# Statistics and Computing

Statistics and Computing (SC) includes monographs and advanced texts on statistical computing and statistical packages.

Giovanni Cerulli

# Fundamentals of Supervised Machine Learning

With Applications in Python, R, and Stata

Springer

Giovanni Cerulli 🆔
IRCRES-CNR, Research Institute for
Sustainable Economic Growth
National Research Council of Italy
Rome, Italy

*To the memory of my grandmother Esterina and aunt Letizia who silently but constantly accompanied the writing of this book. For the good, the wisdom, and the knowledge they pass on to me. For their priceless and everlasting love.*

# Preface

## Scope of the Book

This is my second book published with Springer. The first, *Econometrics of Program Evaluation: Theory and Applications* (second edition), deals with causal inference for policy impact assessment and has a strong inferential character. This second book, dealing with supervised Machine Learning, largely complements the first in many regards, in that it privileges predictive analytics over inferential analysis. Inference and prediction are however both essential for data science and analysis, as they shed light on two different, although increasingly linked, modes of scientific investigation.

This book presents the fundamental theoretical notions of supervised Machine Learning along with a wide range of applications with Python, R, and Stata software. The book is dedicated to Ph.D. students, academics, and practitioners in various fields of study (social sciences, medicine, epidemiology, as well as hard sciences) who intend to learn the fundamentals of supervised Machine Learning to apply it to concrete case studies.

The book assumes the reader to have a good understanding of basic statistics (both descriptive and inferential), the meaning and the writing of algorithms, and a working knowledge of Python, R, or Stata software. Mathematics is used only when strictly necessary, and a large focus is paid to graphical explanations instead of analytic proofs.

Why should the reader consider this book as a valuable source of knowledge for learning theoretical and applied supervised Machine Learning? Here are three compelling reasons:

1. *Comprehensive coverage of theoretical foundations.* The book starts by presenting the fundamental theoretical notions of supervised Machine Learning. It covers key concepts such as regression, classification, ensemble methods, and evaluation metrics, providing a solid foundation for understanding the principles and techniques behind supervised learning. This theoretical grounding ensures that readers gain a deep understanding of the subject matter before moving on to practical applications.

2. *Diverse applications and case studies.* The book goes beyond theory and offers a wide range of applications using three popular software packages: Python, R, and Stata. By including multiple software options, the book caters to the diverse needs and preferences of readers. Python, known for its versatility and extensive libraries like scikit-learn, allows readers to implement Machine Learning algorithms and explore advanced techniques. R, a widely used language in statistics and data analysis, provides readers with additional tools and packages—such as CARET—specifically tailored for statistical learning modeling and results' visualization. Stata, a statistical software widely used in social sciences, medicine, and epidemiology, offers readers an alternative perspective with its unique features and capabilities. By including all three software packages, this book ensures that readers can leverage the software they are most comfortable with or choose the one best suited to their specific domain.

3. *Differentiating factors among other books.* While there are indeed several good books available on Statistical Learning and software, this book stands out by addressing the specific needs of Ph.D. students, academics, and practitioners in various fields of study. It bridges the gap between theory and application, making it particularly relevant for those who wish to apply supervised Machine Learning techniques to concrete case studies in their respective disciplines. The inclusion of multiple software packages further enhances the book's value, as readers can gain exposure to different tools and methodologies and develop a versatile skill set.

In summary, this book offers a comprehensive introduction to supervised Machine Learning, catering to the needs of Ph.D. students, academics, and practitioners in diverse fields of study. Its inclusion of Python, R, and Stata software packages allows readers to choose the tool that aligns with their preferences and domain-specific requirements. By combining theoretical foundations with practical applications and addressing the specific needs of its target audience, this book provides a unique and valuable resource for anyone looking to master supervised Machine Learning for real-world scenarios.

Every chapter presents an initial theoretical part, where the basics of the methodologies are explained, followed by an applicative part, where the methods are applied to real-world datasets. Each chapter is self-contained, but the reader is invited to consider reading the two introductory chapters before going through the study of chapters dealing with single methods.

For ease of reproducibility, the Python, R, and Stata codes used in the book for carrying out the applications, along with the related datasets, are available at this GitHub link: https://github.com/GioCer73/Book_FSML_Ed1.

## Organization of the Book

The book is organized into eight chapters for a structured and cohesive learning experience.

Chapter 1. The first chapter introduces the rationale and ontology of Machine Learning. Its purpose is to pave the way for the subsequent chapters and facilitate the understanding of the material covered.

Chapter 2. This chapter focuses on the statistics of Machine Learning, which serves as a crucial foundation for comprehending the following chapters. Understanding this section will ensure a smooth transition into the presentation and application of individual Machine Learning methods.

Chapter 3. Model selection and regularization are the main topics discussed in this chapter. Special emphasis is placed on lasso and elastic-net regression and classification, subset selection models, and the lasso's role in inferential analysis.

Chapter 4. The fourth chapter delves into discriminant analysis, nearest neighbor methods, and support vector machines. Although primarily utilized for classification purposes, it is worth noting that regression analysis can also be performed using the nearest neighbor and support vector machine methods.

Chapter 5. In this chapter, tree modeling is explored in both classification and regression scenarios. After outlining the tree-building algorithm, the chapter introduces three ensemble methods: bagging, random forests, and boosting machines, which are valuable for prediction and feature-importance detection.

Chapter 6. The sixth chapter provides an introduction to artificial neural networks. It covers fully connected neural networks, explaining their construction logic and demonstrating software applications for intelligent tasks such as image recognition.

Chapter 7. Building upon the foundation laid in the previous chapter, the seventh chapter delves into deep learning modeling. This section focuses on special artificial neural networks that leverage specific data ordering to enhance predictability and computation efficiency. Two deep learning architectures are presented and discussed: convolutional neural networks and recurrent neural networks.

Chapter 8. The concluding chapter serves as a primer on sentiment analysis. This approach involves developing a predictive mapping between human textual documents and the corresponding human polar sentiment associated with those documents. The chapter explores textual feature engineering and its application in predictive analytics.

By structuring the book into these eight chapters, readers can systematically progress from understanding the fundamentals to exploring advanced topics and applications in Machine Learning.

Rome, Italy                                                        Giovanni Cerulli

# Acknowledgments

# Contents

# List of Figures

# List of Tables