Peter Dalgaard

# Introductory Statistics with R

Second Edition

Springer

# Contents

xvi    Contents