# A First Course in
# Machine
# Learning

Simon Rogers
Mark Girolami

# A First Course in
# Machine
# Learning

Simon Rogers
Mark Girolami

MATLAB® is a trademark of The MathWorks, Inc. and is used with permission. The MathWorks does not warrant the accuracy of the text or exercises in this book. This book's use or discussion of MATLAB® software or related products does not constitute endorsement or sponsorship by The MathWorks of a particular pedagogical approach or particular use of the MATLAB® software.

# A First Course in

# Machine Learning

# Chapman & Hall/CRC
## Machine Learning & Pattern Recognition Series

SERIES EDITORS

**Ralf Herbrich and Thore Graepel**
Microsoft Research Ltd.
Cambridge, UK

## AIMS AND SCOPE

This series reflects the latest advances and applications in machine learning and pattern recognition through the publication of a broad range of reference works, textbooks, and handbooks. The inclusion of concrete examples, applications, and methods is highly encouraged. The scope of the series includes, but is not limited to, titles in the areas of machine learning, pattern recognition, computational intelligence, robotics, computational/statistical learning theory, natural language processing, computer vision, game AI, game theory, neural networks, computational neuroscience, and other relevant topics, such as machine learning applied to bioinformatics or cognitive science, which might be proposed by potential contributors.

## PUBLISHED TITLES

MACHINE LEARNING: An Algorithmic Perspective
*Stephen Marsland*

HANDBOOK OF NATURAL LANGUAGE PROCESSING,
Second Edition
*Nitin Indurkhya and Fred J. Damerau*

UTILITY-BASED LEARNING FROM DATA
*Craig Friedman and Sven Sandow*

A FIRST COURSE IN MACHINE LEARNING
*Simon Rogers and Mark Girolami*

# Preface

Machine learning is rapidly becoming one of the most important areas of general practice, research and development activity within computing science. This is reflected in the scale of the academic research area devoted to the subject and the active recruitment of machine learning specialists by major international banks and financial institutions as well as companies such as Microsoft®, Google®, Yahoo® and Amazon®.

This growth can be partly explained by the increase in the quantity and diversity of measurements we are able to make of the world. A particularly fascinating example arises from the wave of new biological measurement technologies that preceded the sequencing of the first genomes. It is now possible to measure the detailed molecular state of an organism in ways that would have been hard to imagine only a short time ago. Such measurements go far beyond our understanding of these organisms and machine learning techniques have been heavily involved in the distillation of useful structures from them.

This book is based on material presented in a machine learning course in the School of Computing Science at the University of Glasgow, UK. The course, presented to final year undergraduates and taught by postgraduates, is made up of 20 hour-long lectures and 10 hour-long laboratory sessions. In such a short teaching period, it is impossible to cover more than a small fraction of the material that now comes under the banner of machine learning. Our intention when teaching this course, therefore, is to present the core mathematical and statistical techniques required to understand some of the most popular machine learning algorithms and then present a few of these algorithms that span the main problem areas within machine learning: classification, clustering and projection. At the end of the course, the students should have the knowledge and confidence to be able to explore machine learning literature to find methods that are more appropriate for them. The same is hopefully true of readers of this book.

Due to the varying mathematical literacy of students taking the course, we assume only very minor mathematical pre-requisites. An undergraduate student from computer science, engineering, physics (or any other numerical subject) should have no problem. This does not preclude those without such experience – additional mathematical explanations appear throughout the text in comment boxes. In addition, important equations have been highlighted – it is worth spending time understanding these equations before proceeding.

Students attending this course often find the practical sessions very useful. Experimenting with the various algorithms and concepts helps transfer them from an abstract set of equations into something that could be used to solve real problems. We have attempted to transfer this to the book through an extensive collection of MATLAB®/Octave[1] scripts, available from the associated web page and referenced throughout the text. These scripts enable the user to recreate plots that appear in the book and investigate changing model specifications and parameter values.

Finally, the machine learning methods that are covered in this book are our choice of those we feel students should understand. In limited space and time, we think that it is more worthwhile to give detailed descriptions and derivations for a small number of algorithms than attempt to cover many algorithms at a lower level of detail – many people will not find their favourite algorithms within this book!

MATLAB® is a registered trademark of The MathWorks, Inc.
For product information, please contact:

The MathWorks, Inc.
3 Apple Hill Drive
Natick MA 01760-2098 USA
Tel: 508-647-7000
Fax: 508-647-7001
E-mail: info@mathworks.com
Web: www.mathworks.com

Simon Rogers and Mark Girolami

---

[1] A free mathematical software environment, available from www.gnu.org/software/octave/

# Contents

# List of Tables

# List of Figures