# Contents

# Chapter 8 The One Row-per-Subject Data Mart 61

# Chapter 9 The Multiple-Rows-per-Subject Data Mart 69

# Chapter 10 Data Structures for Longitudinal Analysis 77

# Chapter 11 Considerations for Data Marts 89

# Chapter 12 Considerations for Predictive Modeling 95

# Part 3    Data Mart Coding and Content

# Part 5   Case Studies