John MacGregor

# Predictive Analysis with SAP®

The Comprehensive Guide

**Galileo Press**

# Contents

## PART III Predictive Analysis Categories

# 13   Classification Analysis—Decision Trees  ...................................  387

## 14 Classification Analysis—K Nearest Neighbor ............................ 427

## PART V  Advanced Predictive Analysis

## 15 Time Series Analysis ................................................................. 441