

# Active Learning



# Synthesis Lectures on Artificial Intelligence and Machine Learning

## Editor

**Ronald J. Brachman**, *Yahoo! Research*  
**William W. Cohen**, *Carnegie Mellon University*  
**Thomas Dietterich**, *Oregon State University*

## Active Learning

Burr Settles  
2012

## Planning with Markov Decision Processes: An AI Perspective

Mausam and Andrey Kolobov  
2012

## Computational Aspects of Cooperative Game Theory

Georgios Chalkiadakis, Edith Elkind, and Michael Wooldridge  
2011

## Representations and Techniques for 3D Object Recognition and Scene Interpretation

Derek Hoiem and Silvio Savarese  
2011

## A Short Introduction to Preferences: Between Artificial Intelligence and Social Choice

Francesca Rossi, Kristen Brent Venable, and Toby Walsh  
2011

## Human Computation

Edith Law and Luis von Ahn  
2011

## Trading Agents

Michael P. Wellman  
2011

### Visual Object Recognition

Kristen Grauman and Bastian Leibe

2011

### Learning with Support Vector Machines

Colin Campbell and Yiming Ying

2011

### Algorithms for Reinforcement Learning

Csaba Szepesvári

2010

### Data Integration: The Relational Logic Approach

Michael Genesereth

2010

### Markov Logic: An Interface Layer for Artificial Intelligence

Pedro Domingos and Daniel Lowd

2009

### Introduction to Semi-Supervised Learning

Xiaojin Zhu and Andrew B. Goldberg

2009

### Action Programming Languages

Michael Thielscher

2008

### Representation Discovery using Harmonic Analysis

Sridhar Mahadevan

2008

### Essentials of Game Theory: A Concise Multidisciplinary Introduction

Kevin Leyton-Brown and Yoav Shoham

2008

### A Concise Introduction to Multiagent Systems and Distributed Artificial Intelligence

Nikos Vlassis

2007

### Intelligent Autonomous Robotics: A Robot Soccer Case Study

Peter Stone

2007

© Springer Nature Switzerland AG 2022

Reprint of original edition © Morgan & Claypool 2012

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means—electronic, mechanical, photocopy, recording, or any other except for brief quotations in printed reviews, without the prior permission of the publisher.

Active Learning

Burr Settles

ISBN:978-3-031-00432-2 paperback

ISBN: 978-3-031-01560-1 ebook

DOI 10.1007/978-3-031-01560-1

A Publication in the Springer series

*SYNTHESIS LECTURES ON ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING*

Lecture #18

Series Editors: Ronald J. Brachman, *Yahoo Research*

William W. Cohen, *Carnegie Mellon University*

Thomas Dietterich, *Oregon State University*

Series ISSN

Synthesis Lectures on Artificial Intelligence and Machine Learning

Print 1939-4608 Electronic 1939-4616

# Active Learning

Burr Settles  
Carnegie Mellon University

*SYNTHESIS LECTURES ON ARTIFICIAL INTELLIGENCE AND MACHINE  
LEARNING #18*

## ABSTRACT

The key idea behind active learning is that a machine learning algorithm can perform better with less training if it is allowed to *choose* the data from which it learns. An active learner may pose “queries,” usually in the form of unlabeled data instances to be labeled by an “oracle” (e.g., a human annotator) that already understands the nature of the problem. This sort of approach is well-motivated in many modern machine learning and data mining applications, where unlabeled data may be abundant or easy to come by, but training labels are difficult, time-consuming, or expensive to obtain.

This book is a general introduction to active learning. It outlines several scenarios in which queries might be formulated, and details many query selection algorithms which have been organized into four broad categories, or “query selection frameworks.” We also touch on some of the theoretical foundations of active learning, and conclude with an overview of the strengths and weaknesses of these approaches in practice, including a summary of ongoing work to address these open challenges and opportunities.

## KEYWORDS

active learning, expected error reduction, hierarchical sampling, optimal experimental design, query by committee, query by disagreement, query learning, uncertainty sampling, variance reduction

*Dedicated to my family and friends,  
who keep me asking questions.*

# Contents

	<b>Preface</b> .....	<b>xi</b>
	<b>Acknowledgments</b> .....	<b>xiii</b>
<b>1</b>	<b>Automating Inquiry</b> .....	<b>1</b>
	1.1 A Thought Experiment .....	1
	1.2 Active Learning .....	3
	1.3 Scenarios for Active Learning .....	5
<b>2</b>	<b>Uncertainty Sampling</b> .....	<b>11</b>
	2.1 Pushing the Boundaries .....	11
	2.2 An Example .....	12
	2.3 Measures of Uncertainty .....	13
	2.4 Beyond Classification .....	16
	2.5 Discussion .....	18
<b>3</b>	<b>Searching Through the Hypothesis Space</b> .....	<b>21</b>
	3.1 The Version Space .....	21
	3.2 Uncertainty Sampling as Version Space Search .....	22
	3.3 Query by Disagreement .....	24
	3.4 Query by Committee .....	28
	3.5 Discussion .....	32
<b>4</b>	<b>Minimizing Expected Error and Variance</b> .....	<b>37</b>
	4.1 Expected Error Reduction .....	37
	4.2 Variance Reduction .....	40
	4.3 Batch Queries and Submodularity .....	44
	4.4 Discussion .....	46
<b>5</b>	<b>Exploiting Structure in Data</b> .....	<b>47</b>
	5.1 Density-Weighted Methods .....	47
	5.2 Cluster-Based Active Learning .....	49



5.3	Active + Semi-Supervised Learning .....	53
5.4	Discussion .....	54
<b>6</b>	<b>Theory .....</b>	<b>55</b>
6.1	A Unified View .....	55
6.2	A PAC Bound for Active Learning .....	57
6.3	Discussion .....	61
<b>7</b>	<b>Practical Considerations .....</b>	<b>63</b>
7.1	Which Algorithm is Best? .....	63
7.2	Real Labeling Costs .....	65
7.3	Alternative Query Types .....	68
7.4	Skewed Label Distributions .....	72
7.5	Unreliable Oracles .....	73
7.6	Multi-Task Active Learning .....	74
7.7	Data Reuse and the Unknown Model Class .....	76
7.8	Stopping Criteria .....	77
<b>A</b>	<b>Nomenclature Reference .....</b>	<b>79</b>
	<b>Bibliography .....</b>	<b>81</b>
	<b>Author's Biography .....</b>	<b>97</b>
	<b>Index .....</b>	<b>99</b>

# Preface

Machine learning is the study of computer systems that improve through experience. Active learning is the study of machine learning systems that improve by asking questions. So why ask questions? (Good question.) The key hypothesis is that if the learner is allowed to choose the data from which it learns — to be active, curious, or exploratory, if you will — it can perform better with less training. Consider that in order for most supervised machine learning systems to perform well they must often be trained on many hundreds or thousands of labeled data instances. Sometimes these labels come at little or no cost, but for many real-world applications, labeling is a difficult, time-consuming, or expensive process. Fortunately in today's data-drenched society, unlabeled data are often abundant (or at least easier to acquire). This suggests that much can be gained by using active learning systems to ask effective questions, exploring the most informative nooks and crannies of a vast data landscape (rather than randomly and expensively sampling data from the domain).

This book was written with students, researchers, and other practitioners of machine learning in mind. It will be most useful to those who are already familiar with the basics of machine learning and are looking for a thorough but gentle introduction to active learning techniques. We will assume a basic familiarity with probability and statistics, some linear algebra, and common supervised learning algorithms. An introductory text in artificial intelligence ([Russell and Norvig, 2003](#)) or machine learning ([Bishop, 2006](#); [Duda et al., 2001](#); [Mitchell, 1997](#)) is probably sufficient background. Ardent students of computational learning theory might find themselves annoyed at the lack of rigorous mathematical analysis in this book. This is partially because, until very recently, there has been little interaction between the sub-communities of theory and practice within active learning. While some discussion of underlying theory can be found in [Chapter 6](#), most of this volume is focused on algorithms at a qualitative level, motivated by issues of practice.

The presentation includes a mix of contrived, illustrative examples as well as benchmark-style evaluations that compare and contrast various algorithms on real data sets. However, I caution the reader not to take any of these results at face value, as there are many factors at play when choosing an active learning approach. It is my hope that this book does a good job of pointing out all the subtleties at play, and helps the reader gain some intuition about which approaches are most appropriate for the task at hand.

This active learning book is the synthesis of a previous literature survey ([Settles, 2009](#)) with material from other lectures and talks I have given on the subject. It is meant to be used as an introduction and reference for researchers, or as a supplementary text for courses in machine learning — supporting a week or two of lectures — rather than as a textbook for a complete full-term course on active learning. (Despite two decades of research, I am not sure that there is enough breadth or

depth of understanding to warrant a full-semester course dedicated to active learning. At least not yet!) Here is a road map:

- Chapter 1 introduces the basic idea of, and motivations for, active learning.
- Chapters 2–5 focus on different “query frameworks,” or families of active learning heuristics. These include several algorithms each.
- Chapter 6 covers some of the theoretical foundations of active learning.
- Chapter 7 summarizes the various pros and cons of algorithms covered in this book. It outlines several important considerations for active learning in practice, and discusses recent work aimed at addressing these practical issues.

I have attempted to wade through the thicket of papers and distill active learning approaches into core conceptual categories, characterizing their strengths and weaknesses in both theory and practice. I hope you enjoy it and find it useful in your work.

Supplementary materials, as well as a mailing list, links to video lectures, software implementations, and other resources for active learning can be found online at <http://active-learning.net>.

Burr Settles  
May 2012

# Acknowledgments

This book has a roundabout history, and there are a lot of people to thank along the way. It grew out of an informal literature survey I wrote on active learning (Settles, 2009) which in turn began as a chapter in my PhD thesis. During that phase of my career I am indebted to my committee, Mark Craven, Jude Shavlik, Xiaojin “Jerry” Zhu, David Page, and Lewis Friedland, who encouraged me to expand on my review and make it publicly available. There has been a lot of work in active learning over the past two decades, from simple heuristics to complex and crazy ideas coming from a variety of subfields in AI and statistics. The survey was my attempt to curate, organize, and make sense of it for myself; to help me understand how my work fit into the overall landscape.

Thanks to John Langford, who mentioned the survey on his popular machine learning blog<sup>1</sup>. As a result, many other people found it and found it helpful as well. Several people encouraged me to write this book. To that end, Jude Shavlik and Edith Law (independently) introduced me to Michael Morgan. Thanks to Michael, William Cohen, Tom Dietterich, and others at Morgan & Claypool for doing their best to keep things on schedule, and for patiently encouraging me through the process of expanding what was a literature review into more of a tutorial or textbook. Thanks also to Tom Mitchell for his support and helpful advice on how to organize and write a book.

Special thanks to Steve Hanneke and Sanjoy Dasgupta for the detailed feedback on both the original survey and the expanded manuscript. Chapter 6 is particularly indebted to their comments as well as their research. I also found Dasgupta’s review of active learning from a theoretical perspective (Dasgupta, 2010) quite helpful. The insights and organization of ideas presented here are not wholly my own, but draw on conversations I have had with numerous people. In addition to the names mentioned above, I would like to thank Josh Attenberg, Jason Baldrige, Carla Brodley, Aron Culotta, Pinar Donmez, Miroslav Dudík, Gregory Druck, Jacob Eisenstein, Russ Greiner, Carlos Guestrin, Robbie Haertel, Ashish Kapoor, Percy Liang, Andrew McCallum, Prem Melville, Clare Monteleoni, Ray Mooney, Foster Provost, Soumya Ray, Eric Ringger, Teddy Seidenfeld, Kevin Small, Partha Talukdar, Katrin Tomanek, Byron Wallace, and other colleagues for turning me on to papers, ideas, and perspectives that I might have otherwise overlooked. I am sure there are other names I have forgotten to list here, but know that I appreciate all the ongoing discussions on active learning (and machine learning in general), both online and in person. Thanks also to Daniel Hsu, Eric Baum, Nicholas Roy, and their coauthors (some listed above) for kindly allowing me to reuse figures from their publications.

I would like to thank my parents for getting me started, and my wife Natalie for keeping me going. She remained amazingly supportive during my long hours of writing (and re-writing).

<sup>1</sup><http://hunch.net>

Whenever I was stumped or frustrated, she was quick to offer a fresh perspective: “Look at you, you’re writing a book!” Lo and behold, I have written a book. I hope you enjoy the book.

While writing this book, I was supported by the Defense Advanced Research Projects Agency (under contracts FA8750-08-1-0009 and AF8750-09-C-0179), the National Science Foundation (under grant IIS-0968487), and Google. The text also includes material written while I was supported by a grant from National Human Genome Research Institute (HGRI). Any opinions, findings and conclusions, or recommendations expressed in this material are mine and do not necessarily reflect those of the sponsors.

Burr Settles

May 2012