A PELICAN BOOK

The Art of Statistics Learning from Data David Spiegelhalter





DAVID SPIEGELHALTER

The Art of Statistics Learning from Data

A PELICAN BOOK



PELICAN an imprint of PENGUIN BOOKS Contents

LIST OF FIGURES

LIST OF TABLES

ACKNOWLEDGEMENTS

INTRODUCTION

CHAPTER 1

Getting Things in Proportion: Categorical Data and Percentages

CHAPTER 2

Summarizing and Communicating Numbers. Lots of Numbers

CHAPTER 3

Why Are We Looking at Data Anyway? Populations and Measurement

CHAPTER 4

What Causes What?

CHAPTER 5

Modelling Relationships Using Regression

CHAPTER 6

Algorithms, Analytics and Prediction

CHAPTER 7

How Sure Can We Be About What Is Going On? Estimates and

Intervals

CHAPTER 8

Probability - the Language of Uncertainty and Variability

CHAPTER 9

Putting Probability and Statistics Together

CHAPTER 10

Answering Questions and Claiming Discoveries

CHAPTER 11

Learning from Experience the Bayesian Way

CHAPTER 12

How Things Go Wrong

CHAPTER 13

How We Can Do Statistics Better

CHAPTER 14

In Conclusion

GLOSSARY

NOTES

INDEX

About the Author

Sir David John Spiegelhalter is a British statistician and Chair of the Winton

Centre for Risk and Evidence Communication in the Statistical Laboratory at

the University of Cambridge. Spiegelhalter is one of the most cited and influential researchers in his field, and was elected as President of the Royal Statistical Society for 2017–18.

To statisticians everywhere, with their endearing

traits of pedantry, generosity, integrity, and desire

to use data in the best way possible

List of Figures

0.1 Age and Year of Death of Harold Shipman's Victims

0.2 Time of Death of Harold Shipman's Patients

0.3 The PPDAC Problem-Solving Cycle

1.1 30-Day Survival Rates Following Heart Surgery

1.2 Proportion of Child Heart Operations Per Hospital

1.3 Percentage of Child Heart Operations Per Hospital

1.4 Risk of Eating Bacon Sandwiches

2.1 Jar of Jelly Beans

2.2 Different Ways of Displaying Jelly Bean Guesses

2.3 Jelly-Bean Guesses Plotted on a Logarithmic Scale

2.4 Reported Number of Lifetime Opposite-Sex Partners

2.5 Survival Rates Against Number of Operations in Child Heart Surgery

2.6 Pearson Correlation Coefficients of 0

- 2.7 World Population Trends
- 2.8 Relative Increase in Population by Country
- 2.9 Popularity of the Name 'David' Over Time
- 2.10 Infographic on Sexual Attitudes and Lifestyles
- 3.1 Diagram of Inductive Inference
- 3.2 Distribution of Birth Weights
- 5.1 Scatter of Sons' Heights v. Fathers' Heights
- 5.2 Fitted Logistic Regression Model for Child Heart Surgery Data
- 6.1 Memorial to Titanic Victim
- 6.2 Summary Survival Statistics for Titanic Passengers
- 6.3 Classification Tree for *Titanic* Data
- 6.4 ROC Curves for Algorithms Applied to Training and Test Sets
- 6.5 Probabilities of Surviving the Titanic Sinking
- 6.6 Over-Fitted Classification Tree for Titanic Data
- 6.7 Post-Surgery Survival Rates for Women with Breast Cancer
- 7.1 Empirical Distribution of Number of Sexual Partners for Varying Sample
- <u>Sizes</u>
- 7.2 Bootstrap Resamples from Original Sample of 50
- 7.3 Bootstrap Distribution of Means at Varying Sample Sizes
- 7.4 Bootstrap Regressions on Galton's Mother-Daughter Data

- 8.1 A Simulation of the Chevalier de Méré's Games
- 8.2 Expected Frequency Tree for Two Coin Flips
- 8.3 Probability Tree for Flipping Two Coins
- 8.4 Expected Frequency Tree for Breast Cancer Screening
- 8.5 Observed and Expected Number of Homicides
- 9.1 Probability Distribution of Left-Handers
- 9.2 Funnel Plot of Bowel-Cancer Death Rates
- 9.3 BBC Plot of Opinion Polls Before the 2017 General Election
- 9.4 Homicide Rates in England and Wales
- 10.1 Sex Ratio for London Baptisms, 1629–1710
- 10.2 Empirical Distribution of Observed Difference in Proportions of
- Left/Right Arm Crossers
- 10.3 Cumulative Number of Death Certificates Signed by Shipman
- 10.4 Sequential Probability Ratio Test for Detection of a Doubling in
- Mortality Risk
- 10.5 Expected Frequencies of the Outcomes of 1,000 Hypothesis Tests
- <u>11.1 Expected Frequency Tree for Three-Coin Problem</u>
- <u>11.2 Expected Frequency Tree for Sports Doping</u>
- 11.3 Reversed Expected Frequency Tree for Sports Doping
- 11.4 Bayes' 'Billiard' Table

12.1 Traditional Information Flows for Statistical Evidence

List of Tables

1.1 Outcomes of Children's Heart Surgery

1.2 Methods for Communicating the Lifetime Risk of Bowel Cancer in

Bacon Eaters

2.1 Summary Statistics for Jelly-Bean Guesses

2.2 Summary Statistics for the Lifetime Number of Sexual Partners

4.1 Outcomes for Patients in the Heart Protection Study

4.2 Illustration of Simpson's Paradox

5.1 Summary Statistics of Heights of Parents and Their Adult Children

5.2 Correlations between Heights of Adult Children and Parent of Same

Gender

5.3 Results of a Multiple Linear Regression Relating Adult Offspring Height

to Mother and Father

6.1 Error Matrix of Classification Tree on *Titanic* Training and Test Data

6.2 Fictional 'Probability of Precipitation' Forecasts

6.3 Results of a Logistic Regression for Titanic Survivor Data

6.4 Performance of Different Algorithms on *Titanic* Test Data

6.5 Breast Cancer Survival Rates Using the Predict 2.1 Algorithm

7.1 Summary Statistics for Lifetime Sexual Partners Reported by Men

7.2 Sample Means of Lifetime Sexual Partners Reported by Men

9.1 Comparison of Exact and Bootstrap Confidence Intervals

10.1 Cross-Tabulation of Arm-Crossing Behaviour by Gender

10.2 Observed and Expected Counts of Arm-Crossing by Gender

10.3 Observed and Expected Days with Each Number of Homicide Incidents

10.4 Results of Heart Protection Study with Confidence Intervals and P-

<u>values</u>

10.5 The Output in *R* of a Multiple Regression Using Galton's Data

10.6 Possible Outcomes of a Hypothesis Test

11.1 Likelihood Ratios for Evidence Concerning Richard III Skeleton

11.2 Recommended Verbal Interpretations of Likelihood Ratios

11.3 Kass and Raftery's Scale for Interpretation of Bayes Factors

12.1 Questionable Interpretation and Communication Practices

13.1 Exit Poll Predictions for Three Recent General Elections

Acknowledgements

Any insights gained from a long career in statistics come from listening to inspiring colleagues. These are too numerous even for a statistician to count, but a shortlist of those I have stolen most from might include Nicky Best, Sheila Bird, David Cox, Philip Dawid, Stephen Evans, Andrew Gelman, Tim Harford, Kevin McConway, Wayne Oldford, Sylvia Richardson, Hetan Shah, Adrian Smith and Chris Wild. I am grateful to them and so many others for encouraging me in a challenging subject.

This book has been a long time in development, entirely due to my chronic

procrastination. So I would primarily like to thank Laura Stickney of Penguin

for not only commissioning the book, but remaining calm as the months, and

years, went by, even when the book was finished and we still could not agree

on a title. And all credit to Jonathan Pegg for negotiating me a fine deal, Jane

Birdsell for showing huge patience when editing, and all the production staff

at Penguin for their meticulous work.

I am very grateful for permission to adapt illustrations, specifically Chris

Wild (Figure 0.3), James Grime (Figure 2.1), Cath Mercer of Natsal (Figures

2.4 and 2.10), Office for National Statistics (Figures 2.9, 8.5 and 9.4), Public

Health England (<u>Figure 6.7</u>), Paul Barden (<u>Figure 9.2</u>), and BBC (<u>Figure 9.3</u>).

UK public sector information is licensed under the Open Government

Licence v3.0.

I am not a good R programmer, and Matthew Pearce and Maria

Skoularidou helped me enormously in doing the analyses and graphics. I also

struggle with writing, and so am indebted to numerous people who read and

commented on chapters, including George Farmer, Alex Freeman, Cameron Brick, Michael Posner, Sander van der Linden and Simone Warr: in particular Julian Gilbey had a fine eye for errors and ambiguity. Above all, I must thank Kate Bull not only for her vital comments on the text, but also for supporting me through times that have been both good (writing in a beach hut in Goa) and not so good (a wet February juggling too many commitments).

I am also profoundly grateful to David and Claudia Harding for both their financial support and their continued encouragement, which has enabled me to do such fun things over the last ten years.

Finally, much as I would like to find someone else to blame, I am afraid I must acknowledge full responsibility for the inevitable remaining inadequacies of this book.

CODE FOR EXAMPLES

R code and data for reproducing most of the analyses and Figures are

available from <u>https://github.com/dspiegel29/ArtofStatistics</u>. I am grateful for the assistance received in preparing this material.

Introduction

The numbers have no way of speaking for themselves. We speak for them. We imbue

them with meaning.

— Nate Silver, *The Signal and the Noise1*

Why We Need Statistics

Harold Shipman was Britain's most prolific convicted murderer, though he does not fit the archetypal profile of a serial killer. A mild-mannered family doctor working in a suburb of Manchester, between 1975 and 1998 he injected at least 215 of his mostly elderly patients with a massive opiate overdose. He finally made the mistake of forging the will of one of his victims so as to leave him some money: her daughter was a solicitor, suspicions were aroused, and forensic analysis of his computer showed he had been retrospectively changing patient records to make his victims appear sicker than they really were. He was well known as an enthusiastic early adopter of technology, but he was not tech-savvy enough to realize that every

change he made was time-stamped (incidentally, a good example of data revealing hidden meaning).

Of his patients who had not been cremated, fifteen were exhumed and lethal levels of diamorphine, the medical form of heroin, were found in their bodies. Shipman was subsequently tried for fifteen murders in 1999, but chose not to offer any defence and never uttered a word at his trial. He was found guilty and jailed for life, and a public inquiry was set up to determine what crimes he might have committed apart from those for which he had been tried, and whether he could have been caught earlier. I was one of a number of statisticians called to give evidence at the public inquiry, which concluded that he had definitely murdered 215 of his patients, and possibly 45 more. $\underline{2}$

This book will focus on using <u>statistical sciencefn1</u> to answer the kind of questions that arise when we want to better understand the world – some of these questions will be highlighted in a box. In order to get some insight into Shipman's behaviour, a natural first question is:

What kind of people did Harold Shipman murder, and when did they die? The public inquiry provided details of each victim's age, gender and date of death. <u>Figure 0.1</u> is a fairly sophisticated visualization of this data, showing a

scatter-plot of the age of victim against their date of death, with the shading of the points indicating whether the victim was male or female. Bar-charts have been superimposed on the axes showing the pattern of ages (in 5–year bands) and years.

Some conclusions can be drawn by simply taking some time to look at the figure. There are more black than white dots, and so Shipman's victims were mainly women. The bar-chart on the right of the picture shows that most of his victims were in their 70s and 80s, but looking at the scatter of points reveals that although initially they were all elderly, some younger cases crept

in as the years went by. The bar-chart at the top clearly shows a gap around 1992 when there were no murders. It turned out that before that time Shipman had been working in a joint practice with other doctors but then, possibly as he felt under suspicion, he left to form a single-handed general practice. After this his activities accelerated, as demonstrated by the top barchart.



Figure 0.1

A scatter-plot showing the age and the year of death of Harold Shipman's 215 confirmed victims. Bar-charts have been added on

the axes to reveal the pattern of ages and the pattern of years in which he committed murders.

This analysis of the victims identified by the inquiry raises further

questions about the way he committed his murders. Some statistical evidence

is provided by data on the time of day of the death of his supposed victims, as

recorded on the death certificate. Figure 0.2 is a line graph comparing the

times of day that Shipman's patients died to the times that a sample of

patients of other local family doctors died. The pattern does not require subtle

analysis: the conclusion is sometimes known as 'inter-ocular', since it hits

you between the eyes. Shipman's patients tended overwhelmingly to die in

the early afternoon.

The data cannot tell us *why* they tended to die at that time, but further investigation revealed that he performed his home visits after lunch, when he was generally alone with his elderly patients. He would offer them an



injection that he said was to make them more comfortable, but which was in fact a lethal dose of diamorphine: after a patient had died peacefully in front of him, he would change their medical record to make it appear as if this was an expected natural death. Dame Janet Smith, who chaired the public inquiry,

later said, 'I still do feel it was unspeakably dreadful, just unspeakable and

unthinkable and unimaginable that he should be going about day after day pretending to be this wonderfully caring doctor and having with him in his bag his lethal weapon ... which he would just take out in the most matter-offact way.'

Figure 0.2

The time at which Harold Shipman's patients died, compared to the times at which patients of other local general practitioners died. The pattern does not require sophisticated statistical analysis. He was taking some risk, since a single post-mortem would have exposed him, but given the age of his patients and the apparent natural causes of death, none were performed. And his reasons for committing these murders have never been explained: he gave no evidence at his trial, never spoke about his misdeeds to anyone, including his family, and committed suicide in prison, conveniently just in time for his wife to collect his pension. We can think of this type of iterative, exploratory work as 'forensic' statistics, and in this case it was literally true. There is no mathematics, no theory, just a search for patterns that might lead to more interesting questions.

The details of Shipman's misdeeds were determined using evidence specific to each individual case, but this data analysis supported a general understanding of how he went about his crimes. Later in this book, in <u>Chapter 10</u>, we will see whether formal statistical analysis could have helped catch Shipman earlier<u>.fn2</u> In the meantime, the Shipman story amply demonstrates the great potential of using data to help us

understand the world and make better judgements. This is what statistical science is all about.

Turning the World Into Data

A statistical approach to Harold Shipman's crimes required us to stand back

from the long list of individual tragedies for which he was responsible. All

those personal, unique details of people's lives, and deaths, had to be reduced

to a set of facts and numbers that could be counted and drawn on graphs. This

might at first seem cold and dehumanizing, but if we are to use statistical science to illuminate the world, then our daily experiences have to be turned into data, and this means categorizing and labelling events, recording measurements, analysing the results and communicating the conclusions. Simply categorizing and labelling can, however, present a serious challenge. Take the following basic question, which should be of interest to everyone concerned with our environment:

How many trees are there on the planet?

Before even starting to think about how we might go about answering this

question, we first have to settle a rather basic issue. What is a 'tree'? You may feel you know a tree when you see it, but your judgement may differ considerably from others who might consider it a bush or a shrub. So to turn experience into data, we have to start with rigorous definitions. It turns out that the official definition of a 'tree' is a plant with a woody

stem that has a sufficiently large diameter at breast height, known as the DBH. The US Forest Service demands a plant has a DBH of greater than 5 inches (12.7 cm) before officially declaring it a tree, but most authorities use a DBH of 10 cm (4 inches).

But we cannot wander round the entire planet individually measuring each woody-stemmed plant and counting up those that meet this criterion. So the researchers who investigated this question took a more pragmatic approach: they first took a series of areas with a common type of landscape, known as a

biome, and counted the average number of trees found per square kilometre. They then used satellite imaging to estimate the total area of the planet covered by each type of biome, carried out some complex statistical modelling, and eventually came up with an estimated total of 3.04 trillion (that is 3,040,000,000,000) trees on the planet. This sounds a lot, except they reckoned there used to be twice this number.<u>fn3 3</u>

If authorities differ about what they call a tree, it should be no surprise that

more nebulous concepts are even more challenging to pin down. To take an extreme example, the official definition of 'unemployment' in the UK was changed at least thirty-one times between 1979 and 1996. <u>4</u> The definition of Gross Domestic Product (GDP) is continually being revised, as when trade in

illegal drugs and prostitution was added to the UK GDP in 2014; the estimates used some unusual data sources – for example Punternet, a review website that rates prostitution services, provided prices for different

activities. 5

Even our most personal feelings can be codified and subjected to statistical

analysis. In the year ending September 2017, 150,000 people in the UK were

asked as part of a survey: 'Overall, how happy did you feel yesterday? <u>'6</u>

Their average response, on a scale from zero to ten, was 7.5, an improvement

from 2012 when it was 7.3, which might be related to economic recovery since the financial crash of 2008. The lowest scores were reported for those aged between 50 and 54, and the highest between 70 and 74, a typical pattern

for the UK<u>.fn4</u>

Measuring happiness is hard, whereas deciding whether someone is alive or dead should be more straightforward: as the examples in this book will demonstrate, survival and mortality is a common concern of statistical science. But in the US each state can have its own legal definition of death, and although the Uniform Declaration of Death Act was introduced in 1981 to try to establish a common model, some small differences remain. Someone

who had been declared dead in Alabama could, at least in principle, cease to be legally dead were they across the state border in Florida, where the registration must be made by two qualified doctors. <u>7</u>

These examples show that statistics are always to some extent constructed on the basis of judgements, and it would be an obvious delusion to think the full complexity of personal experience can be unambiguously coded and put into a spreadsheet or other software. Challenging though it is to define, count and measure characteristics of ourselves and the world around us, it is still just information, and only the starting point to real understanding of the world.

Data has two main limitations as a source of such knowledge. First, it is almost always an imperfect measure of what we are really interested in: asking how happy people were last week on a scale from zero to ten hardly encapsulates the emotional wellbeing of the nation. Second, anything we choose to measure will differ from place to place, from person to person, from time to time, and the problem is to extract meaningful insights from all this apparently random <u>variability</u>. For centuries, statistical science has faced up to these twin challenges, and played a leading role in scientific attempts to understand the world. It has provided the basis for interpreting data, which is always imperfect, in order to

distinguish important relationships from the background variability that makes us all unique. But the world is always changing, as new questions are asked and new sources of data become available, and statistical science has had to change too.

People have always counted and measured, but modern statistics as a discipline really began in the 1650s when, as we shall see in <u>Chapter 8</u>, probability was properly understood for the first time by Blaise Pascal and Pierre de Fermat. Given this solid mathematical basis for dealing with variability, progress was then remarkably rapid. When combined with data on

the ages at which people die, the theory of probability provided a firm basis for calculating pensions and annuities. Astronomy was revolutionized when scientists grasped how probability theory could handle variability in measurements. Victorian enthusiasts became obsessed with collecting data about the human body (and everything else), and established a strong connection between statistical analysis and genetics, biology and medicine. Then in the twentieth century statistics became more mathematical and, unfortunately for many students and practitioners, the topic became synonymous with the mechanical application of a bag of statistical tools, many named after eccentric and argumentative statisticians that we shall meet

later in this book.

This common view of statistics as a basic 'bag of tools' is now facing major challenges. First, we are in an age of <u>data science</u>, in which large and complex data sets are collected from routine sources such as traffic monitors, social media posts and internet purchases, and used as a basis for technological innovations such as optimizing travel routes, targeted advertising or purchase recommendation systems – we shall look at algorithms based on <u>'big data</u>' in <u>Chapter 6</u>. Statistical training is increasingly seen as just one necessary component of being a data scientist, together with skills in data management, programming and algorithm development, as well as proper knowledge of the subject matter. Another challenge to the traditional view of statistics comes from the huge rise in the amount of scientific research being carried out, particularly in the biomedical and social sciences, combined with pressure to publish in highranking journals. This has led to doubts about the reliability of parts of the scientific literature, with claims that many 'discoveries' cannot be reproduced

by other researchers – such as the continuing dispute over whether adopting

an assertive posture popularly known as a 'power pose' can induce hormonal and other changes. <u>8</u> The inappropriate use of standard statistical methods has

received a fair share of the blame for what has become known as the reproducibility or replication crisis in science.

With the growing availability of massive data sets and user-friendly analysis software, it might be thought that there is less need for training in statistical methods. This would be naïve in the extreme. Far from freeing us from the need for statistical skills, bigger data and the rise in the number and complexity of scientific studies makes it even more difficult to draw appropriate conclusions. More data means that we need to be even more aware of what the evidence is actually worth.

For example, intensive analysis of data sets derived from routine data can increase the possibility of false discoveries, both due to systematic bias inherent in the data sources and from carrying out many analyses and only reporting whatever looks most interesting, a practice sometimes known as 'data-dredging'. In order to be able to critique published scientific work, and even more the media reports which we all encounter on a daily basis, we should have an acute awareness of the dangers of selective reporting, the need for scientific claims to be replicated by independent researchers, and the danger of over-interpreting a single study out of context.

All these insights can be brought together under the term <u>data literacy</u>, which describes the ability to not only carry out statistical analysis on realworld problems, but also to understand and critique any conclusions drawn by others on the basis of statistics. But improving data literacy means changing the way statistics is taught.

Teaching Statistics

Generations of students have suffered through dry statistics courses based on learning a set of techniques to be applied in different situations, with more regard to mathematical theory than understanding both why the formulae are being used, and the challenges that arise when trying to use data to answer questions.

Fortunately this is changing. The needs of data science and data literacy demand a more problem-driven approach, in which the application of specific statistical tools is seen as just one component of a complete cycle of investigation. The **PPDAC** structure has been suggested as a way of representing a problem-solving cycle, which we shall adopt throughout this book. <u>9 Figure 0.3</u> is based on an example from New Zealand, which has been a world-leader in statistics education in schools.

The first stage of the cycle is specifying a Problem; statistical inquiry

always starts with a question, such as our asking about the pattern of Harold Shipman's murders or the number of trees in the world. Later in this book we

shall focus on problems ranging from the expected benefit of different therapies immediately following breast cancer surgery, to why old men have big ears.

It is tempting to skip over the need for a careful Plan. The Shipman question simply required the collection of as much data as possible on his victims. But the people counting trees paid meticulous attention to precise definitions and how to carry out the measurements, since confident conclusions can only be drawn from a study which has been appropriately



designed. Unfortunately, in the rush to get data and start analysis, attention to

design is often glossed over.

Figure 0.3

The PPDAC problem-solving cycle, going from Problem, Plan, Data, Analysis to Conclusion and communication, and starting again on another cycle.

Collecting good Data requires the kind of organizational and coding skills that are being seen as increasingly important in data science, particularly as data from routine sources may need a lot of cleaning in order to get it ready to be analysed. Data collection systems may have changed over time, there may be obvious errors, and so on – the phrase 'found data' neatly communicates that it may be rather messy, like something picked up in the street.

The Analysis stage has traditionally been the main emphasis of statistics courses, and we shall cover a range of analytic techniques in this book; but sometimes all that is required is a useful visualization, as in Figure 0.1.

Finally, the key to good statistical science is drawing appropriate Conclusions

that fully acknowledge the limitations in the evidence, and communicating them clearly, as in the graphical illustrations of the Shipman data. Any conclusions generally raise more questions, and so the cycle starts over again,

as when we started looking at the time of day when Shipman's patients died.

Although in practice the PPDAC cycle laid out in <u>Figure 0.3</u> may not be followed precisely, it underscores that formal techniques for statistical analysis play only one part in the work of a statistician or data scientist. Statistical science is a lot more than a branch of mathematics involving esoteric formulae with which generations of students have (often reluctantly) struggled.

This Book

When I was a student in Britain in the 1970s, there were just three TV channels, computers were the size of a double wardrobe, and the closest thing

we had to Wikipedia was on the imaginary handheld device in Douglas Adams' (remarkably prescient) *Hitchhiker's Guide to the Galaxy*. For selfimprovement we therefore turned to Pelican books, and their iconic blue spines were a standard feature of every student bookshelf.

Because I was studying statistics, my Pelican collection featured Facts

from Figures by M. J. Moroney (1951) and How to Lie with Statistics by

Darrell Huff (1954). These venerable publications sold in the hundreds of

thousands, reflecting both the level of interest in statistics and the dismal lack

of choice at that time. These classics have stood up remarkably well to the

intervening sixty-five years, but the current era demands a different approach

to teaching statistics based on the principles laid out above.

This book therefore uses real-world problem-solving as a starting point for introducing statistical ideas. Some of these ideas may seem obvious, but some are more subtle and may require some mental effort, although mathematical skills will not be needed. Compared to traditional texts, this book focuses on conceptual issues rather than technicalities, and features only

a few, fairly innocuous equations supported by a Glossary. Software is a vital part of any work in data science and statistics but it is not a focus of this book

 tutorials are readily available for freely available environments such as R and Python.

The questions featured in the boxes can all, to a certain extent, be answered through statistical analysis, although they differ widely in their scope. Some are important scientific hypotheses, such as whether the Higgs boson exists, or if there really is convincing evidence for extra-sensory perception (ESP). Others are questions about health care, such as whether busier hospitals have higher survival rates, and if screening for ovarian cancer is beneficial. Sometimes we just want to estimate quantities, such as the cancer risk from bacon sandwiches, the number of sexual partners people in Britain have in their lifetime, and the benefit of taking a daily statin. And some questions are just interesting, such as identifying the luckiest survivor from the *Titanic*; whether Harold Shipman could have been caught earlier; and assessing the probability that a skeleton found in a Leicester car park really was that of Richard III.

This book is intended for both students of statistics who are seeking a nontechnical introduction to the basic issues, and general readers who want to be more informed about the statistics they encounter both in their work and in everyday life. My emphasis is on handling statistics skilfully and with care: numbers may appear to be cold, hard facts, but the attempts to measure trees, happiness and death have already shown that they need to be treated with delicacy.

Statistics can bring clarity and insight into the problems we face, but we are all familiar with the way they can be abused, often to promote an opinion or simply to attract attention. The ability to assess the trustworthiness of statistical claims seems a key skill in the modern world, and I hope that this book may help to empower people to question the numbers that they encounter in their daily life.

Summary

Turning experiences into data is not straightforward, and data is inevitably limited in its

capacity to describe the world.

Statistical science has a long and successful history, but is now changing in the light of

increased availability of data.

Skill in statistical methods plays an important part of being a data scientist.

Teaching statistics is changing from a focus on mathematical methods to one based on an

entire problem-solving cycle.

The PPDAC cycle provides a convenient framework: Problem – Plan – Data – Analysis –

Conclusion and communication.

Data literacy is a key skill for the modern world.

CHAPTER1

Getting Things in Proportion: Categorical

Data and Percentages

What happened to children having heart surgery in Bristol between 1984 and 1995?

Joshua L was 16 months old and had transposition of the great arteries, a

severe form of congenital heart disease in which the main vessels coming

from the heart are attached to the wrong ventricle. He needed an operation to

'switch' the arteries, and just after 7 a.m. on 12 January 1995 his parents said

goodbye to him and watched as he was taken for his surgery in Bristol Royal

Infirmary. But Joshua's parents were unaware that stories about the poor

surgical survival rates at Bristol had been circulating since the early 1990s. Nobody told them that nurses had left the unit rather than continue telling parents that their child had died, or that the previous evening there had been a

late-night meeting at which it had been debated whether to cancel Joshua's operation. **1**

Joshua died on the operating table. The following year the General Medical Council (the medical regulator) launched an investigation after complaints from Joshua's and other bereaved parents, and in 1998 two surgeons and the ex-chief executive were found guilty of serious medical misconduct. Public concern did not die down, and an official inquiry was ordered: this brought in a team of statisticians who were given the grim task of comparing the survival rates in Bristol with elsewhere in the UK between 1984 and 1995. I led this team.

We first had to determine how many children had had heart surgery, and how many had died. This sounds like it should be straightforward but, as shown in the previous chapter, simply counting events can be challenging. What is a 'child'? What counts as 'heart surgery'? When can death be attributed to surgery? And even when these definitions have been decided, could we determine how many of each there had been? We took a 'child' as anyone under 16, and focused on 'open' surgery in which the heart had been stopped and its function replaced by cardiopulmonary bypass. There can be multiple operations per admission, but these were considered as one event. Deaths were counted if they occurred within 30 days of the operation, whether or not in hospital or due to the surgery. We knew that death was an imperfect measure of the quality of the outcome, as it

ignored children who were brain-damaged or otherwise disabled as a result of

the surgery, but we did not have the data on longer-term outcomes.

The main source of data was national Hospital Episode Statistics (HES), which were derived from administrative data entered by low-paid coders. HES had a poor reputation among doctors, but this source had the great advantage that it could be linked to national death records. There was also a parallel system of data submitted directly to a Cardiac Surgical Registry (CSR) established by the surgeons' professional society.

These two sources of data, though they were supposed to be about exactly the same practice, showed considerable disagreement: for 1991–1995, HES said there had been 62 deaths out of 505 open operations (14%), whereas CSR said there had been 71 deaths out of 563 operations (13%). No less than five additional local sources of data were available, from anaesthetic records to the surgeons' own personal logs. Bristol was awash with data, but none of the data sources could be considered the 'truth', and nobody had taken responsibility for analysing and acting on the surgical outcomes.

We calculated that if patients at Bristol had the average risk prevailing elsewhere in the UK, Bristol would have expected to have had 32 deaths over

this period, instead of the 62 recorded in HES, which we reported as '30 excess deaths' between 1991 and 1995. <u>fn1</u> The exact numbers varied according to the data sources, and it may seem extraordinary that we could not even establish the basic facts about the number of operations and their outcome, although current record systems should be better.

These findings had wide press coverage, and the Bristol inquiry led to a major change in attitudes to monitoring clinical performance: no longer was the medical profession trusted to police itself. Mechanisms to publicly report hospital survival data were established, although, as we shall now see, the way in which that data is displayed can itself influence the perception of audiences.

Communicating Counts and Proportions

Data that records whether individual events have happened or not is known as <u>binary data</u>, as it can only take on two values, generally labelled as yes and no. Sets of binary data can be summarized by the number of times and

the percentage of cases in which an event occurred.

The theme of this chapter is that the basic presentation of statistics is important. In a sense we are jumping to the last step of the PPDAC cycle in which conclusions are communicated, and while the form of this communication has not traditionally been considered an important topic in statistics, rising interest in data visualization reflects a change in this attitude.

So both in this chapter and the next we shall concentrate on ways of displaying data so that we can quickly get the gist of what is going on without

detailed analysis, starting with a look at alternative ways of displaying data that, largely because of the Bristol inquiry, are now publicly available. <u>Table 1.1</u> shows the outcomes of nearly 13,000 children who had heart surgery in the UK and Ireland between 2012 and 2015. <u>2</u> Two hundred and sixty-three babies died within 30 days of their operation, and every one of these deaths is a tragedy to the family involved. It will be little consolation

to

them that survival rates have improved hugely from the time of the Bristol inquiry, and now average 98%, and so there is a more hopeful prospect for families of children facing heart surgery.

A table can be considered as a type of graphic, and requires careful design choices of colour, font and language to ensure engagement and readability. The audience's emotional response to the table may also be influenced by the

choice of which columns to display. <u>Table 1.1</u> shows the results in terms of both survivors and deaths, but in the US *mortality* rates from child heart surgery are reported, while the UK provides *survival* rates. This is known as negative or positive <u>framing</u>, and its overall effect on how we feel is intuitive and well-documented: '5% mortality' sounds worse than '95% survival'. Reporting the actual number of deaths as well as the percentage can also increase the impression of risk, as this total might then be imagined as a crowd of real people.

Hospital	Number of babies having surgery	Number surviving for at least 30 days after surgcry	Number dying within 30 days of surgery	Percentage surviving	Percentage dying
London, Harley Street	418	413	5	98.8	1.2
Leicester	607	593	14	97.7	2.3
Newcastle	668	653	15	97.8	2.2
Glasgow	760	733	27	96.3	3.7
Southampton	829	815	14	98.3	1.7
Bristol	835	821	14	98.3	1.7
Dublin	983	960	23	97.7	2.3
Leeds	1,038	1,016	22	97.9	2.1
London, Brompton	1,094	1,075	19	98.3	1.7
Liverpool	1,132	1,112	20	98.2	1.8
London, Evelina	1,220	1,185	35	97.1	2.9
Birmingham	1,457	1,421	36	97.5	2.5
London, Great Ormond Street	1,892	1,873	19	99.0	1.0
Total	12,933	12,670	263	98.0	2.0

Table 1.1

Outcomes of children's heart surgery in UK and Irish hospitals between 2012 and 2015, in terms of

survival or not, 30 days after surgery.

A classic example of how alternative framing can change the emotional impact of a number is an advertisement that appeared on the London Underground in 2011, proclaiming that '99% of young Londoners do not commit serious youth violence'. These ads were presumably intended to reassure passengers about their city, but we could reverse its emotional impact with two simple changes. First, the statement means that 1% of young

Londoners *do* commit serious violence. Second, since the population of London is around 9 million, there are around 1 million people aged between 15 and 25, and if we consider these as 'young', this means there are 1% of 1 million or a total of 10,000 seriously violent young people in the city. This does not sound at all reassuring. Note the two tricks used to manipulate the impact of this statistic: convert from a positive to a negative frame, and then turn a percentage into actual numbers of people.

Ideally both positive and negative frames should be presented if we want to provide impartial information, although the order of columns might still influence how the table is interpreted. The order of the rows of a table also needs to be considered carefully. <u>Table 1.1</u> shows the hospitals in order of the

number of operations in each, but if they had been presented, say, in order of mortality rates with the highest at the top of the table, this might give the impression that this was a valid and important way of comparing hospitals. Such league tables are favoured by the media and even some politicians, but can be grossly misleading: not only because the differences could be due to chance variation, but because the hospitals may be taking in very different types of cases. In <u>Table 1.1,</u> for example, we might suspect that Birmingham,

one of the biggest and most well-known children's hospitals, takes on the most severe cases, and so it would be unfair, to put it mildly, to highlight their apparently unimpressive overall survival rates.<u>fn2</u>

The survival rates can be presented in a horizontal bar-chart such as the one shown in Figure 1.1. A crucial choice is where to start the horizontal axis: if the values start from 0%, all the bars will be almost the full length of the graphic, which will clearly show the extraordinarily high survival rates, but the lines will be indistinguishable. But the oldest trick of misleading graphics is to start the axis at say 95%, which will make the hospitals look



extremely different, even if the variation is in fact only what is attributable to chance alone.

Choosing the start of the axis therefore presents a dilemma. Alberto Cairo, author of influential books on data visualization, <u>3</u> suggests you should

always begin with a 'logical and meaningful baseline', which in this situation

appears difficult to identify – my rather arbitrary choice of 86% roughly represents the unacceptably low survival in Bristol twenty years previously.

Figure 1.1

Horizontal bar-chart of 30–day survival rates for thirteen hospitals. The choice of the start of the horizontal axis, here 86%, can have a crucial effect on the impression given by the graphic. If the axis starts at 0%, all the hospitals will look indistinguishable, whereas if we started at 95% the differences would look misleadingly

dramatic.

I began this book with a quotation from Nate Silver, the founder of data-

based platform *FiveThirtyEight* and first famous for accurately predicting the

2008 US presidential election, who eloquently expressed the idea that numbers do not speak for themselves – we are responsible for giving them meaning. This implies that communication is a key part of the problemsolving cycle, and I have shown in this section how the message from a set of

simple proportions can be influenced by our choices of presentation.

We now need to introduce an important and convenient concept that will help us get beyond simple yes/no questions.

Categorical Variables

A variable is defined as any measurement that can take on different values in different circumstances; it's a very useful shorthand term for all the types of observations that comprise data. Binary variables are yes/no questions such as whether someone is alive or dead and whether they are female or not: both

of these vary between people, and can, even for gender, vary within people at

different times. <u>Categorical variables</u> are measures that can take on two or more categories, which may be

Unordered categories: such as a person's country of origin, the colour of a car, or the hospital in which an operation takes place.

Ordered categories: such as the rank of military personnel.

Numbers that have been grouped: such as levels of obesity, which is

often defined in terms of thresholds for the body mass index (BMI). fn3

When it comes to presenting categorical data, pie charts allow an impression

of the size of each category relative to the whole pie, but are often visually

confusing, especially if they attempt to show too many categories in the same

chart, or use a three-dimensional representation that distorts areas. <u>Figure 1.2</u> shows a fairly hideous example modelled on the kind offered by Microsoft <u>Excel, showing the proportions of the 12,933 child heart patients from Table</u>

1.1 that are treated in each hospital.

Multiple pie charts are generally not a good idea, as comparisons are

hampered by the difficulty in assessing the relative sizes of areas of different

shapes. Comparisons are better based on height or length alone in a bar chart.

Figure 1.3 shows a simpler, clearer example of a horizontal bar chart of the

proportions being treated in each hospital.

Comparing a Pair of Proportions



We have seen how a set of proportions can be elegantly compared using a bar

chart, and so it would be reasonable to think that comparing two proportions would be a trivial matter. But when these proportions represent estimates of the risks of experiencing some harm, then the way in which those risks are compared becomes a serious and contested issue. Here is a typical question: What's the cancer risk from bacon sandwiches?

Figure 1.2

The proportion of all child heart operations being carried out in each hospital, displayed in a 3D pie chart from Excel. This deeply unpleasant chart makes categories near the front look bigger, and so makes it impossible to make visual comparisons between hospitals.



Figure 1.3

Percentage of all child heart operations being carried out in each hospital: a clearer representation using a horizontal bar chart. We're all familiar with hyperbolic media headlines that warn us that