# DATA ANALYSIS

A gentle introduction for future data scientists

GRAHAM UPTON | DAN BRAWN

# Data Analysis

*A gentle introduction for future data scientists*

## Graham Upton

*Former Professor of Applied Statistics, University of Essex*

## and

## Dan Brawn

*Lecturer, Department of Mathematical Sciences, University of Essex*

# Contents