# Data Mining and Data Visualization

**Edited by**

# C.R. Rao

**Center for Multivariate Analysis**
**Department of Statistics, The Pennsylvania State University**
**University Park, PA, USA**

# E.J. Wegman

**Center for Computational Statistics**
**George Mason University**
**Fairfax, VA, USA**

# J.L. Solka

**Naval Surface Warfare Center, DD**
**Dahlgren, VA, USA**

# Table of contents

## Ch. 16.  Interactive Statistical Graphics: the Paradigm of Linked Views  437
### *Adalbert Wilhelm*

## Ch. 17.  Data Visualization and Virtual Reality  539
### *Jim X. Chen*