# Statistics for Veterinary and Animal Science

## THIRD EDITION

Aviva Petrie & Paul Watson

with website

# Statistics for Veterinary and Animal Science

# Statistics for Veterinary and Animal Science

## Third Edition

### Aviva Petrie, BSc, MSc, CStat, CSci, FHEA

Senior Lecturer in Statistics and Head of the Biostatistics Unit
UCL Eastman Dental Institute
University College London

Honorary Lecturer in Medical Statistics
London School of Hygiene and Tropical Medicine
University of London
London
UK

### Paul Watson, BSc, BVetMed, PhD, DSc, FRCVS

Emeritus Professor of Reproductive Cryobiology
The Royal Veterinary College
University of London
London
UK

# Contents

*Colour plate section can be found facing
page 240*

# Preface to third edition

The continuing interest in our textbook together with the ongoing development of statistical applications in veterinary and animal science has encouraged us to prepare this third edition of *Statistics for Veterinary and Animal Science*. We have introduced some new material but we want to reassure all readers that our original intention of this being an introductory text still stands. Again, you will find everything that you need to begin to understand statistics and its application to your scientific and clinical endeavours; it remains an introduction for the novice with emphasis on understanding the application, rather than exhibiting mathematical competence in the calculations. Readily available statistical software packages, which provide the mechanics of the calculations, have become more extensive in the range of procedures they offer. Accordingly, we have augmented our text, within the bounds of an introductory exposition, to match their development.

As in previous editions, we use two commonly employed statistical software packages, SPSS and Stata, to analyse the data in our examples. We believe that by presenting you with different forms of computer output, you will have the confidence and proficiency to interpret output from other statistical packages. The previous edition of the book had an accompanying CD which contained the data sets (in ASCII, Excel, SPSS and Stata) used as examples in the text. These data sets are now available at www.wiley.com/go/petrie/statisticsforvets, and will be helpful if you wish to get to grips with various statistical techniques by attempting the analyses yourselves. You will find a website icon next to the examples for which the data are available on the website.

Please note that, although we have provided details of a considerable number of websites that you may find useful, we cannot guarantee that these website addresses will remain correct over the course of time because of the mutability of the internet.

Some sections of the book are, as in previous editions, in small print and are accompanied by a jumping horse symbol. These sections contain information that relates to more advanced or obscure topics, and you may skip (jump over) them without loss of continuity. Our teaching experience has demonstrated that one of hardest tasks for the novice when analysing his or her own data set is deciding which test or procedure is most appropriate. To overcome this difficulty, we provide two flow charts (Figure E.2 for binary data and Figure E.3 for numerical data) which lead you through the various questions that need to be asked to aid that decision. Another flow chart (Figure E.1) organizes the tests and procedures into relevant groups and indicates the particular section of the book where each is located: you can find these flow charts in the Appendix as well as on the inside back/front covers for easy reference.

Many of the chapters in this third edition are similar to those in the second edition, apart from some minor modifications and additional exercises. However, Chapter 5 has been extended to include techniques for recognizing and dealing with confounding, and this chapter now provides a description of the different types of missing data that might be encountered. We have added a section on checking the assumptions underlying a logistic regression model to Chapter 11, and have included modifications of the sample

size estimation process to take account of different group sizes and losses to follow-up in Chapter 13. Chapter 14 has been expanded considerably by extending the sections on diagnostic tests, measuring agreement and survival analysis as well as Bayesian analysis. Chapter 15 is entirely new, bringing together a group of specialist topics – ethical issues of animal investigation (some of which was in Chapter 5 of the second edition), spatial statistics, surveillance and its importance, and statistics in molecular and quantitative genetics. While none of these is intended as more than an introduction, you will find references to help you explore the topics more fully should you so desire. The section on evidence-based veterinary medicine (EBVM) in Chapter 16 is unchanged from that in the second edition's Chapter 15, but in the third edition this chapter no longer provides guidelines for reporting results. Instead, we have devoted the new Chapter 17 to this topic by presenting different published guidelines relevant to veterinary medicine (i.e. for reporting of livestock trials, research using laboratory animals, diagnostic accuracy studies, observational studies in epidemiology, and systematic reviews and meta-analyses) as a ready reference for those wanting to follow best practice both in planning and in writing up their research. Lastly, in Chapter 18, which is entirely new, we bring together the concepts of EBVM and the guidelines provided in Chapter 17 by proffering a template for the critical appraisal of randomized controlled trials and observational studies. We use this template to critically appraise two published papers, both of which are reproduced in full, and hope that by providing these examples, we will help you develop your own skills in what is an essential, but frequently overlooked, component of statistics.

We are indebted as always to those who, for earlier editions of this book, have offered their data to us to use for examples or exercises, have assisted with the presentation of the illustrations and tables, and have provided critical advice on the text. These colleagues are all identified in the prefaces to the first and second editions. As in earlier editions, we have occasionally taken summary data or abstracts from published papers and have used them to develop exercises or to illustrate techniques: we extend our thanks to the authors and the publishers for the use of this material. For this third edition, we are most grateful to Dr Geoff Pollott and Professor Dirk Pfeiffer (both of the Royal Veterinary College, University of London) for their critical reading and suggestions for sections of Chapter 15. We wish to record our particular thanks to Professor Garry Anderson (University of Melbourne) for his critique of much of the new text. His suggestions have drawn our attention to errors and have considerably improved the presentation. Nonetheless, we remain responsible for all contained herein, and offer it, with all its shortcomings, to our readership.

This preface would not be complete without acknowledging our marriage partners, Gerald and Rosie, and our children, Nina, Andrew and Karen, and Oliver and Anna, who have allowed us once again to engage with this task to their inevitable exclusion, and offer them our most grateful thanks.

*Aviva Petrie*
*Paul Watson*
*2013*

# Preface to second edition

It is six years since this book was first available, and we are glad to acknowledge the positive responses we have received to the first edition and the evident uptake of the text for a number of courses around the world. In the intervening period much has happened to encourage us to update and expand our initial text. However, many of the chapters which were in the first edition of the book are changed only slightly, if at all, in this second edition. To these chapters, we have added some exercises and further explanations (for example, on equivalence studies, confounding, interactions and bias, Bayesian analysis and Cox survival analysis) to make the book more comprehensive. We have nevertheless retained our original intent of this being an introductory text starting with very basic concepts for the complete novice in statistics. You will still find sections marked for skipping unless you have a particular need to explore them, and these include the newer more complex analysis methods. This edition also contains the glossaries of notation and of terms, but we have expanded them to reflect the enhanced content of the text. For easy reference, the flow charts for choosing the correct statistical analyses in different situations are now found immediately before the index, and we hope these will serve to guide you to the appropriate procedures and text relating to their use.

Computer software to deal with increasingly sophisticated analytical tools has been developed in recent years in such a way that the associated methodology is more readily accessible to those who previously believed such techniques were out of their reach. As a consequence, we have substantially enhanced the material relating to regression analysis and created a new chapter (Chapter 11) to describe some advanced regression techniques. The latter incorporates the sections on multiple regression and an expanded section on logistic regression from Chapter 10 of the first edition, and introduces Poisson regression, different regression methods which can be used to analyse clustered data, maximum likelihood estimation and the concept of the generalized linear model. Because we have inserted this new Chapter 11, the numbering of the chapters which follow does not accord with that of the corresponding chapters in the first edition.

Chapter 15 is an entirely new chapter which is devoted in large part to introducing the concepts of evidence-based veterinary medicine (EBVM), stressing the role of statistical knowledge as a basis for its practice. The methodology of EBVM describes the processes for integrating, in a systematic way, the results of scientifically conducted studies into day-to-day clinical practice with the aim of improving clinical outcome. This requires the practitioner to develop the skills to evaluate critically the efforts of others in respect of the design of studies, and of the presentation, analysis and interpretation of results. The recognition of the value of the evidence-based approach to veterinary medicine has followed a similar emphasis in human clinical medicine, and is influencing the whole veterinary profession. Accordingly, it is also very much a part of the mainstream veterinary curriculum. Whether you are a practitioner of veterinary medicine or of one of the allied sciences, you will now more than ever need to be conversant with modern biostatistical analysis. Knowing how best to report your own results is also vital if you are to impart

knowledge correctly, and so, to this end, we include in Chapter 15 a section on the CONSORT Statement, designed to standardize clinical trial reporting.

Although we refer only to two common statistical packages in the text, SPSS and Stata, sufficient information is given to interpret output from other packages, even though the layout and content may differ to some degree. We have also mentioned a number of websites containing useful information, and which were correct at the time of printing. Given the mutability of the internet, we cannot guarantee that such sites will stay available.

Also included with this edition is a CD containing the data sets used as examples in the text. You can use these data sets to consolidate the learning process. It is only when you attempt the analyses yourself that you are fully able to get to grips with the techniques. Each data set is presented in four different formats (ASCII, Excel, SPSS and Stata), so you should be able to access the data and use the software that is available to you.

We would like to acknowledge the generosity of the late Dr Penny Barber, Mark Corbett, Dr J. E. Edwards, Professor Jonathan Elliott, Professor Gary England, Dr Oliver Garden, Dr Ilke Klaas, Dr Teresa Martinez, Dr Anne Pearson, Dr P. D. Warriss, Professor Avril Waterman-Pearson and Dr Susannah Williams who shared their original data with us, and to others who have allowed us to use their published data. In places, we have taken published summary data and constructed a primary data set to suit our own purposes; if we have misrepresented our colleagues' data, we accept full responsibility. We are particularly grateful to Alex Hunte who lent us his skills in refining the illustrations in the first edition, and to Dr David Moles who assisted with the preparation of the statistical tables. We especially thank Dr Ben Armstrong, Professor Caroline Sabin and Dr Ian Martin who kindly gave us their critical advice as the text of the first edition was developed, and Professor John Smith who was instrumental in getting us to consider writing the book in the first place. In addition, we acknowledge our debt to a host of other colleagues who have helped with discussions over the telephone, with their expertise in areas we are lacking, and in their encouragement to complete what we hope will be a useful contribution to the field of veterinary and animal science. We are particularly indebted to those of our colleagues who have graciously pointed us to our errors, which we hope are now corrected.

Lastly, we again acknowledge with gratitude the patience and encouragement of our marriage partners, Gerald and Rosie, and our children, Nina, Andrew and Karen, and Oliver and Anna, who have once more graciously allowed us to become absorbed in the book and have had to suffer neglect in the process. We trust that they still appreciate the worthiness of the cause!

*Aviva Petrie*
*Paul Watson*

# Preface to first edition

Although statistics is anathema to many, it is, unquestionably, an essential tool for those involved in animal health and veterinary science. It is imperative that practitioners and research workers alike keep abreast with reports on animal production, new and emerging diseases, risk factors for disease and the efficacy of the ever-increasing number of innovations in veterinary care and of developments in training methods and performance. The most cogent information is usually contained in the appropriate journals; however, the usefulness of these journals relies on the reader having a proper understanding of the statistical methodology underlying study design and data analysis. The modern animal scientist and veterinary surgeon therefore need to be able to handle numerical data confidently and properly. Often, for us, as teachers, there is little time in busy curricula to introduce the subject slowly and systematically; students find they are left bewildered and dejected because the concepts seem too difficult to grasp. While there are many excellent introductory books on medical statistics and on statistics in other disciplines such as economics, business studies and engineering, these books are unrelated to the world of animal science and health, and students soon lose heart. It is our intention to provide a guide to statistics relevant to the study of animal health and disease. In order to illustrate the principles and methods, the reader will find that the text is well endowed with real examples drawn from companion and agricultural animals. Although veterinary epidemiology is closely allied to statistics, we have concentrated only on statistical issues as we feel that this is an area which, until now, has been neglected in veterinary and animal health sciences.

Our book is an introductory text on statistics. We start from very simple concepts, assuming no previous knowledge of statistics, and endeavour to build up an understanding in such a way that progression on to advanced texts is possible. We intend the book to be useful for those without mathematical expertise but with the ability to utilize simple formulae. We recognize the influence of the computer and so we avoid the description of complex hand calculations. Instead, emphasis is placed on understanding of concepts and interpretation of results, often in the context of computer output. In addition to acquiring an ability to perform simple statistical techniques on original data, the reader will be able critically to evaluate the efforts of others in respect of the design of studies, and of the presentation, analysis and interpretation of results. The book can be used either as a self-instructional text or as a basis for courses in statistics. In addition, those who are further on in their studies will be able to use the text as a reference guide to the analysis of their data, whether they be postgraduate students, veterinary practitioners or animal scientists in various other settings. Every section contains sufficient cross referencing for the reader to find the relevant background to the topic.

We are particularly grateful to Alex Hunte who lent us his skills in preparing the illustrations, and to Dr David Moles who assisted with the preparation of the statistical tables. We especially thank Dr Ben Armstrong, Dr Caroline Sabin and Dr Ian Martin who kindly gave us their critical advice as the text was developed. Professor John Smith was instrumental in getting us to consider writing the text in the first place, and we thank him for his continual encouragement. In addition, we acknowledge our debt to a host of other colleagues who have helped with discussions over the telephone, with their expertise in areas we are lacking, and in general encouragement to complete what we hope will be a useful contribution to the field of veterinary and animal science.

Lastly, we acknowledge with gratitude the patience and encouragement of our families. Our marriage partners, Gerald and Rosie, have endured with fortitude our neglect of them while this work was in preparation. In particular, our children, Nina, Andrew and Karen, and Oliver and Anna, have had to cope with our absorption with the project and lack of involvement in their activities. We trust they will recognize that it was in a good cause.

*Aviva Petrie*
*Paul Watson*

# About the companion website

This book is accompanied by a companion website:
**www.wiley.com/go/petrie/statisticsforvets**

The website includes:
- Data files which relate to some of the examples in the text. Each data file is provided for download in four different formats: ASCII, Excel, SPSS and Stata.
- Examples relating to the data files are indicated in the text using the following icon:

# 1 The whys and wherefores of statistics

## 1.1 Learning objectives

By the end of this chapter, you should be able to:

- State what is meant by the term 'statistics'.
- Explain the importance of a statistical understanding to the animal scientist.
- Distinguish between a qualitative/categorical and a quantitative/numerical variable.
- List the types of scales on which variables are measured.
- Explain what is meant by the term 'biological variation'.
- Define the terms 'systematic error' and 'random error', and give examples of circumstances in which they may occur.
- Distinguish between precision and accuracy.
- Define the terms 'population' and 'sample', and provide examples of real (finite) and hypothetical (infinite) populations.
- Summarize the differences between descriptive and inferential statistics.

## 1.2 Aims of the book

### 1.2.1 What will you get from this book?

All the biological sciences have moved on from simple qualitative description to concepts founded on numerical measurements and counts. The proper handling of these values, leading to a correct understanding of the phenomena, is encompassed by statistics. This book will help you appreciate how the theory of statistics can be useful to you in veterinary and animal science. Statistical techniques are an essential part of communicating information about health and disease of animals, and their agricultural productivity, or value as pets, or in the sporting or working environment. We, the authors, aim to introduce you to the subject of statistics, giving you a sound basis for managing straightforward study design and analysis. Where necessary, we recommend that you extend your knowledge by reference to more specialized texts. Occasionally, we advocate that you seek expert statistical advice to guide you through particularly tricky aspects.

You can use this book in two ways:

1. The chapter sequence is designed to develop your understanding systematically and we therefore recommend that, initially, you work through the chapters in order. You will find certain sections marked in small type with a symbol, which indicates that you can skip these, at a first read through, without subsequent loss of continuity. These marked sections contain information you will find useful as your knowledge develops. Chapters 11, 14 and 15 deal with particular types of analyses which, depending on your areas of interest, you may rarely need.
2. When you are more familiar with the concepts, you can use the book as a reference manual;

you will find sufficient cross-referenced information in any section to answer specific queries.

## 1.2.2 What are learning objectives?

Each chapter has a set of **learning objectives** at the beginning. These set out in task-oriented terms what you should be able to 'do' when you have mastered the concepts in the chapter. You can therefore test your growing understanding; if you are able to perform the tasks in the learning objectives, you have understood the concepts.

## 1.2.3 Should you use a computer statistics package?

We encourage you to use available computer statistics packages, and therefore we do not dwell on the development of the equations on which the analyses are based. We do, however, present the equations (apart from when they are very complex) for completeness, but you will normally not need to become familiar with them since computer packages will provide an automatic solution. We provide computer output, produced when we analyse the data in the examples, from two statistical packages, mostly from SPSS (IBM SPSS Version 20 (www-01.ibm.com/software/analytics/spss, accessed 9 October 2012)) and occasionally from Stata (Stata 12, StataCorp, 2011, *Stata Statistical Software: Release 12*. College Station, TX: StataCorp LP (www.stata.com/products, accessed 9 October 2012)). Although the layout of the output is particular to each individual package, from our description you should be able to make sense of the output from any other major statistical package.

## 1.2.4 Will you be able to decide when and how to use a particular procedure?

Our main concern is with the *understanding* that underlies statistical analyses. This will prevent you falling into the pitfalls of misuse that sur-

round the unwitting user of statistical packages. We present the subject in a form that we hope is accessible, using examples showing the application of the subject to veterinary and animal science. A brief set of exercises is provided at the end of each chapter, based on the ideas presented within. These exercises should be used to check your understanding of the concepts and procedures; solutions to the exercises are given at the back of the book. The two exceptions are Chapter 17, which provides reporting guidelines and Chapter 18 in which we ask you to critically appraise two published articles, preferably before looking at the 'model answers' provided in the chapter.

## 1.2.5 Use of the glossaries of notation and terms

Statistical nomenclature is often difficult to remember. We have gathered the most common symbols and equations used throughout this book into a Glossary of notation in Appendix C. This gives you a readily accessible reminder of the meaning of the terminology.

You will find a Glossary of terms in Appendix D. In this glossary, we define common statistical terms which are used in this book. They are also defined at the appropriate places in relevant chapters, but the glossary provides you with a ready reference if you forget the meaning of a term. Terms that are in the glossary are introduced in the text in bold type. Note, however, that there are some instances where bold is purely used for extra emphasis.

## 1.3 What is statistics?

The number of introductory or elementary texts on the subject of statistics indicates how important the subject has become for everyone in the biological sciences. However, the fact that there are many texts might also suggest that we have yet to discover a foolproof method of presenting what is required.

The problem confronted in biological statistics is as follows. When you make a set of numerical

observations in biology, you will usually find that the values are scattered. You need to know whether the values differ because of factors you are interested in (e.g. treatments) or because they are part of a 'background' natural variation. You need to evaluate what the numbers actually mean, and to represent them in a way that readily communicates their meaning to others.

The subject of statistics embraces:

- The design of the study in order that it will reveal the most information efficiently.
- The collection of the data.
- The analysis of the data.
- The presentation of suitably summarized information, often in a graphical or tabular form.
- The interpretation of the analyses in a manner that communicates the findings accurately.

Strictly, this broad numerical approach to biology is correctly termed '**biometry**' but we shall adopt the more generally used term '**statistics**' to cover all aspects. Statistics (meaning this entire process) has become one of the essential tools in modern biology.

## 1.4 Statistics in veterinary and animal science

One of the common initial responses of both veterinary students and animal science students is: Why do I need to study statistics? The mathematical basis of the subject causes much uncertainty, and the analytical approach is alien. However, in professional life, there are many instances of the relevance of statistics:

- The **published scientific literature** is full of studies in which statistical procedures are employed. Look in any of the relevant scientific journals and notice the number of times reference is made to mean ± SEM (standard error of mean), to statistical significance, to $P$-values or to $t$-tests or Chi-squared analysis or analysis of variance or multiple regression analysis. The information is presented in the usual brief form and, without a working knowledge of statistics, you are left to accept the conclusions of the author, unable to examine the strength of the supporting data. Indeed, with the advent of computer-assisted data handling, many practitioners can now collect their own observations and summarize them for the advantage of their colleagues; to do this, they need the benefit of statistical insights.

- The subject of **epidemiology** (see Section 5.2) is gaining prominence in veterinary and animal science, and the concepts of **evidence-based veterinary medicine** (see Section 1.5 and Chapter 16) are being explicitly introduced into clinical practice. As never before, there is an essential need for you to understand the types of trials and investigations that are carried out and to know the meaning of the terms associated with them.

- In the animal health sciences, there are an increasing number of independent **diagnostic services** that will analyse samples for the benefit of health monitoring and maintenance. Those running such laboratory services must always be concerned about quality control and accuracy in measurements made for diagnostic purposes, and must be able to supply clear guidelines for the interpretation of results obtained in their laboratories.

- The **pharmaceutical and agrochemical industries** are required to demonstrate both the safety and the efficacy of their products in an indisputable manner. Such data invariably require a statistical approach to establish and illustrate the basis of the claim for both these aspects. Those involved in pharmaceutical product development need to understand the importance of study design and to ensure the adequacy of the numbers of animals used in treatment groups in order to perform meaningful experiments. Veterinary product licensing committees require a thorough understanding of statistical science so that they can appreciate the data presented to substantiate the claims for a novel therapeutic substance. Finally, practitioners and animal carers are faced with the blandishments of sales representatives with competing claims, and must evaluate the literature which is offered in support of specific agents, from licensed drugs to animal nutrition supplements.

- Increasingly, there is concern about the regulation of **safety and quality of food for human consumption**. Where products of animal origin are involved, the animal scientist and the veterinary profession are at the forefront. Examples are: pharmaceutical product withdrawal times before slaughter based on the pharmacokinetics and pharmacodynamics of the products, the withholding times for milk after therapeutic treatment of the animal, tissue residues of herbicides and insecticides, and the possible contamination of carcasses by antibiotic-resistant bacteria. In every case, advice and appropriate regulations are established by experimental studies and statistical evaluation. The experts need to be aware of the appropriate statistical procedures in order to play their proper roles.

In all these areas, a common basic vocabulary and understanding of biometrical concepts is assumed to enable scientists to communicate accurately with one another. It is important that you gain mastery of these concepts if you are to play a full part in your chosen profession.

## 1.5 Evidence-based veterinary medicine

The veterinary profession is following the medical profession in introducing a more objective basis to its practice. Under the term **evidence-based veterinary medicine** (EBVM) – by which we mean the conscientious, explicit and judicious use of current best evidence to inform clinical judgements and decision-making in veterinary care (see Cockcroft and Holmes, 2003) – we are now seeing a move towards dependence upon good scientific studies to underpin clinical decisions. In many ways, practice has implicitly been about using clinical experience to make the best decisions, but what has changed is the explicit use of the accessible information. No longer do clinicians have to depend on their own clinical experience and judgement alone; now they can benefit from other studies in a formalized manner to assist their work. The clinician has to know what information is relevant and how to access this

evidence, and be able to use rigorous methods to assess it. Generally, this requires a familiarity with the terminology used and an understanding of the principles of statistical analysis. Moreover, the wider world of animal science is finding a need to understand these ideas as the evidence-based concepts are being applied not only in the treatment of clinical disease but also in aspects of production and performance.

One of the differences between the application of EBVM in veterinary science and in human medicine is that in the latter the body of literature is now very large, and this makes finding relevant information easier. In the veterinary field, EBVM is still hampered by the relatively small amount and variable quality of the evidence available. Nevertheless, EBVM is gaining momentum, and we have devoted Chapter 16 to its concepts. One of the key requirements of EBVM is reliably reported information and, as in the human medical field, the veterinary publishing field is in the process of consolidating a set of guidelines for good reporting. We have addressed this in Chapter 17, outlining the information that is available at the time of writing. As critical appraisal of the published literature is invariably an essential component of evaluating evidence, we have devoted Chapter 18 to it. In this chapter, we provide templates for critically appraising randomized controlled trials and observational studies, and invite you to develop your skills by critically appraising two published articles.

## 1.6 Types of variable

A **variable** is a characteristic that can take values which *vary* from individual to individual or group to group, e.g. height, weight, litter size, blood count, enzyme activity, coat colour, percentage of the flock which are pregnant, etc. Clearly some of these are more readily quantifiable than others. For some variables, we can assign a number to a category and so create the appearance of a numerical scale, but others have a true numerical scale on which the values lie. We take **readings** of the variable which are measurements of a biological characteristic, and these become

the **values** which we use for the statistical procedures. Both these terms are in general use, and both refer to the original measurements, the **raw data**.

Numerical data take various forms; a proper understanding of the nature of the data and the classification of variables is an important first step in choosing an appropriate statistical approach. The flow charts shown in Appendix E, and on the inside front and back covers, illustrate this train of thought, which culminates in a suitable choice of statistical procedure to analyse a particular data set.

We distinguish the main types of variable in a systematic manner by determining whether the variable can take 'one of two distinct values', 'one of several distinct values' or 'any value' within the given range. In particular, the variable may be one of the following:

1. **Categorical (qualitative) variable** – an individual belongs to any one of two or more distinct categories for this variable. A *binary* or dichotomous variable is a particular type of categorical variable defined by only *two* categories; for example, pregnant or non-pregnant, male or female. We customarily summarize the information for the categorical variable by determining the number and percentage (or proportion) of individuals in each category in the sample or population. Particular scales of a categorical variable are:
   - **Nominal scale** – the distinct categories that define the variable are unordered and each can be assigned a name, e.g. coat colours (piebald, roan or grey).
   - **Ordinal scale** – the categories that constitute the variable have some intrinsic order; for example, body condition scores, subjective intensity of fluorescence of cells in the fluorescence microscope, degree of vigour of motility of a semen sample. These 'scales' are often given numerical values 1 to *n*.
2. **Numerical (quantitative) variable** – consisting of numerical values on a well-defined scale, which may be:
   - **Discrete (discontinuous) scale**, i.e. data can take only particular integer values, typically counts, e.g. litter size, clutch size, parity (number of pregnancies within an animal).
   - **Continuous scale**, for which all values are theoretically possible (perhaps limited by an upper and/or lower boundary), e.g. height, weight, speed, concentration of a chemical constituent of the blood or urine. Theoretically, the number of values that the continuous variable can take is infinite since the scale is a continuum. In practice, continuous data are restricted by the degree of accuracy of the measurement process. By definition, the interval between two adjacent points on the scale is of the same magnitude as the interval between two other adjacent points, e.g. the interval on a temperature scale between 37°C and 38°C is the same as the interval between 39°C and 40°C.

## 1.7 Variations in measurements

It is well known that if we repeatedly observe and quantify a particular biological phenomenon, the measurements will rarely be identical. Part of the variability is due to an inherent variation in the biological material being measured. For example, not all cows eat the same quantity of grass per day even if differences both in body weight and water content of the feed are taken into account. We shall use the term '**biological variation**' for this phenomenon, although some people use the term 'biological error'. (Biological error is actually a misleading term since the variability is not in any sense due to a mistake.)

By the selection of individuals according to certain characteristics in advance of the collection of data, we may be able to *reduce* the range of biological variation but we cannot eliminate it. Selection is often based on animal characteristics (e.g. species, strain, age, sex, degree of maturity, body weight, show-jumpers, milking herds, hill sheep, etc.), the choice of which depends upon the particular factors under investigation. However, the result is then only valid for that *restricted population* and we are not justified in extrapolating beyond that population. For example, we should not assume that a study

based on beef cattle applies to other types of cattle.

In addition to biological variation, there will most likely be differences in repeated measurements of the same subject within a very short period of time. These are **technical variations or errors**, due to a variety of instrumental causes and to human error. We may properly consider them to be errors since they represent departures from the true values.

## 1.7.1 Biological variation

The causes of biological variation, which makes one individual differ from the next or from one time to another, may be obvious or subtle. For example, variations in any characteristic may be attributable to:

- Genetics – e.g. greater variability in the whole cow population compared with just Friesians.
- Environment – e.g. body weight varies with diet, housing, intercurrent disease, etc.
- Gender – sexual dimorphism is common.
- Age – many biological data are influenced by age and maturity, e.g. the quantity of body fat.

In a heterogeneous population, the biological variation may be considerable and may mask the variation due to particular factors under investigation. Statistical approaches must take account of this inherent variability. The problem for the scientist, having measured a range of results of a particular feature in a group of individuals, is to distinguish between the sources of variation.

Here are two examples of problems created by biological variation:

- Two groups of growing cattle have been fed different diets. The ranges of the recorded weights at 6 months of age show an overlap in the two groups. Is there a real difference between the groups?
- You have the results of an electrolyte blood test which shows that the serum potassium level is elevated. By how much must it be elevated before you regard it as abnormal?

## 1.7.2 Technical errors

A technical or measurement error is defined as the difference between an observed reading and its 'true' value. Measurement errors are due to factors which are, typically, **human** (e.g. variations within and between observers) or **instrumental**, but may also be attributed to differences in conditions (e.g. different laboratories).

Technical errors may be systematic or random. A **systematic error** is one in which the observed values have a tendency to be above (or below) the true value; the result is then said to be *biased*. When the observed values are evenly distributed above and below the true value, **random errors**, due to unexplained sources, are said to be occurring. Random variation can be so great as to obscure differences between groups but this problem may be minimized by taking repeated observations.

### (a) Human error

Human error can occur whenever a person is performing either an unfamiliar task or a routine or monotonous task; fatigue increases the chances of error. Errors due to these factors are usually random, and providing steps are taken to minimize them (e.g. practice to acquire a proper level of skill, avoiding long periods of monotonous labour, and checking results as measurements are made), they are generally not of great concern.

Other sorts of human error can arise because of data handling. **Rounding errors** can introduce inaccuracies if performed too early in an analysis. If you use a computer to manage your data, you need not be concerned about this, since computer algorithms generally avoid rounding errors by carrying long number strings even if these are not displayed.

Another recognized human error is called **digit preference**. Whenever there is an element of judgement involved in making readings from instruments (as in determining the last digit of a number on a scale), certain digits between 0 and 9 are more commonly chosen than others to represent the readings; such preferences differ

between individuals. This may introduce either a random or a systematic error, the magnitude of which will depend on the importance of the last digit to the results.

## (b) Instrumental error

Instrumental errors arise for a number of reasons (Figure 1.1). Providing we are aware of the potential problem, the causes are often correctable or reducible.

- With a *systematic offset* or *zero error*, a 'blank' sample consistently reads other than zero. It is common in colorimetry and radioisotope measurements (Figure 1.1a).
- *Non-linearity* is a systematic error, commonly seen in the performance of strain gauges, thermocouples and colorimeters (Figure 1.1b).
- *Proportional* or *scale error* is usually due to electronic gain being incorrectly adjusted or altered after calibration; it results in a systematic error (Figure 1.1c).
- *Hysteresis* is a systematic error commonly encountered in measurements involving galvanometers. It may require a standard measurement procedure, e.g. always adjusting input *down* to desired level (Figure 1.1d).
- *Instability* or *drift* – electronic gain calibration may drift with temperature and humidity giving rise to an intermittent but systematic error, resulting in an unstable baseline (Figure 1.1e).
- *Random errors* are commonly seen in attempts to measure with a sensitivity beyond the limits of resolution of an instrument (Figure 1.1f). Most instruments carry a specification of their accuracy, for example it is no use attempting to measure to the nearest gram with a balance accurate only to 10 g.

Two or more of these sources of error may occur simultaneously. Technical errors of all kinds can be minimized by careful experimentation. This is the essence of **quality control** and is of paramount importance in a diagnostic laboratory. Quality control in the laboratory is about ensuring that processes and procedures are carried out in a consistently satisfactory manner

so that the results are trustworthy. We introduce some additional terms in order to understand these concepts more fully.

## 1.8 Terms relating to measurement quality

Two terms which are of major importance in understanding the principles of biological measurement are **precision** and **accuracy**. It is essential they are understood early in a consideration of the nature of data measurement.

- **Precision** refers to how well repeated observations agree with one another.
- **Accuracy** refers to how well the observed value agrees with the true value.

To understand these terms consider the diagrams in Figure 1.2, in which the bull's-eye represents the true value: in Figure 1.2a there is poor accuracy and poor precision, in Figure 1.2b there is poor accuracy and good precision, while in Figure 1.2c there is both good accuracy and good precision.

It is possible to have a diagnostic method (e.g. blood enzyme estimation) that gives good precision but poor accuracy (Figure 1.2b) because of systematic error. In an enzyme activity estimation, such an error might be due to variation in temperature.

Several other terms, all of which describe aspects of **reliability**, are in use and these are defined as follows:

- **Repeatability** is concerned with gauging the similarity of replicate, often duplicate, measurements of a particular technique or instrument or observer under identical conditions, e.g. measurements made by the same observer in the same laboratory. It assesses technical errors (see Section 14.4).
- **Reproducibility** (sometimes called **method agreement**) is concerned with determining how well two or more approaches to measuring the same quantity agree with one another, e.g. measurements made by the same observer but using different methods, or by different

(a) Zero error

(b) Non-linearity

(c) Scale error

(d) Hysteresis

(e) Instability

(f) Random error

**Figure 1.1** Types of instrumental error. 'Input' refers to the true value of the measurements being recorded, 'output' refers to the recorded response, and the solid line refers to the situation when the output values equal the input values. Errors in measurements are represented by dots or dashed lines.

(a)          (b)          (c)

**Figure 1.2** Diagram representing the concepts of accuracy and precision: (a) represents poor accuracy and precision, (b) represents poor accuracy but good precision, and (c) represents both good accuracy and precision.

observers using the same method, or by observers using the same method but in different laboratories (see Section 14.4).

- **Stability** concerns the long-term repeatability of measurement. Diagnostic laboratories will usually have reference material kept for checking stability over time.
- **Validity** is concerned with determining whether the measurement is actually measuring what it purports to be measuring. In the clinical context, the measurement is compared with a 'gold standard' (see Section 14.2).

## 1.9 Populations and samples

The concept of a **population** from which our measurements are a **sample** is fundamental. A population includes all representatives of a particular group, whereas a sample is a subgroup drawn from the population. We aim to choose a sufficiently large sample in such a manner that it is representative (i.e. is typical) of the population (see Sections 1.9.2, 4.2 and 13.3).

### 1.9.1 Types of population

In this book we usually use the word 'animal' to suggest the unit of investigation, but we also use other terms such as 'individual' or 'case'. We want you to become familiar with different terminology. A population of animals may be represented by:

- The individuals, e.g. all cattle, all beef cattle, all Herefords, all the herd.
- The measurements of a particular variable on every animal, e.g. liver weight, bone length, blood hormone or enzyme level.
- Numbers of items (in a given area, volume or time), e.g. blood cell counts or faecal egg counts, counts of radioactive particle emissions.

The population may be either a **real (or finite) group** or a **hypothetical (or infinite) group**. For example, if we are interested in the growth rate of pigs in Suffolk, then the population is all pigs in Suffolk. This is a real or finite population. If,

however, we want to know the effect of an experimental diet in these pigs, we will feed the test diet to a sample of pigs which now comprises the only representatives of a hypothetical population fed on the test diet. Theoretically, at least, we could actually measure the entire population in finite cases, but infinite populations are represented *only* by the sample.

### 1.9.2 Random sampling and random allocation

We examine a sample with a view to making statements about the population. The sample must therefore be *representative* of the population from which it is taken if it is to give useful results applicable to the population at large. In order for the sample to be representative, strictly, there should be **random selection** from *all* possible members of the entire population, implying that the individuals should be selected using a method based on chance (see Section 13.6). However, in reality, random selection is generally not feasible (for example, in an observational study (see Section 5.2.1) or in a clinical study when the disease under investigation is rare). In that case, it is important that we try to ensure that the individuals in the sample are a true reflection of those in the population of interest, and that, if groups are to be compared, we check that the individuals in the different groups are comparable with similar baseline characteristics.

It is essential to use an objective method to achieve random sampling, and a method based on a random number sequence is the method of choice. The sequence may be obtained from a table of random numbers (see Table A.11) or be generated by a computer random number generator or, if only a small sequence, it could be generated by a mechanical method such as rolling a die, although the latter approach is not recommended.

Note that for allocating individuals into treatment groups in an experimental situation, principles of **random allocation** (**randomization**) should also be employed to avoid subjective influence and ensure that the groups are comparable (see Section 5.6). Again, a random number

sequence is recommended to provide objective allocation of individuals or treatments so that the causes of any subsequent differences in performance between the groups can be properly identified.

## 1.10 Types of statistical procedures

Statistical procedures can be divided into descriptive statistics and inferential statistics.

- **Descriptive statistics**. We use these techniques to reduce a data set to manageable proportions, summarizing the trends and tendencies within it, in order to represent the results clearly. From these procedures we can produce diagrams, tables and numerical descriptors. Numerical descriptors include measures that convey where the centre of the data set lies, like the arithmetic mean or median, and measures of the scatter or dispersion of the data, such as the variance or range. These are described more fully in Chapter 2.
- **Inferential statistics**. Statistical inference is the process of generalizing from the sample to the population: it enables us to draw conclusions about certain features of the population when only a subset of it, the sample, is available for investigation. One aspect of inferential statistics is the **estimation** of population parameters using sample data. A parameter, such as the mean or proportion, describes a particular feature of the distribution of a variable in the entire population (see Section 4.3.2). Usually, estimation is followed by a procedure called **hypothesis testing**, another aspect of inferential statistics that investigates a particular theory about the data. Hypothesis tests allow conclusions relating to the population to be drawn from the information in a sample. You can only use these tests properly, and so avoid the pitfalls of misinterpretation of the data, when you have a knowledge of their inherent assumptions. Some of these techniques are simple and require little expertise to master, while others are complex and are best left to the qualified statistician. Details of these procedures can be found in Chapters 6–14; the

flow charts in Appendix E provide a quick guide to the choice of the correct test.

## 1.11 Conclusion

We develop the ideas presented in this chapter in subsequent chapters. As we have said, the concepts are introduced building on one another, and you will need a sound understanding of the earlier theory in order to appreciate the material presented later.

The best incentive for wrestling with statistical concepts is the need to know the meaning of a data set of your own. Remember – statistical procedures cannot enhance poor data. Providing the data have been acquired with sufficient care and in sufficient number, the statistical procedures can supply you with sound summary statements and interpretative guidelines; the interpretation is still down to you! In the chapters that follow, the emphasis is on developing your understanding of the procedures and their limitations to aid your interpretation. We hope you find the experience of getting to grips with your data rewarding, and discover that statistics can be both satisfying and fun!

## Exercises

The statements in questions 1.1–1.3 are either TRUE or FALSE.

**1.1**   Biological variation:
(a) Is the main cause of differences between animals.
(b) Is the term given to differences between animals in a population.
(c) Is the reason why statistics is necessary in animal science.
(d) Makes it impossible to be sure of any aspect of animal science.
(e) Is the term given to the variation in ability of a technician performing a monotonous task throughout the day.

**1.2**   A sample is *randomly* drawn from a population:
(a) To reduce the study to a manageable size.

(b) To ensure that the full range of possibilities is included.
(c) To obtain 'normal' animals.
(d) To obtain a representative group.
(e) To avoid selector preferences.

**1.3** A nominal scale of measurement is used for data that:
(a) Comprise categories which cannot be ordered.
(b) Are not qualitative.
(c) Take many possible discrete quantitative values.
(d) Are evaluated as percentages.
(e) Are ranked.

**1.4** Decide whether the following errors are likely to be systematic or random (S or R):
(a) The water bath that holds samples for an enzyme assay fails during incubation.
(b) A clinician reading a clinical thermometer has a digit preference for the numbers 0 and 5.
(c) The calibration on a colorimeter was not checked before use.
(d) Scales for measuring the weight of animal feed packs are activated sometimes before the sack is put on and sometimes after, depending on the operator.
(e) A chemical balance weighing to 100 mg is used to weigh quantities of 2550 mg.

**1.5** Decide whether the following are either real or hypothetical populations (R or H):
(a) Milking cows in a trial for the effectiveness of a novel mastitis treatment.
(b) Horses in livery stables in the southeast of England.
(c) Fleas on dogs in urban Liverpool.
(d) Fleas on dogs treated with an oral monthly ectoparasite treatment.
(e) Blood glucose levels in diabetic dogs.

**1.6** Identify the appropriate type of variable (nominal, ordinal, discrete or continuous: N, O, D or C) for the following data:

(a) Coat colour of cats: in a colony of 35 cats there were one white, three black, seven ginger, seven agouti, 11 tortoiseshell and six of other colours.
(b) Percentages of motile spermatozoa in the ejaculates of six bulls at an artificial insemination centre collected on a single day during March: they were 73%, 81%, 64%, 76%, 69% and 84%.
(c) Spectrophotometer measurements of maximum light absorbance at a wavelength of 280 nm of solutions of egg yolk proteins: they were 0.724, 0.591 and 0.520 arbitrary units.
(d) The motility of a series of frozen and thawed samples of spermatozoa estimated on an arbitrary scale of 0–10 (0 indicating a completely immotile sample).
(e) Plasma progesterone levels (ng/ml) measured monthly in pregnant sheep throughout gestation by means of radioimmunoassay.
(f) Kittens classified 1 week post-natally as either flat-chested (abnormal) or normal.
(g) The optical density of negative micrographs of fluorescent cells calculated from measurements obtained with a densitometer: the results for groups A, B and C were 0.814, 0.986 and 1.103 units, respectively.
(h) Litter sizes of rabbits during an investigation of behavioural disturbances about the time of implantation.
(i) Body condition scores of goats.
(j) Numbers of deaths due to particular diseases in a year studied in an epidemiological investigation.
(k) Radioactivity determined by scintillation counts per minute in a $\beta$-counter.
(l) The gestation length (days) in cattle carrying twins and in those carrying singletons.

# 2 Descriptive statistics

## 2.1 Learning objectives

By the end of this chapter, you should be able to:

- Explain, with diagrams, the concepts of frequency distributions.
- Interpret diagrams of the frequency distributions of both categorical and numerical data.
- Identify frequency distributions that are skewed to the right and skewed to the left.
- Describe and conduct strategies to compare frequency distributions that have different numbers of observations.
- List the essential attributes of good tables and good diagrams.
- Interpret a pie chart, bar chart, dot diagram and histogram and state their appropriate uses.
- Interpret a stem-and-leaf diagram and a box-and-whisker plot, and state their appropriate uses.
- Interpret a scatter diagram and explain its usage.
- List different measures of location and identify their strengths and limitations.
- List different measures of dispersion and identify their strengths and limitations.
- Summarize any given data set appropriately in tabular and/or diagrammatic form to demonstrate its features.

## 2.2 Summarizing data

We collect data with the intention of gleaning information which, usually, we then convey to interested parties. This presents little problem when the data set comprises relatively few observations made on a small group of animals. However, as the quantity of information grows, it becomes increasingly difficult to obtain an overall 'picture' of what is happening.

The first stage in the process of obtaining this picture is to organize the data to establish how often different values occur (see frequency distributions in Section 2.3). Then it is helpful to further condense the information, reducing it to a manageable size, and so obtain a snapshot view as an aid to understanding and interpretation. There are various stratagems that we adopt; most notably, we can use:

- **Tables** to exhibit features of the data (see Section 2.4).
- **Diagrams** to illustrate patterns (see Section 2.5).
- **Numerical measures** to summarize the data (see Section 2.6).

## 2.3 Empirical frequency distributions

### 2.3.1 What is a frequency distribution?

A **frequency distribution** shows the frequencies of occurrence of the observations in a data set. Often the distribution of the observed data is