# The Analysis of Biological Data

## WHITLOCK · SCHLUTER

# Contents

## PART **5**    MODERN STATISTICAL METHODS