# Business Intelligence: Data Mining and Optimization for Decision Making

**Carlo Vercellis**

*Politecnico di Milano, Italy.*

## WILEY

A John Wiley and Sons, Ltd., Publication

# Contents