

Adversarial Machine Learning

Written by leading researchers, this complete introduction brings together all the theory and tools needed for building robust machine learning in adversarial environments. Discover how machine learning systems can adapt when an adversary actively poisons data to manipulate statistical inference, learn the latest practical techniques for investigating system security and performing robust data analysis, and gain insight into new approaches for designing effective countermeasures against the latest wave of cyberattacks. Privacy-preserving mechanisms and near-optimal evasion of classifiers are discussed in detail, and in-depth case studies on email spam and network security highlight successful attacks on traditional machine learning algorithms. Providing a thorough overview of the current state of the art in the field and possible future directions, this groundbreaking work is essential reading for researchers, practitioners, and students in computer security and machine learning and for those wanting to learn about the next stage of the cybersecurity arms race.

Anthony D. Joseph is a Chancellor's Professor in the Department of Electrical Engineering and Computer Sciences at the University of California, Berkeley. He was formerly the Director of Intel Labs Berkeley.

Blaine Nelson is a Software Engineer in the Counter-Abuse Technology (CAT) team at Google. He previously worked at the University of Potsdam and the University of Tübingen.

Benjamin I. P. Rubinstein is an Associate Professor in Computing and Information Systems at the University of Melbourne. He has previously worked at Microsoft Research, Google Research, Yahoo! Research, Intel Labs Berkeley, and IBM Research.

J. D. Tygar is a Professor at the University of California, Berkeley, and he has worked widely in the field of computer security. At Berkeley, he holds appointments in both the Department of Electrical Engineering and Computer Sciences and the School of Information.

Cambridge University Press

978-1-107-04346-6 — Adversarial Machine Learning

Anthony D. Joseph , Blaine Nelson , Benjamin I. P. Rubinstein , J. D. Tygar

Frontmatter

[More Information](#)

“Data Science practitioners tend to be unaware of how easy it is for adversaries to manipulate and misuse adaptive machine learning systems. This book demonstrates the severity of the problem by providing a taxonomy of attacks and studies of adversarial learning. It analyzes older attacks as well as recently discovered surprising weaknesses in deep learning systems. A variety of defenses are discussed for different learning systems and attack types that could help researchers and developers design systems that are more robust to attacks.”

Richard Lippmann, *Lincoln Laboratory, MIT*

“This is a timely book. Right time and right book, written with an authoritative but inclusive style. Machine learning is becoming ubiquitous. But for people to trust it, they first need to understand how reliable it is.”

Fabio Roli, *University of Cagliari*

Cambridge University Press

978-1-107-04346-6 — Adversarial Machine Learning

Anthony D. Joseph , Blaine Nelson , Benjamin I. P. Rubinstein , J. D. Tygar

Frontmatter

[More Information](#)

Adversarial Machine Learning

ANTHONY D. JOSEPH

University of California, Berkeley

BLAINE NELSON

Google

BENJAMIN I. P. RUBINSTEIN

University of Melbourne

J. D. TYGAR

University of California, Berkeley



CAMBRIDGE
UNIVERSITY PRESS

Cambridge University Press

978-1-107-04346-6 — Adversarial Machine Learning

Anthony D. Joseph , Blaine Nelson , Benjamin I. P. Rubinstein , J. D. Tygar

Frontmatter

[More Information](#)

CAMBRIDGE UNIVERSITY PRESS

University Printing House, Cambridge CB2 8BS, United Kingdom

One Liberty Plaza, 20th Floor, New York, NY 10006, USA

477 Williamstown Road, Port Melbourne, VIC 3207, Australia

314–321, 3rd Floor, Plot 3, Splendor Forum, Jasola District Centre,
New Delhi - 110025, India

79 Anson Road, #06-04/06, Singapore 079906

Cambridge University Press is part of the University of Cambridge.

It furthers the University's mission by disseminating knowledge in the pursuit of
education, learning and research at the highest international levels of excellence.

www.cambridge.org

Information on this title: www.cambridge.org/9781107043466

DOI:10.1017/9781107338548

© Cambridge University Press 2019

This publication is in copyright. Subject to statutory exception
and to the provisions of relevant collective licensing agreements,
no reproduction of any part may take place without the written
permission of Cambridge University Press.

First published 2019

Printed and bound in Great Britain by Clays Ltd, Elcograf S.p.A.

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

Library of Congress Cataloging-in-Publication data

Names: Joseph, Anthony D., author. | Nelson, Blaine, author. | Rubinstein, Benjamin I. P., author. |
Tygar, J. D., author.

Title: Adversarial machine learning / Anthony D. Joseph, University of California, Berkeley, Blaine Nelson,
Google, Benjamin I.P. Rubinstein, University of Melbourne, J.D. Tygar, University of California, Berkeley.

Description: Cambridge, United Kingdom ; New York, NY : Cambridge University Press, 2019. |

Includes bibliographical references and index.

Identifiers: LCCN 2017026016 | ISBN 9781107043466 (hardback)

Subjects: LCSH: Machine learning. | Computer security. | BISAC: COMPUTERS / Security / General.

Classification: LCC Q325.5 .J69 2017 | DDC 006.3/1 – dc23

LC record available at <https://lccn.loc.gov/2017026016>

ISBN 978-1-107-04346-6 Hardback

Cambridge University Press has no responsibility for the persistence or accuracy
of URLs for external or third-party internet websites referred to in this publication,
and does not guarantee that any content on such websites is, or will remain,
accurate or appropriate.

Contents

<i>List of Symbols</i>	page xi
<i>Acknowledgments</i>	xiii

Part I Overview of Adversarial Machine Learning

1	Introduction	3
	1.1 Motivation	4
	1.2 A Principled Approach to Secure Learning	11
	1.3 Chronology of Secure Learning	14
	1.4 Overview	17
2	Background and Notation	20
	2.1 Basic Notation	20
	2.2 Statistical Machine Learning	20
	2.2.1 Data	22
	2.2.2 Hypothesis Space	24
	2.2.3 The Learning Model	24
	2.2.4 Supervised Learning	25
	2.2.5 Other Learning Paradigms	28
3	A Framework for Secure Learning	29
	3.1 Analyzing the Phases of Learning	29
	3.2 Security Analysis	31
	3.2.1 Security Goals	32
	3.2.2 Threat Model	32
	3.2.3 Discussion of Machine Learning Applications in Security	33
	3.3 Framework	34
	3.3.1 Taxonomy	34
	3.3.2 The Adversarial Learning Game	36
	3.3.3 Characteristics of Adversarial Capabilities	36
	3.3.4 Attacks	38
	3.3.5 Defenses	38

3.4	Exploratory Attacks	38
3.4.1	The Exploratory Game	39
3.4.2	Exploratory Integrity Attacks	40
3.4.3	Exploratory Availability Attacks	44
3.4.4	Defending against Exploratory Attacks	45
3.5	Causative Attacks	49
3.5.1	The Causative Game	50
3.5.2	Causative Integrity Attacks	51
3.5.3	Causative Availability Attacks	53
3.5.4	Defending against Causative Attacks	54
3.6	Repeated Learning Games	58
3.6.1	Repeated Learning Games in Security	60
3.7	Privacy-Preserving Learning	61
3.7.1	Differential Privacy	62
3.7.2	Exploratory and Causative Privacy Attacks	64
3.7.3	Utility despite Randomness	65

Part II Causative Attacks on Machine Learning

4	Attacking a Hypersphere Learner	69
4.1	Causative Attacks on Hypersphere Detectors	70
4.1.1	Learning Assumptions	71
4.1.2	Attacker Assumptions	72
4.1.3	Analytic Methodology	73
4.2	Hypersphere Attack Description	73
4.2.1	Displacing the Centroid	75
4.2.2	Formal Description of the Attack	79
4.2.3	Characteristics of Attack Sequences	80
4.3	Optimal Unconstrained Attacks	83
4.3.1	Optimal Unconstrained Attack: Stacking Blocks	83
4.4	Imposing Time Constraints on the Attack	84
4.4.1	Stacking Blocks of Variable Mass	86
4.4.2	An Alternate Formulation	87
4.4.3	The Optimal Relaxed Solution	88
4.5	Attacks against Retraining with Data Replacement	91
4.5.1	Average-out and Random-out Replacement Policy	92
4.5.2	Nearest-out Replacement Policy	94
4.6	Constrained Attackers	96
4.6.1	Greedy Optimal Attacks	98
4.6.2	Attacks with Mixed Data	100
4.6.3	Extensions	102
4.7	Summary	103

5	Availability Attack Case Study: SpamBayes	105
5.1	The SpamBayes Spam Filter	106
5.1.1	SpamBayes' Training Algorithm	106
5.1.2	SpamBayes' Predictions	108
5.1.3	SpamBayes' Model	109
5.2	Threat Model for SpamBayes	112
5.2.1	Attacker Goals	113
5.2.2	Attacker Knowledge	113
5.2.3	Training Model	114
5.2.4	The Contamination Assumption	115
5.3	Causative Attacks against SpamBayes' Learner	115
5.3.1	Causative Availability Attacks	116
5.3.2	Causative Integrity Attacks—Pseudospam	119
5.4	The Reject on Negative Impact (RONI) Defense	119
5.5	Experiments with SpamBayes	120
5.5.1	Experimental Method	120
5.5.2	Dictionary Attack Results	122
5.5.3	Focused Attack Results	124
5.5.4	Pseudospam Attack Experiments	126
5.5.5	RONI Results	128
5.6	Summary	131
6	Integrity Attack Case Study: PCA Detector	134
6.1	PCA Method for Detecting Traffic Anomalies	137
6.1.1	Traffic Matrices and Volume Anomalies	137
6.1.2	Subspace Method for Anomaly Detection	138
6.2	Corrupting the PCA Subspace	140
6.2.1	The Threat Model	140
6.2.2	Uninformed Chaff Selection	141
6.2.3	Locally Informed Chaff Selection	141
6.2.4	Globally Informed Chaff Selection	142
6.2.5	Boiling Frog Attacks	143
6.3	Corruption-Resilient Detectors	144
6.3.1	Intuition	145
6.3.2	PCA-GRID	147
6.3.3	Robust Laplace Threshold	148
6.4	Empirical Evaluation	149
6.4.1	Setup	149
6.4.2	Identifying Vulnerable Flows	152
6.4.3	Evaluation of Attacks	153
6.4.4	Evaluation of ANTIDOTE	156
6.4.5	Empirical Evaluation of the Boiling Frog Poisoning Attack	159
6.5	Summary	162

Part III Exploratory Attacks on Machine Learning

7	Privacy-Preserving Mechanisms for SVM Learning	167
7.1	Privacy Breach Case Studies	167
7.1.1	Massachusetts State Employees Health Records	167
7.1.2	AOL Search Query Logs	168
7.1.3	The Netflix Prize	168
7.1.4	Deanonymizing Twitter Pseudonyms	169
7.1.5	Genome-Wide Association Studies	169
7.1.6	Ad Microtargeting	170
7.1.7	Lessons Learned	170
7.2	Problem Setting: Privacy-Preserving Learning	170
7.2.1	Differential Privacy	171
7.2.2	Utility	173
7.2.3	Historical Research Directions in Differential Privacy	174
7.3	Support Vector Machines: A Brief Primer	176
7.3.1	Translation-Invariant Kernels	178
7.3.2	Algorithmic Stability	179
7.4	Differential Privacy by Output Perturbation	179
7.5	Differential Privacy by Objective Perturbation	183
7.6	Infinite-Dimensional Feature Spaces	185
7.7	Bounds on Optimal Differential Privacy	192
7.7.1	Upper Bounds	193
7.7.2	Lower Bounds	194
7.8	Summary	197
8	Near-Optimal Evasion of Classifiers	199
8.1	Characterizing Near-Optimal Evasion	202
8.1.1	Adversarial Cost	203
8.1.2	Near-Optimal Evasion	204
8.1.3	Search Terminology	206
8.1.4	Multiplicative vs. Additive Optimality	208
8.1.5	The Family of Convex-Inducing Classifiers	210
8.2	Evasion of Convex Classes for ℓ_1 Costs	211
8.2.1	ϵ -IMAC Search for a Convex \mathcal{X}_f^+	212
8.2.2	ϵ -IMAC Learning for a Convex \mathcal{X}_f^-	221
8.3	Evasion for General ℓ_p Costs	225
8.3.1	Convex Positive Set	225
8.3.2	Convex Negative Set	231
8.4	Summary	231
8.4.1	Open Problems in Near-Optimal Evasion	232
8.4.2	Alternative Evasion Criteria	234
8.4.3	Real-World Evasion	236

Part IV Future Directions in Adversarial Machine Learning

9	Adversarial Machine Learning Challenges	241
9.1	Discussion and Open Problems	245
9.1.1	Unexplored Components of the Adversarial Game	245
9.1.2	Development of Defensive Technologies	247
9.2	Review of Open Problems	250
9.3	Concluding Remarks	252

Part V Appendixes

Appendix A: Background for Learning and Hyper-Geometry	255
Appendix B: Full Proofs for Hypersphere Attacks	266
Appendix C: Analysis of SpamBayes	276
Appendix D: Full Proofs for Near-Optimal Evasion	285
<i>Glossary</i>	295
<i>References</i>	307
<i>Index</i>	322

Cambridge University Press
978-1-107-04346-6 — Adversarial Machine Learning
Anthony D. Joseph , Blaine Nelson , Benjamin I. P. Rubinstein , J. D. Tygar
Frontmatter
[More Information](#)

Symbols

$A(\cdot)$: The adversary's cost function on \mathcal{X} (see Section 8.1.1). See 203–209, 211, 212, 214, 216, 219–221, 228, 231, 234, 235

\mathbb{D} : A set of data points (see also: dataset). See 23–26, 183

N : The number of data points in the training dataset used by a learning algorithm; i.e., $N \triangleq |\mathbb{D}^{(\text{train})}|$. See 21, 23, 25–27, 36, 38–40, 46, 47, 50, 54, 56, 183, 184, 256

$\mathbb{D}^{(\text{train})}$: A dataset used by a training algorithm to construct or select a classifier (see also: dataset). See 21, 25, 26, 36, 39, 40, 48, 50, 107, 120, 128

$\mathbb{D}^{(\text{eval})}$: A dataset used to evaluate a classifier (see also: dataset). See 21, 22, 25, 27, 36, 39, 40, 46, 48, 50, 51, 128

\triangleq : Symbol used to provide a definition. See 23, 24, 26, 57, 58, 60, 107, 108, 137, 139, 141, 142, 149, 153, 204, 206, 256–259, 265, 276, 278, 279, 281, 282

ϵ -*IMAC*: The set of objects in \mathcal{X}_f^- within a cost of $1 + \epsilon$ of the *MAC*, or any of the members of this set (see also: $MAC(f, A)$). See 204–206, 209–214, 216, 219–221, 225, 229, 231–235, 237, 251

$f(\cdot)$: The classifier function or hypothesis learned by a training procedure $H^{(N)}$ from the dataset $\mathbb{D}^{(\text{train})}$ (see also: classifier). See 21, 24–27, 39, 40, 48–51, 54, 71, 74, 102, 120, 139, 174, 176–183, 189, 195, 196, 202–207, 209, 210, 212, 215, 217–220, 234–237, 251

L_ϵ : The number of steps required by a binary search to achieve ϵ -optimality (see Section 8.1.3). See 205, 210–212, 214–218, 220, 221, 225, 228, 232

$MAC(f, A)$: The largest lower bound on the adversary's cost A over \mathcal{X}_f^- (see also: Equation 8.2). See 204, 206, 207, 210, 212–214, 216, 219, 221, 229, 230, 234, 235

\mathfrak{N} : The set of natural numbers, $\{1, 2, 3, \dots\}$. See 77–79, 81, 83, 88, 137, 256, 257, 276

\mathfrak{N}_0 : The set of all whole numbers, $\{0, 1, 2, \dots\}$. See 73, 77–79, 81, 82, 86, 256

$\|\cdot\|$: A non-negative function defined on a vector space that is positive homogeneous and obeys the triangle inequality (see also: norm). See 145, 147, 149, 152, 153, 203, 226, 227, 257

- ℓ_p ($p > 0$) : A norm on a multidimensional real-value space defined in Appendix A.1 by Equation (A.1) and denoted by $\| \cdot \|_p$. See 11, 18, 200, 203, 204, 208, 210–214, 216–218, 220, 221, 223, 225–234, 244, 245, 260, 261, 264, 265
- $m_{\mathbb{C}} (\cdot)$: A function that defines a distance metric for a convex set \mathbb{C} relative to some central element $\mathbf{x}^{(c)}$ in the interior of \mathbb{C} (see also: Minkowski metric). See 210, 211
- $N^{(h)}$: The total number of ham messages in the training dataset. See 107, 108, 111, 277–280
- $n_j^{(h)}$: The number of occurrences of the j^{th} token in training ham messages. See 107, 108, 111, 277–280
- $N^{(s)}$: The total number of spam messages in the training dataset. See 107, 108, 111, 119, 277–280
- $n_j^{(s)}$: The number of occurrences of the j^{th} token in training spam messages. See 107, 108, 111, 119, 277–280
- Q** : The matrix of network flow data. See 137, 138, 152
- R** : The routing matrix that describes the links used to route each OD flow. See 138, 142, 152
- \Re : The set of all real numbers. See 23–25, 27, 142, 256–259
 - \Re_{0+} : The set of all real numbers greater than or equal to zero. See 26, 203, 211, 256, 276
 - \Re_+ : The set of all real numbers greater than zero. See 27, 216, 256, 257
 - \Re^D : The D -dimensional real-valued space. See 24, 139, 142, 143, 147, 202, 216, 226, 257, 259
- x** : A data point from the input space \mathcal{X} (see also data point). See 22–24, 138, 139, 141, 145, 147, 148, 200, 203–206, 208–211, 213, 223, 224, 226, 255
 - \mathbf{x}^A : A (malicious) data point that the adversary would like to sneak past the detector. See 70, 203–206, 209–215, 217, 218, 221–223, 225, 226, 231, 233–235, 261, 264, 265
- \mathcal{X} : The input space of the data (see also: input space). See 22–25, 49, 202, 203, 206, 208, 210, 211, 224, 233, 235, 236, 259, 260
 - D : The dimensionality of the input space \mathcal{X} . See 22, 23, 202–205, 212–218, 220, 223–232, 235, 237, 259
 - \mathcal{X}_f^- : The negative class for the deterministic classifier f (see also: negative class). See 203–206, 208, 210–213, 219, 221–224, 226, 231, 233, 235, 236
 - \mathcal{X}_f^+ : The positive class for the deterministic classifier f (see also: positive class). See 203, 205, 210–214, 216, 218, 233, 234
- y : A label from the response space \mathcal{Y} (see also: label). See 23, 26, 27, 107
- \mathcal{Y} : The response space of the data (see also response space). See 23–27, 59, 203
- \mathbb{Z} : The set of all integers. See 23, 256, 258

Acknowledgments

We gratefully acknowledge the contributions and assistance of our colleagues in making this book possible, who include but are not limited to Sadia Afroz, Scott Alfeld, Tansu Alpcan, Rekha Bachwani, Marco Barreno, Adam Barth, Peter Bartlett, Battista Biggio, Chris Cai, Fuching Jack Chi, David Fifield, Laurent El Ghaoui, Barbara Goto, Rachel Greenstadt, Yi Han, Ling Huang, Michael Jordan, Alex Kantchelian, Hideaki Kawabata, Marius Kloft, Pavel Laskov, Shing-hon Lau, Chris Leckie, Steven Lee, Justin Ma, Steve Martin, Brad Miller, Satish Rao, Fabio Roli, Udam Saini, Tobias Scheffer, Russell Sears, Anil Sewani, Arunesh Sinha, Dawn Song, Nedim Šrndić, Charles Sutton, Nina Taft, Anthony Tran, Michael Tschantz, Kai Xai, Takumi Yamamoto, and Qi Zhong. We additionally thank Matthias Bussas and Marius Kloft for their careful proofreading of Chapter 4 and the staff at Cambridge University Press including Heather Brolly and Julie Lancashire for their help in preparing this manuscript. We would also like to thank the many colleagues with whom we have had fruitful discussions at the Dagstuhl Perspectives Workshop on Machine Learning Methods for Computer Security (Joseph, Laskov, Roli, Tygar, & Nelson 2013) and at the ACM Workshop on Artificial Intelligence and Security (AISec) and other workshops and conferences.

The authors are currently at the University of California, Berkeley, the University of Melbourne, and Google. We thank these institutions. While we were writing this book, some of the authors were at Universität Tübingen, Universität Potsdam, Università di Cagliari, IBM Research, and Microsoft Research, and we also thank those institutions. We offer special thanks to our support staff, including Angie Abbatecola, Kattt Atchley, Carlyn Chinen, Barbara Goto, Damon Hinson, Michaela Iglesia, Shane Knapp, Jey Kottalam, Jon Kuroda, Lena Lau-Stewart, Christian Legg, and Boban Zarkovich.

We are grateful for the financial sponsors of this research. We received U.S. government funding from the Air Force Office of Scientific Research, Homeland Security Advanced Research Projects Agency, National Science Foundation, and State Department DRL, and in some cases through UC Berkeley laboratories (DETERlab and TRUST). Some authors received additional support from the Alexander von Humboldt Foundation, the Australian Research Council (DE160100584), the Center for Long-Term Cybersecurity, the Future of Life Institute, Oak Ridge National Laboratory, and the Open Technology Fund. The opinions expressed in this book are solely those of the authors and do not necessarily reflect the views of any funder.

The authors could not have written this book without the support, encouragement, and patience of their friends and families.

Cambridge University Press
978-1-107-04346-6 — Adversarial Machine Learning
Anthony D. Joseph , Blaine Nelson , Benjamin I. P. Rubinstein , J. D. Tygar
Frontmatter
[More Information](#)
