Alexander N. Gorban  Balázs Kégl
Donald C. Wunsch   Andrei Zinovyev
(Eds.)

# Principal Manifolds
# for Data Visualization
# and Dimension Reduction

With 82 Figures and 22 Tables

# Contents

## 3 Learning Nonlinear Principal Manifolds
## by Self-Organising Maps

## 4 Elastic Maps and Nets for Approximating Principal
## Manifolds and Their Application to Microarray Data
## Visualization