An Introduction to Data Analysis using
Aggregation Functions in R

Simon James

# An Introduction to Data Analysis using Aggregation Functions in R

Springer

Simon James
School of Information Technology
Deakin University
Burwood, VIC, Australia

# Preface

This introduction to the area of aggregation and data analysis is intended for computer science or business students who are interested in tools for summarizing and interpreting data, but who do not necessarily have a formal mathematics background. It was motivated by our need for such a text in postgraduate data analytics subjects at my university; however it could also be of interest to undergraduate data science students. By reading through the chapters and completing the questions, you will be introduced to some of the issues that can arise when we try to make sense of data. It is hoped that you will gain an appreciation for the usefulness of the results established in the field of data aggregation and their wide applicability. I personally find the study of aggregation functions to be very enjoyable and fulfilling, as it allows for a nice mix of theoretically interesting results and practical applications. The rise of data analytics and the growing need for companies to learn more from their data also make the need for intelligent data aggregation techniques more indispensable than ever.

The overall aim of this book is to allow future data analysts to become aware of aggregation functions theory and methods in an accessible way, focusing on a fundamental understanding of the data and summarization tools that complements any study in statistical or machine learning techniques. To this end, included in each of the chapters is an R tutorial giving an introduction to commands and techniques relevant to the topics studied. These tutorials assume no programming background and will give you some exposure to one of the most widely adopted (and freely available) computing languages used in data analysis.

For a more in-depth overview of aggregation functions, there are some great existing monographs, including the following to name a few.

- Beliakov, G., Bustince, H. and Calvo, T.: A Practical Guide to Averaging Functions. Springer, Berlin, New York (2015)
- Beliakov, G., Pradera, A. and Calvo, T.: Aggregation Functions: A Guide for Practitioners. Springer, Heidelberg (2007)
- Gagolewski, M.: Data Fusion. Theory, Methods and Applications. Institute of Computer Science, Polish Academy of Sciences (2015)

- Grabisch, M., Marichal, J.-L., Mesiar, R. and Pap, E.: Aggregation Functions. Cambridge University press, Cambridge, (2009)
- Torra, V. and Narukawa, Y.: Modeling Decisions. Information Fusion and Aggregation Operators. Springer, Berlin, Heidelberg (2007)

Whilst focusing on established results, I hope that the following pages will offer a fresh perspective on the recent trends in aggregation research. In this respect I am indebted to a number of colleagues and researchers who lead this field, many of whom I have been fortunate enough to meet or collaborate with over the span of my short career. In particular, I would like to acknowledge Humberto Bustince, Radko Mesiar, and Michal Baczynski, who have all made me feel very welcome when visiting their institutions and participating in conferences overseas. Special mention should be made regarding the works of Michel Grabisch and Ronald R. Yager on fuzzy integrals and OWA operators respectively, which have been a huge influence in my understanding of these areas (and so I hope I do them justice in the short introductions provided). I am grateful to my friend and ecology research expert, Dale Nimmo, who has helped identify a number of interesting applications aggregation theory. A huge thank-you to Marek Gagolewski for our recent discussions and for his extensive comments on drafts of this text; I am looking forward to our future research endeavors! Finally (on the academic side), I would like to thank Gleb Beliakov for continuing to be a mentor to my research. Certainly the majority of these chapters have been influenced by what you have taught me and our work together.

Lastly I would like to acknowledge the support of family (in particular my Mum, who proofread some chapters of this book), friends (especially Rachel, who has been a great support while I've been working on this book in addition to her usual honest and constructive comments), and my work colleagues (a very busy 2016 was made much more bearable by having people like Lauren, Elicia, Guy, Tim, Lei, Crystal, Menuri, Gang, Shui, Vicky, Sutharshan, and Michelle to talk to).

Melbourne, VIC, Australia                                                              Simon James
August 2016

# Contents