# CONTENTS

## PART II  Analytics: Modeling Data

## PART IV  Applications: Using Data