# Contents

## Chapter 3

# Data Preparation as a Process   89

## Chapter 6
## Handling Nonnumerical Variables 191

## Chapter 9

### Series Variables   299

## Chapter 10

### Preparing the Data Set   351