# Machine Learning:

**An Introductory Handbook Using R**

## Dr. Karen Mazidi

# Contents

Machine Learning: An Introductory Handbook Using R ©Karen Mazidi

| II | **Part Two: Linear Models** |
|----|------------------------------|

## IV    Part Four: Neural Networks

| **V** | **Part Five: Modeling the World** |
|-------|-----------------------------------|

# Preface to the Book

I accidentally started writing this book during Spring Break 2018. I was driving down a long Texas road out in the country, thinking about a problem with the undergraduate machine learning class I had been teaching. Chiefly, there was no book suitable for my course. There are several excellent machine learning books for graduate students, but most are too heavy on the math and too light on practical issues for undergraduate students. There are a few books for undergraduate students, but they didn't cover all the topics I was required to teach for my course at The University of Texas at Dallas. The solution I came up with was to format my own notes in latex. I scribbled an outline on a piece of paper as I drove down the road. When I came home I started looking for a suitable latex template and found this one which I really liked. It was designed for a book and I thought, why not just make this a book? And so the book began.

Many years ago I authored several textbooks on microprocessors and microcontrollers with my husband. These books were quite successful and provided a nice second income for our family. I was able to write the books at home while taking care of our two rambunctious sons. However, the days when one can write textbooks as a job are gone. As soon as a book is published, it is pirated, and the pdf is available for free to anyone in the world. That is not entirely a bad thing in my opinion. Yes it is stealing intellectual property from the owner but on the other hand it makes information available to people who want to learn but perhaps can't pay the exhorbitant price that publishing companies charge for textbooks. And so my goals are not monetary. The pdf of this book will remain free, although there is a print version available on Amazon for those who prefer to read from paper. The Amazon print book is printed in grayscale so refer to the pdf when color is important in graphs.

My goals in writing this book are the following:
- To create an undergraduate book that covers the most popular machine learning algorithms in a hands-on approach that reinforces understanding of the algorithms.
- To provide an accessible introduction to the field to professionals who want to get started with machine learning.
- To increase the number of practitioners in the field. I believe that machine learning is to the present and future what computer programming was in the 1970s when I started: a revolutionary new way of solving problems.
- To help people enhance their skill set. If I can introduce people to these skills it will help them in their career goals.

I hope you enjoy the book. This first edition covers what I teach in a one-semester undergraduate Introduction to Machine Learning course. However, I consider it to be a work in progress and I will include additional material in the future, as well as some short video vignettes. Links will be posted on my website: `www.karenmazidi.com` and my blog: `http://karenmazidi.blogspot.com/`

The book has a companion github site for code samples: `https://github.com/kjmazidi/Machine_Learning/`