# Index

Page numbers in *italics* are recommended to be consulted first. Page numbers in **bold** contain boxed algorithms.

Adaptive Computation and Machine Learning

Francis Bach, Editor

*Bioinformatics: The Machine Learning Approach*, Pierre Baldi and Søren Brunak

*Reinforcement Learning: An Introduction*, Richard S. Sutton and Andrew G. Barto

*Graphical Models for Machine Learning and Digital Communication*, Brendan J. Frey

*Learning in Graphical Models*, Michael I. Jordan

*Causation, Prediction, and Search*, second edition, Peter Spirtes, Clark Glymour, and Richard Scheines

*Principles of Data Mining*, David Hand, Heikki Mannila, and Padhraic Smyth

*Bioinformatics: The Machine Learning Approach*, second edition, Pierre Baldi and Søren Brunak

*Learning Kernel Classifiers: Theory and Algorithms*, Ralf Herbrich

*Learning with Kernels: Support Vector Machines, Regularization, Optimization, and Beyond*, Bernhard Schölkopf and Alexander J. Smola

*Introduction to Machine Learning*, Ethem Alpaydin

*Gaussian Processes for Machine Learning*, Carl Edward Rasmussen and Christopher K.I. Williams

*Semi-Supervised Learning*, Olivier Chapelle, Bernhard Schölkopf, and Alexander Zien, Eds.

*The Minimum Description Length Principle*, Peter D. Grünwald

*Introduction to Statistical Relational Learning*, Lise Getoor and Ben Taskar, Eds.

*Probabilistic Graphical Models: Principles and Techniques*, Daphne Koller and Nir Friedman

*Introduction to Machine Learning*, second edition, Ethem Alpaydin

*Machine Learning in Non-Stationary Environments: Introduction to Covariate Shift Adaptation*, Masashi Sugiyama and Motoaki Kawanabe

*Boosting: Foundations and Algorithms*, Robert E. Schapire and Yoav Freund

*Machine Learning: A Probabilistic Perspective*, Kevin P. Murphy

*Foundations of Machine Learning*, Mehryar Mohri, Afshin Rostami, and Ameet Talwalker

*Introduction to Machine Learning*, third edition, Ethem Alpaydin

*Deep Learning*, Ian Goodfellow, Yoshua Bengio, and Aaron Courville

*Elements of Causal Inference*, Jonas Peters, Dominik Janzing, and Bernhard Schölkopf

*Machine Learning for Data Streams, with Practical Examples in MOA*, Albert Bifet, Ricard Gavaldà, Geoffrey Holmes, Bernhard Pfahringer