# Applied
# Data Mining

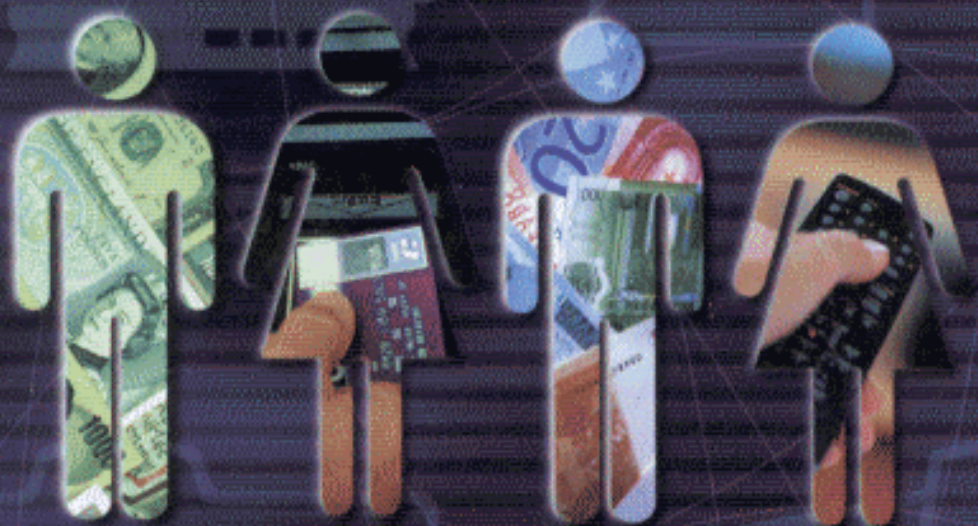## Statistical Methods for Business and Industry

### Paolo Giudici

# Contents