	Relate	ed Work	152
	7	Exper	riments	154
		7.1 7 2	Training Times for Reuters, WebKB, and Ohsumed How does Training Time Scale with the Number of	154
		1.2	Training Examples?	154
		7.3	What is the Influence of the Working-Set-Selection Strategy?	160
		7.4	What is the Influence of Caching?	161
		7.5	What is the Influence of Shrinking?	161
	8	Summ	nary and Conclusions	162
9.	TRA	INING	G TRANSDUCTIVE SUPPORT VECTOR MACHINES	163
	1	Proble	em and Approach	163
	2	The T	SVM Algorithm	165
	3	Analy	sis of the Algorithm	166
		3.1	How does the Algorithm work?	166
		3.2	Convergence	168
	4	Exper	iments	169
		4.1	Does the Algorithm Effectively Maximize Margin?	169
		4.2	Training Times for Reuters, WebKB, and Ohsumed	170
		4.3	How does Training Time Scale with the Number of Training Examples?	170
		4.4	How does Training Time Scale with the Number of Test Examples?	172
	5	Relate	ed Work	172
	6	Summ	ary and Conclusions	174
10	. CON	ICLUS	IONS	175
	1	Open	Question	177
Bi	bliogr	aphy		180
Ar	pend	ices		197
1	SVM	l-Light	Commands and Options	197
Inc	lex			
				203

ix

Foreword

"A good theory yields the best practice" – this saying characterizes the work in this book by Thorsten Joachims. Practitioners need to understand whether and how well their problem can be solved, and which approaches are likely to be most effective. This book reports on machine learning research that has produced both strong experimental results and a theory of text classification that can inform the practitioner interested in applying similar learning methods.

The problem addressed in this book is that of text classification: automatically classifying text documents based on their content. The appearance of billions of online documents in the Internet has created a need for automated methods for classification, where rapid changes and growth of the collection of documents prohibit the use of manual techniques. Automatic text classification is useful for many services such as search, filtering, and routing relevant email to the appropriate addressees. The problem is characterized by very high dimensional data - every word in the document is treated as an attribute - and by little training data - there are typically fewer training documents than there are attributes! This book describes a model of automatic text classification that characterizes the difficulty of a given instance of the problem. One models an instance of the problem by counting the words that occur in the documents each of the classes under consideration. Learnability results refer to this model, and predict how hard it will be to train a classifier for this instance of the problem. Although the proofs of the learnability of TCat concepts are subtle, their application to real data sets is straightforward, as illustrated in this book.

While many learning algorithms have been studied for text, one of the most effective is the Support Vector Machine (SVM) which is the focus of this book. Support vector machines (SVM's) were suggested some time ago by Vladimir Vapnik, but Joachims' work is among the first to seriously explore their use for classifying text. In fact, Joachims created an implementation of SVM's, called *SVM*^{light}, which has been used by many researchers worldwide, because it is an efficient and easy to use implementation of SVMs that is well suited

to text classification. In this book, Joachims explains and explores the use of SVM's for text classification, reporting experimental results applying SVM's to a variety of real-world text classification tasks. He goes on to explore the use of SVM's for transduction, in which the unlabeled examples that are to be classified are used as part of the procedure for learning a good decision boundary, and demonstrates experimentally that this transduction process significantly improves classification accuracy on several text classification problems.

In addition to developing and experimenting with practical algorithms for text classification, the book goes on to develop components of a general theory of text classification. What is the question one would like such a theory to answer? Perhaps the most important question is "can we predict the accuracy that will be achieved by our trained classifier, as a function of the number of training examples it is provided?" Joachims provides in this book the first steps toward a general theory to answer this question for the problem of text classification. This theory builds on the known statistical distributions of words as they occur in natural language, and posits several other properties of text and text classification problems. Within this set of assumptions, Joachims succeeds in relating accuracy to the number of training examples. In fact, the theory is framed in terms of parameters of the text corpus, such as parameters of the word distribution that can be efficiently measured for any text corpus, then used to characterize the properties of the corpus that influence this relation between accuracy and number of training examples. As with most theories, this one formulates the theoretical problem as a somewhat simplified version of the actual problem. Nevertheless, Joachims demonstrates that the theory successfully predicts the relative difficulty of text learning over three different real-world data sets and classification problems.

In short, the work in this book represents an important step in our understanding of text learning. It describes a state-of-the-art machine learning algorithm for practical text classification, with a freeware implementation accessible over the web. In addition, it provides the first theoretical characterization of text classification learning problems, relating the expected error of a classifier to measurable properties of the text data, plus the number of training examples provided to the learner. This research, combining solid experimental work with highly relevant theory, represents an important contribution to our understanding of text learning, and represents a model for future work combining theory and practice. We believe that students, lecturers, computational theorists, and practitioners will enjoy reading the book as much as we enjoyed accompanying its formation.

Prof. Tom Mitchell Carnegie Mellon University Prof. Katharina Morik Universität Dortmund

Preface

Text classification, or the task of automatically assigning semantic categories to natural language text, has become one of the key methods for organizing online information. Since hand-coding such classification rules is costly or even impractical, most modern approaches employ machine learning techniques to automatically learn text classifiers from examples. However, none of these conventional approaches combines good prediction performance, theoretical understanding, and efficient training algorithms.

Based on ideas from Support Vector Machines (SVMs), this book presents a new approach to learning text classifiers from examples. It provides not only learning methods that empirically have state-of-the-art predictive performance, but also a theory that connects the properties of text-classification tasks with the generalization accuracy of the learning methods, as well as algorithms that efficiently implement the methods.

In particular, the results show that the SVM approach to learning text classifiers is highly effective without greedy heuristic components. To explain these empirical findings, this book analyzes the statistical properties of text-classification tasks and presents a theoretical learning model that leads to bounds on the expected error rate of an SVM. The bounds are based on improved results about leave-one-out estimators for SVMs. These results also lead to a new group of performance estimators for SVMs, called $\xi\alpha$ -estimators, and to an improved algorithm for computing the leave-one-out error of an SVM. While all results mentioned so far were for the inductive setting, this book also introduces the idea of transduction to text classification. It shows how and why exploiting the location of the test points during learning can improve predictive performance. To make the SVM approach to learning text classifiers applicable in practice, this book also derives new algorithms for training SVMs. For both the inductive and the transductive setting, these algorithms substantially extend the scalability of SVMs to large-scale problems.

xiv LEARNING TO CLASSIFY TEXT USING SUPPORT VECTOR MACHINES

While the work presented in this book is driven by the application, the techniques it develops are not limited to text classification, but can be applied in other contexts as well. In addition, not only can individual techniques be transferred to other tasks, this work as a whole can be useful as a case study for how to approach high-dimensional learning tasks.

This book is based on my dissertation with the title "The Maximum-Margin Approach to Learning Text Classifiers — Methods, Theory, and Algorithms", defended in February 2001 at the Fachbereich Informatik, Universität Dortmund, Germany. From the beginning, it was not written purely for the thesis committee, but with a broader audience in mind. So, the book is designed to be self-contained and I believe it can be useful for all readers interested in text classification — both for research, as well as for product development.

After an introduction to Support Vector Machines and an overview of the state-of-the-art in text classification, this book is divided into three parts: theory, methods, and algorithms. Each part is self-contained, facilitating selective reading. Furthermore, all methods and algorithms proposed in this book are implemented in the software SVM^{light} , making it easy to replicate and extend the results presented in the following.

I hope you will find this book useful and enjoy reading it.

THORSTEN JOACHIMS

Acknowledgments

Without each of the following people, my dissertation and this book would have turned out differently. Or maybe, it would never have been started nor finished.

First, I would like to thank the two people who had the largest immediate influence on this work — Prof. Katharina Morik and Prof. Tom Mitchell. Prof. Morik was the person who got me excited about machine learning. It was her enthusiasm that made and still makes me want to understand the nature of learning. In particular, I thank her for teaching me about good and important research and for the guidance that shaped my dissertation without constraining it. Prof. Mitchell was the person who introduced me to text learning during my time at Carnegie Mellon University. So the particular topic of my dissertation is due to him. I am very grateful for his advice and for the motivation he managed to get across despite the distance. I also thank Prof. Fuhr for his comments and his suggestions about this work.

Particular thanks go to Dr. Phoebe Sengers. She kept me happy and appropriately focused throughout this work. Despite her aversion towards Greek letters, she read this whole document twice, found many logical inconsistencies, and corrected my English.

Many of my colleagues deserve my thanks for helpful discussions, in particular Ralf Klinkenberg, Stefan Rüping, and Peter Brockhausen from the AI Unit. With the DFG Collaborative Research Center on Complexity Reduction in Multivariate Data (SFB475) supporting this work, I got much help from Thomas Fender, Ursula Sondhauß, Prof. Gather, and Prof. Weihs from the statistics department. While more difficult due to the distance, discussions and collaboration with Dr. Tobias Scheffer, Prof. Sebastian Thrun, Dr. Andrew McCallum, Dr. Dunja Mladenić, Marko Grobelnik, Dr. Justin Boyan, Dr. John Platt, Dr. Susan Dumais, Dr. Dayne Freitag, Dr. Carlotta Domeniconi, Dr. Maria Wolters, Dr. Mehran Sahami, Prof. Mark Craven, Prof. Lyle Ungar, Dr. Alex Smola, John Langford, Sean Slattery, and Rosie Jones substantially influenced this work. Of particular inspiration were the discussions with Prof. Vladimir Vapnik during his visit in Dortmund.

Last but not least, I would like to thank my parents Bernhard Joachims and Hiltrud Joachims, as well as my sister Martina Joachims. They supported me throughout my whole life in every respect and made me interested in science and learning about the world — sorry, Mom, it turned out not to be medicine.

Notation

$\vec{x_i}$	input patterns
y_i	target values (classes)
X	feature space
\boldsymbol{n}	number of training examples
k	number of test examples
Ν	dimensionality of the input space
L	learner
\mathcal{H}	hypothesis space
h	hypothesis from the hypothesis space
$h_{\mathcal{L}}$	hypothesis that the learner \mathcal{L} returns
R(h)	(expected) risk of the hypothesis h
$R_{emp}(h)$	empirical risk of the hypothesis h on a training sample
L	loss function
$L_{0/1}$	0/1-loss function
-,-	,
$ec{w}$	weight vector of a hyperplane $\langle \vec{w}, b \rangle$
Ь	constant offset (or threshold) of a hyperplane $\langle \vec{w}, b \rangle$
δ	margin of a hyperplane
R	diameter of a ball containing the data, usually approximated by $\max \vec{x} _2$
$lpha_i$	Lagrange multiplier
ã	vector of all Lagrange multipliers
ξi	slack variables
$(ec{x}_1\cdotec{x}_2)$	dot product between vectors $\vec{x_1}$ and $\vec{x_2}$
κ	Mercer kernel
Q	Hessian of the quadratic program
Err	error rate
Rec	recall
Prec	precision
F_{ρ}	F_{a} -measure
PRBEP	precision/Recall breakeven point
PRAVG	arithmetic average of precision and recall
$ec{x}^T$	transpose of the vector \vec{x}
R	the set of real numbers
N	the set of natural numbers
X	cardinality of set X
abs(a)	absolute value of a
. 1	L_1 -norm, $\ \vec{x}\ _1 := \sum abs(x_i)$
. 2 or .	L_2 -norm (Euclidian distance), $\ \vec{x}\ := \sqrt{(\vec{x} \cdot \vec{x})}$
exp(a)	2.7182818 ^a
ln	logarithm to base 2.7182818