# Statistical and Machine-Learning Data Mining

## Techniques for Better Predictive Modeling and Analysis of Big Data

### Second Edition

# Bruce Ratner

# Contents