Joaquim P. Marques de Sá

# Applied Statistics
Using SPSS, STATISTICA, MATLAB and R

With 195 Figures and a CD

Springer

# Contents

**Appendix F - Tools**                                                    **487**

**References**                                                            **491**

**Index**                                                                **499**