

Critical Appraisal for FCEM

DUNCAN BOOTLAND • EVAN COUGHLAN
ROBERT GALLOWAY • STEPHANIE GOUBET



CRC Press
Taylor & Francis Group

Critical Appraisal for FCEM



Critical Appraisal for FCEM

DUNCAN BOOTLAND

MBBS, BSc, FCEM
Consultant in Emergency Medicine,
Brighton and Sussex University Hospitals NHS Trust

EVAN COUGHLAN

MB BCh NUI, MRCSEd, DSEM, Diptox, FCEM
Emergency Medicine Consultant, Brighton and
Sussex University Hospitals NHS Trust

ROBERT GALLOWAY

MBBS, BSc, MRCP, FCEM, PGcMedED
Emergency Medicine Consultant, Brighton and
Sussex University Hospitals NHS Trust; Year Five
Sub-Dean and Lead for Undergraduate Emergency
Medicine, Brighton and Sussex Medical School

STEPHANIE GOUBET

MSc Medical Statistician for Brighton and
Sussex University Hospitals NHS Trust and
Brighton and Sussex Medical School



CRC Press

Taylor & Francis Group
Boca Raton London New York

CRC Press is an imprint of the
Taylor & Francis Group, an **informa** business

CRC Press
Taylor & Francis Group
6000 Broken Sound Parkway NW, Suite 300
Boca Raton, FL 33487-2742

© 2015 by Taylor & Francis Group, LLC
CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works

Printed on acid-free paper
Version Date: 20140902

International Standard Book Number-13: 978-1-4441-8648-2 (Paperback)

This book contains information obtained from authentic and highly regarded sources. While all reasonable efforts have been made to publish reliable data and information, neither the author[s] nor the publisher can accept any legal responsibility or liability for any errors or omissions that may be made. The publishers wish to make clear that any views or opinions expressed in this book by individual editors, authors or contributors are personal to them and do not necessarily reflect the views/opinions of the publishers. The information or guidance contained in this book is intended for use by medical, scientific or health-care professionals and is provided strictly as a supplement to the medical or other professional's own judgement, their knowledge of the patient's medical history, relevant manufacturer's instructions and the appropriate best practice guidelines. Because of the rapid advances in medical science, any information or advice on dosages, procedures or diagnoses should be independently verified. The reader is strongly urged to consult the relevant national drug formulary and the drug companies' printed instructions, and their websites, before administering any of the drugs recommended in this book. This book does not indicate whether a particular treatment is appropriate or suitable for a particular individual. Ultimately it is the sole responsibility of the medical professional to make his or her own professional judgements, so as to advise and treat patients appropriately. The authors and publishers have also attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright.com (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Library of Congress Cataloging-in-Publication Data

Bootland, Duncan, author.

Critical appraisal for FCEM / Duncan Bootland, Evan Coughlan, Robert Galloway, Stephanie Goubet.
p. ; cm.

Critical appraisal for Fellowship of the College of Emergency Medicine

Includes bibliographical references and index.

ISBN 978-1-4441-8648-2 (Hardcover : alk. paper)

I. Coughlan, Evan, author. II. Galloway, Robert, author. III. Goubet, Stephanie, author. IV. Title. V. Title:

Critical appraisal for Fellowship of the College of Emergency Medicine.

[DNLM: 1. Evidence-Based Emergency Medicine--Great Britain--Examination Questions. 2.

Evidence-Based Emergency Medicine--Ireland--Examination Questions. 3. Statistics as Topic--Great

Britain--Examination Questions. 4. Statistics as Topic--Ireland--Examination Questions. 5. Test Taking

Skills--Great Britain. 6. Test Taking Skills--Ireland. WB 18.2]

RC86.9

616.02'5076--dc23

2014034192

Visit the Taylor & Francis Web site at
<http://www.taylorandfrancis.com>

and the CRC Press Web site at
<http://www.crcpress.com>

Contents

Contents v
About the Authorsvii
Forewordix
Preface.....xi

Section 1 *Fundamental knowledge needed for critical appraisal*

1 Types of papers and grades of evidence 3
2 Understanding basic statistics 9
3 Statistics in papers – p-values, confidence intervals and error 23
4 Probability and odds 35

Section 2 *Critically appraising papers*

5 How to dissect a paper 47
6 Therapeutic studies 57
7 Diagnostic studies 75
8 Meta-analysis 95

Section 3 Passing FCEM

9 Passing the exam	101
10 The Critical Topic Review (CTR)	109
Appendix A: Checklists	115
Appendix B: Glossary	119
Appendix C: Practice Papers	133

About the Authors

Duncan Bootland MBBS, BSc, FCEM

Emergency Medicine Consultant, Brighton and Sussex University Hospitals NHS Trust

Duncan trained at University College London and has worked in the South East since then. His specialist training was in emergency medicine and intensive care as well as working for Kent, Surrey, Sussex HEMS. He was appointed as an emergency medicine consultant in 2012. Duncan has an excellent reputation in teaching, especially for College of Emergency Medicine examinations. He is a leading part of the Bromley Emergency Courses, leading for the MCEM B and FCEM SAQ revision courses. He has also authored chapters in two emergency medicine textbooks.

Evan Coughlan MB Bch NUI, MRCSEd, DSEM, FCEM, DipTox.

Emergency Medicine Consultant, Brighton and Sussex University Hospitals NHS Trust

Evan qualified in Ireland and finished his registrar training in the United Kingdom. He was appointed as a consultant in 2010. He has an interest in ultrasound, trauma and toxicology and is a passionate teacher. He is an established MCEM teacher on the Bromley Emergency Courses and has also written chapters for several textbooks.

Rob Galloway MBBS, BSc, MRCP, FCEM, PGcMedED

Emergency Medicine Consultant, Brighton and Sussex University Hospitals NHS Trust, Year Five Sub-Dean and Lead for Undergraduate Emergency Medicine, Brighton and Sussex Medical School

Rob trained at St George's Hospital Medical School and did his junior rotations in the South East. He dual trained in emergency medicine and intensive care and was appointed as an emergency medicine consultant in 2011. He has a passion for education, having done a PGC in medical education and leads for emergency medicine undergraduate education at Brighton and Sussex Medical School. He has interests in a wide range of academic topics including the effects of good rostering on the sustainability of emergency medicine careers, resuscitation and human factors.

Stephanie Goubet MSc

Medical Statistician for Brighton and Sussex University Hospitals NHS Trust, and Brighton and Sussex Medical School

Stephanie came to the Brighton and Sussex University Hospital National Institute for Health Research Clinical Research Facility from the Department of Primary Care and Population Health, University College London. While there she was a member of their Evidence Based Medicine workshops. She shares her time between Brighton and Sussex University Hospitals NHS Trust and Brighton and Sussex Medical School.

Foreword

Emergency medicine in the twenty-first century has developed from a service that sorted and directed patients to a surgical or medical ‘dresser’ into a respected, effective and accomplished specialty with its own curriculum and academic infrastructure. To be an emergency physician today requires us to ensure our patients have the very best available treatment in the most effective time frame and in the best location. In this way we can ensure that we minimize the disruption for our patients that arises from a hospital admission and maximize the utility of the health resources at our disposal. Emergency physicians must therefore be reflective and analytical, assessing the available evidence and ensuring that data is interpreted correctly to support appropriate decision making. It is in this context that the College of Emergency Medicine continues to include in the professional examinations the critical appraisal paper, which requires the candidate to demonstrate these analytical skills and to demonstrate the practical application by the evaluation of a single published paper.

This book sets out to provide the reader with those skills, in a simple and easy to read style and by covering the essential elements of critical appraisal. With the book’s very useful glossary and practice papers, the reader has a practical primer for success. However, the book continually reinforces the College view that critical appraisal is for life – and that if you can master the skills in this book in the same way as we master history and examination techniques in direct clinical care, we will be able to handle all situations – be it the chief executive who has heard about an unusual solution to performance problems or the consultant endocrinologist who has a great new idea to prevent his registrar from seeing patients in ‘casualty’. Remember, critical appraisal is a life skill – not just for the exam!

Dr Ruth Brown FCEM FIFEM
International Lead and Director of Examinations
College of Emergency Medicine

Preface

This book has one main aim – to help you pass the critical appraisal element of the Fellowship of the College of Emergency Medicine (FCEM) examination. However, the ability to assess the literature and make evidence-based decisions on patient care is a core skill for all doctors and we hope this book will also be of value to those who are not planning to take or have already taken the examination.

Three of us are emergency physicians with a predominant interest in clinical practice and not statistics. This book is not designed to be a medical statistics textbook but a guide on how to become competent at critical appraisal. This difference is important and we feel it is a major strength of the book which allows us to focus on only those aspects of critical appraisal that are most needed for practising emergency physicians and especially for those sitting the FCEM examination. Our fourth author, Stephanie Goubet, is a medical statistician who works within our hospital and medical school research team.

The book is based on the successful, and College of Emergency Medicine approved, Critical Appraisal for FCEM course that we run (www.criticalappraisal-forfcm.com). However, the book is designed to be read independently of the course and the course does not necessarily require any prior reading.

We have endeavoured to keep only what is needed for a working knowledge of critical appraisal within the main text but for those who want to extend their knowledge we have included ‘Spod’s Corner’ throughout the book. Additionally, at the end of each chapter there is a ‘Take Home Message’ as a review of what you must know.

At the end of the book we have included a chapter to guide those writing their Clinical Topic Review (CTR) as we strongly believe it is counter-productive to postpone learning about critical appraisal until after you have written your CTR.

We very much hope you enjoy the book and that it helps you become successful in passing FCEM and getting a good grounding in critical appraisal.

Section 1

*Fundamental knowledge
needed for critical appraisal*

1 Types of papers and grades of evidence

It's nearing the end of a Friday evening in the emergency department (ED) and lying on the trolley in front of me is a 47-year-old man who has attended complaining of breathlessness and a sharp left-sided chest pain. I'm tired and I want to go home, but even as a consultant in emergency medicine, I'm made anxious by the knowledge that chest pain is one of the 'banana skins' of our profession.

As much as any other, this situation reflects the challenges we face many times a day: what clinical information do I need to glean from my history and examination, what clinical decision tools should I use, what diagnostic tests should I subject my patient to and how confident can I be in them? When I finally make the diagnosis, what is the best treatment?

The ability to interpret the evidence and the guidelines that already exist and the new ones that are published every month, and to do this in a way that leaves enough time to actually see some patients is challenging. We need to remember that it is our desire to know how to treat our patients better that should drive our aim to understand evidence based medicine (EBM) and then an understanding of EBM should drive improvements in patient care.

Before we begin trying to work out how to dissect various papers we need to understand how research is divided up into different types, and which types are best for answering which questions. Although in the FCEM Critical Appraisal examination you are likely to get either a therapeutic or diagnostic study, when you write your CTR you may well need to include a number of other types of studies and thus knowledge of them is important. In this chapter we look at all of this before moving on to the more detailed aspects of critical appraisal in later chapters.

TYPES OF PAPERS

Reading through the medical literature, it won't take you long to realize that there are a number of different types of studies out there – the EBM world is not just randomized controlled trials (RCTs), diagnostic studies and meta-analyses.

The way different study types are grouped together depends in some part on whom you speak to but we find the following to be useful in ordering our thoughts when it comes to the different types. Think about three questions:

1. What sort of environment (a clinical laboratory, a highly controlled but still clinical setting, a normal ED etc.) is the trial being performed in?
2. What sort of result are the researchers interested in – is it a result like mortality rate or blood pressure that will give you a numerical value or is it a

subjective feeling like ‘wellbeing’ or ‘comfort’ that will be given in a less well-defined response?

- 3. Are the researchers planning to do something to the patient, e.g. give them a new drug or do a new diagnostic test, or are they just observing what happens to them?

The answers to these three questions should then allow us to describe the trial as explanatory or pragmatic, qualitative or quantitative and experimental or observational (Figure 1.1). Remember that these labels are not exclusive and so it is possible to have, for example, a pragmatic, qualitative, experimental trial. The labels simply allow us to begin to order our thoughts about what has gone on within the trial.

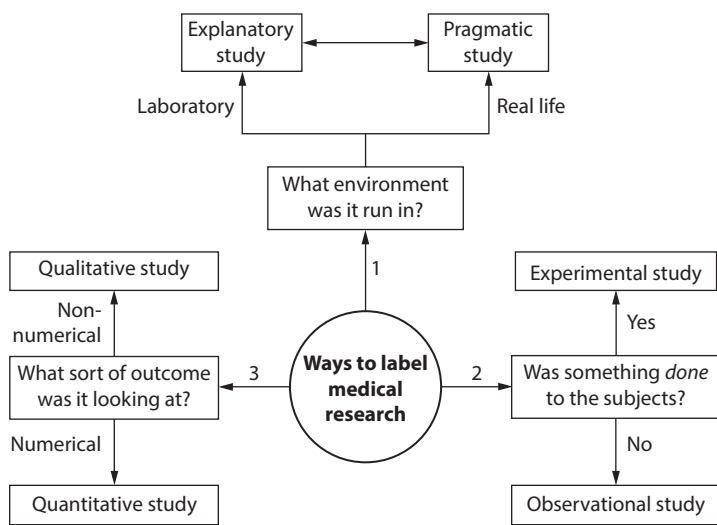


FIGURE 1.1 Different ways to label a trial – three questions to ask.

QUESTION 1: IN WHAT SORT OF ENVIRONMENT WAS THE TRIAL RUN?

As a rough rule clinical trials can be run in two sorts of places: highly controlled settings such as research laboratories (often with carefully selected, specific patients) or in actual clinical settings that are more prone to the normal variations in atmosphere, patient, time of the day and staffing levels. Explanatory research is conducted in the more highly controlled environment and gives the best chance of seeing if the intervention actually works if as many as possible outside influencing factors are taken out of the equation. Pragmatic research, however, is performed in a setting more like a clinical environment (or actually *in* a clinical environment) and gives a better idea of whether the intervention will work in the real world. Clearly there are a variety of settings in which, and patient types on whom, studies can be performed. It is wise

to consider the type of study as being on a spectrum between the explanatory and pragmatic; the more tightly controlled the conditions of the study the closer to being explanatory and vice versa.

QUESTION 2: IN WHAT SORT OF RESULT ARE THE RESEARCHERS INTERESTED?

Whether the result that is generated is a hard clinical number such as blood pressure, height or mortality rate or is instead something less quantifiable like sense of wellbeing or happiness is our second way of dividing up research. Quantitative research will generate numbers such as difference in blood pressure or mortality which can then be analyzed to produce the statistical results that we talk about later in the book and which form much of the published results from medical literature. Studies that produce a result given as a description of feelings (happiness, comfort, etc.) and looking at our changes in opinions are termed qualitative research and can be harder to succinctly analyze and summarize. In many studies looking at what you might think of as quantitative outcomes (quality of life being a good example), scoring systems have been developed so that a quantitative result can be recorded (and hence analyzed).

QUESTION 3: ARE THE RESEARCHERS PLANNING TO DO SOMETHING TO THE PATIENT?

The final way for us to define the type of research is by whether or not the researchers intervene in the patients' care in any way. In RCTs and in diagnostic studies (both of which we talk about in detail later) but also in other types of study we take a group of subjects or patients and apply a new treatment, intervention or diagnostic test to their care. These sorts of studies are termed experimental studies. Studies where the researchers look at what happens to a group of subjects or patients but the care that they receive is influenced only by normal practice (and normal variations in practice) and not by the researchers are called observational studies. Surveys, case reports, cohort studies, case-controlled studies and cross-sectional studies are all observational studies whereas RCTs and diagnostic studies are both experimental studies.

There are a number of different kinds of each type of study, whether experimental/observational, pragmatic/explanatory or qualitative/quantitative. A description of the main features of the important types of studies can be found in the Glossary (Appendix B) and an overview is given in Tables 1.1 and 1.2. The type of study that the researchers perform should be decided before they start and should be based on the strengths and weaknesses of each type of possible study.

Finally, we need to mention secondary research. Although single papers can be highly valuable, often greater value can be derived from combining the results or knowledge from a number of papers. Results of different studies can be brought together in review articles, systematic reviews or meta-analyses, all of which we discuss later in the book.

TABLE 1.1
Different types of experimental studies

Randomized controlled study (RCT)	<ul style="list-style-type: none">• The study group is divided into two or more groups by a randomization process.• One group is given one intervention, and the other group(s) another.• The outcomes in the different groups are then compared.• The randomization works to balance and thereby remove the influence of confounding factors.• The RCT is considered to be the best sort of experimental study and commonly is given in the FCEM examination.• They are however normally expensive to run.
Diagnostic study	<ul style="list-style-type: none">• The diagnostic test is performed on the whole study group alongside the gold standard test.• The performance of the diagnostic test under investigation is then compared to the gold standard test.• Diagnostic studies appear less commonly than RCTs in the literature but appear commonly in the FCEM examination.
Crossover study	<ul style="list-style-type: none">• The study is designed to compare two interventions (usually drugs).• Instead of one group receiving one drug and the other group the second drug, both groups receive one drug for a period of time followed by the alternate drug for a period of time.• The difference between the two groups is the order in which they receive the two drugs.• Essentially each patient is acting as their own control.• Importantly the drugs being studied must not interfere with the process of the underlying disease and the disease must be unlikely to alter in nature during the time course of the trial.
Cluster study	<ul style="list-style-type: none">• Another form of trial that compares two or more interventions.• Instead of individual patients being randomized to the different arms of the study, groups of patients are randomized.• An example is all ED patients from one hospital being randomized to having tea and coffee provided free of charge in the waiting area, and all ED patients in another hospital not receiving that intervention.

GRADES OF EVIDENCE

The main aim of the EBM is to deliver study results that further our understanding about how to diagnose and treat patients. As thousands of studies are published every week, we need a way to judge both the importance of the result and the strength we should ascribe to the recommendations given in the paper’s conclusion (or indeed to any guideline produced by an august professional body). Although there are variations in how different people rate different types of study and how strongly they make their recommendations, most follow a similar pattern to that shown in Tables 1.3 and 1.4.

Often people find the ranking of evidence and recommendations a little confusing. Although it’s fairly standard that an RCT is considered of more value than a cohort study, and the grading system reflects that, remember it is only a tool to help us when summarizing the evidence that we’re reading. To give one famous example, the British

TABLE 1.2
Different types of observational studies

Case control studies	<ul style="list-style-type: none">• A group of patients with an outcome variable (like the presence of a particular disease) are compared to a group without that outcome.• The medical and social histories of both groups can be compared to look for preceding risk factors that may have made the patients with the outcome more susceptible.• Case control studies are by definition retrospective and therefore suffer from recall bias.• They are good for looking at patients with rare diseases.• They are not very good at looking at exposures to rare risk factors.
Cohort study	<ul style="list-style-type: none">• A group of patients with an exposure to an intervention or risk factor are followed in time and compared to those without the exposure to that intervention or risk factor.• The researchers do not alter what exposure any patient gets and the variation between the two groups is just due to normal differences in care or exposure.• The outcomes of both groups are then compared.• Although these are normally prospective, if a large database of patients has been created for something else, the database can be used retrospectively to perform cohort studies.• Cohort studies are valuable as they allow us to look at a number of different outcomes in relation to one exposure and the time sequence of events can be assessed.• Problems include being potentially costly to run and that you need a large sample size if you want to look at rare outcomes.
Cross-sectional studies	<ul style="list-style-type: none">• Disease and exposure status of a population are studied at the same time.• They are good for establishing prevalence or association.• Unfortunately a large number of subjects is normally needed and they cannot be used to demonstrate causation.

Doctors Study was a cohort study performed in the 1950s that proved to be a major factor in demonstrating the link between smoking and lung cancer. If you applied Tables 1.3 and 1.4 literally it would only be considered level 3 evidence and a grade B recommendation not to smoke. From this example you can see that there are situations where an RCT is neither the best nor only way to provide a world changing result.

TABLE 1.3
Ranking the types of study

Grade	Type of paper
1	Systematic review (including at least one RCT) or a single good quality RCT
2	Studies without randomization
3	Well-designed, controlled observational studies such as cohort studies or case control studies
4	Observational studies such as case series or case reports
5	Expert opinion

TABLE 1.4
Ranking the strength of the recommendation made from the EBM available

Grade	Usual basis for strength of recommendation
A	Based upon meta-analyses, systematic reviews or randomized controlled trial
B	Based upon level 2 or 3 evidence (see above)
C	Based on level 4 evidence
D	Based on expert advice

TAKE HOME MESSAGE

1. Studies can be described as explanatory or pragmatic depending on whether the trial was run in highly controlled conditions or real life conditions; experimental or observational depending on whether the researchers intervene in the care of the patients; and qualitative or quantitative depending on whether the results are the patients’ descriptions of feelings or numerical data.
2. There are a number of different study designs for both observational and experimental studies and which ones the researchers choose will depend on what they are studying and the type of result they are looking for.
3. The different types of papers are graded by how valuable within the EBM world they are, with meta-analyses and RCTs at the top.
4. The strength of recommendation that is made by groups such as NICE or Royal Colleges based on the available EBM is usually also graded, with recommendations made on the basis of results from meta-analyses and RCTs considered the strongest.

2 Understanding basic statistics

An in-depth knowledge of statistics is not needed for you to be able to read a paper and decide if you think it's any good. However, to be able to justify your opinion and to be sure that you're not rejecting or accepting the paper's conclusions incorrectly a basic grounding in general statistics *is* needed and, over the next three chapters we endeavour to give you that. In some cases we expand beyond what we feel is core or 'pass/fail' knowledge, but this will be covered in 'Spod's Corner' rather than in the main text. This information will add to your understanding of the subject as a whole and so is valuable, but don't worry if you struggle with it or don't want to commit it to memory. If you don't feel up to 'Spod's Corner' at all, you are quite welcome to move on to the next part in the main text. In this chapter we cover three main concepts: types of data, ways to describe the data including types of distribution and measures of spread, and using sample data to make inferences about whole populations.

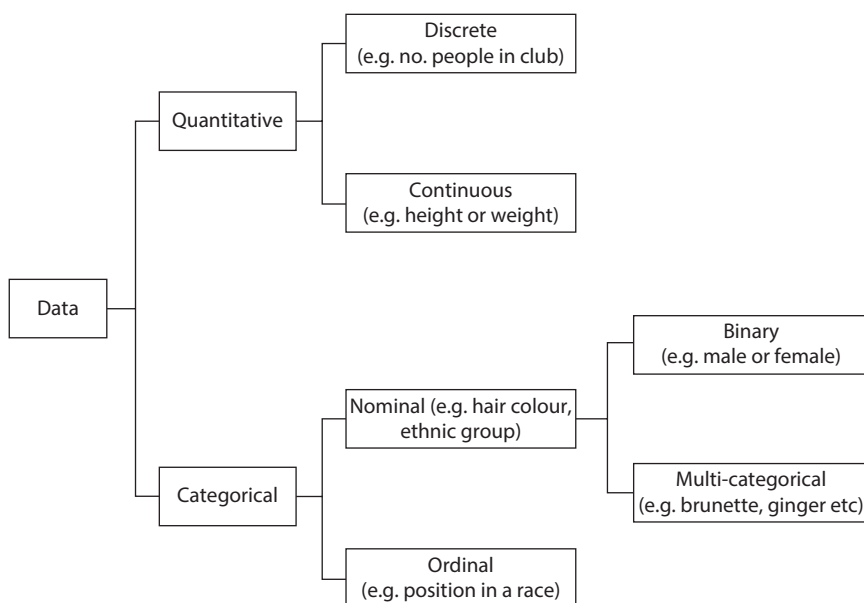


FIGURE 2.1 The main types of data.

TYPES OF DATA

The aim of most studies is to collect data from investigations performed on a sample of the population and then analyze that data to come to a conclusion. Obviously the type of data that is collected during a study will vary depending on what is being studied – for example when you’re looking at hair colour the results you get are fundamentally and instinctively different from the type of result you get when looking at weight or height. From a Critical Appraisal point of view this is important because the way in which data about hair colour is analyzed requires different statistical tests than data about weight or height. For the FCEM, in-depth knowledge about what statistical tests are done is not essential (although we mention this later in a ‘Spod’s Corner’) but knowledge about how different data is described is needed. We find the flow chart in Figure 2.1 to be the easiest and best way to demonstrate the different ways to describe groups of data.

Quantitative data is where a numerical value can be given to the data and this can be either discrete (i.e. with the number of people in a club it can be 3 or 4 but not 3.5) or continuous, where it can be anywhere on a scale. Categorical data (also called qualitative data) is where there is a fixed number of defined categories into which data is organized and this can either be nominal (i.e. given a name) such as hair colour or ordinal (i.e. where there is a rank to the categories) such as first, second and third place in a race.

For completeness and while we’re talking about data it seems reasonable to add a description of the accepted definitions of various types of epidemiological data (Table 2.1). As ever, we would strongly recommend that an understanding of what each term means rather than rote learning of a definition is the most valuable approach.

TABLE 2.1
Epidemiological data

Term	Explanation
Incidence	The rate that new cases appear over a period of time divided by the population size
Standardized (mortality) rate	The (mortality) rate where an adjustment has been made to compensate for confounding factors
Point prevalence	The number of people with the disease at a given point in time divided by the size of the population at that time
Period prevalence	The number of people with the disease in a period of time divided by the size of the population during this period
Lifetime prevalence	The proportion of a population that either has or has had the disease up to the point of the study

DESCRIBING THE DATA – QUANTITATIVE DATA

A mass of values collected as part of a study need to be summarized or described in a way that the reader can easily understand and use. It is pointless to simply list each individual value and expect the reader to interpret what it all means – most of us

need the help of the author and a statistician to make heads or tails of it all. To help explain the data, there are various things that may be done including calculating an average, describing how the data is distributed and describing the spread of the data.

DESCRIBING THE AVERAGE VALUE

An average is a single value that incorporates all the values and so can be helpful in representing the data set in a quick and simple way (Table 2.2).

TABLE 2.2
Averages

Type of average	Description	Advantages	Disadvantages
Mean	The sum of all the observations divided by the number of observations	Easy to calculate and understand	Influenced by outliers so less useful with skewed data
Median	The middle number of the group when they are ranked in order; the mean of the two middle values is taken if there is an even number of values	Reduces the effect of outliers and so is useful in skewed data	Exact 95% confidence intervals of the median require specialist knowledge and software to calculate
Mode	The most frequently occurring value in a list	Can be used with categorical data	Less intuitively useful

DESCRIBING THE DISTRIBUTION OF THE DATA

If we were to plot all the data from our study on a graph it might fall into one of a number of shapes, called distributions. Knowledge of distributions is important because it affects which average we use (mean or median) and which measure of spread (inter-quartile range or standard deviation). It also influences what statistical tests we can do on the data.

A number of names are given to various patterns of distribution but of all the ways that data can be distributed, the ones you must know are the last three shown in Figure 2.2.

Normal distribution

Figure 2.3 is without doubt the most important type of distribution and is referred to as ‘normal’, ‘Bell-shaped’ or ‘Gaussian’. Due to its fundamental importance in medical statistics and its value to us, we will spend some time discussing it before moving on briefly to the two other types of distribution you should be aware of. A normal distribution has a number of important mathematical characteristics, the first of which is that the mean, median and mode all lie at the same point.

Having a data set that is normally distributed allows certain types of statistical test to be performed on the data such as t-tests and analysis of variance. These are

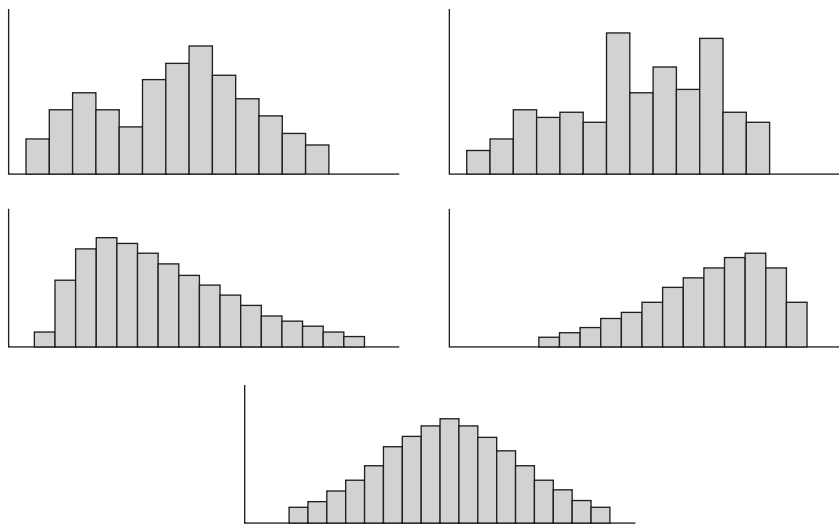


FIGURE 2.2 Various types of distribution.

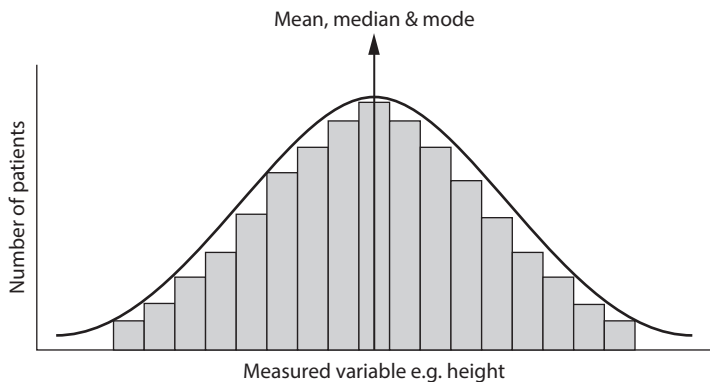


FIGURE 2.3 Normal distribution.

called parametric tests and they are preferred over non-parametric because they are more robust, or powerful (meaning you need less people in the study). The decision about whether data is actually normally distributed can be done visually by plotting the data on a chart and checking to see if it looks like a bell curve. In small studies this can be unreliable but one can quickly see whether the data is symmetrical by checking if the mean is equal to the median. There are more formal ways to confirm the normality of your data set but they are beyond the scope of this book. What is interesting is that the importance of having normally distributed data is such that when faced with a data set that is not normally distributed, many authors will convert their data to a normal distribution through mathematical manipulation so that it can be analyzed with parametric statistical tests.

There are two further important features of normal distributions that allow us to estimate what the result we would get in the whole population would be, based on the results we obtain from a single sample. We discuss these in the next section of this chapter.

Skewed data

Normally distributed data will be symmetrically distributed about the middle (remember: mean = median = mode). Some data sets may instead have their values more biased to the left or right of the graph. These are described as skewed data sets and the long tail of the graph indicates whether the data is skewed to the right (positive) (Figure 2.4a) or to the left (negative) (Figure 2.4b). Unlike normally distributed data, the mean, mode and median do not lie at the same point and this is why we often decide that the mean is not the most appropriate average to use when talking about skewed data.

It is more powerful to analyze normally distributed data than skewed data, and analysts will transform their skewed data into a normal distribution using mathematical tools such as logarithmic transformation.

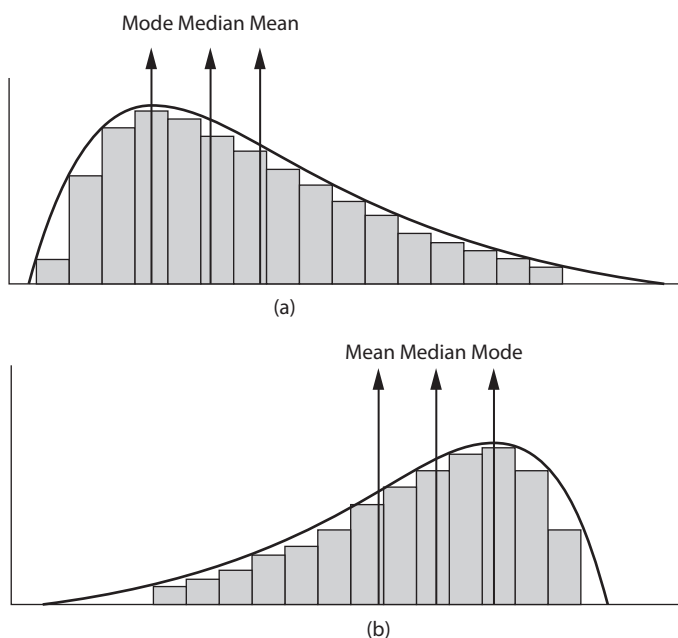


FIGURE 2.4 Skewed data. (a) Right (positive) skewed; (b) left (negative) skewed.

DESCRIPTIONS OF HOW SPREAD OUT THE DATA IS

As we've seen, the shape of the spread of the data is important to us but so is how spread out the data is. An average, on its own, does not give sufficient information to get a real feel for the data set and so an average should be combined with a 'measure

of spread’ to allow a better understanding of the data. Which measure of spread and average to use depends on how the data is distributed, i.e. whether it is normally distributed or skewed (Table 2.3).

TABLE 2.3
Measures of spread and location for normal (or symmetric) and skewed data

Distribution	Measure of spread	Explanation	Average used
Normal or symmetric	Standard deviation	A value to describe the average variation of each value in a data set from the mean of that set	Mean
Skewed	Inter-quartile range (IQR)	The range between the bottom quarter and the upper quarter of the data set; it is therefore the range of the middle 50% of the data set	Median

The importance of looking at the spread or variance of a pair of samples is that when we try to judge whether there is a genuine difference between the two groups, variance is one of three factors that we must consider. The other two are sample size and the actual difference between the mean of the two groups. We come back to this point when we talk about the power of a study in Chapter 3.

An inquisitive reader might now ask for more information about standard deviations, and this is given in ‘Spod’s Corner’. Remember that you don’t necessarily need to learn what’s in the ‘Spod’s Corner’ but an appreciation of the points will help your overall understanding. For those who really struggle with the maths, all you really need to know at the moment is that standard deviation is a measure of the spread of the data and it can be calculated from the data set.

SPOD’S CORNER

Descriptions of the variance of the data

Consider two lists of weights (kilograms) of patients from two emergency departments:

Department A: 77.5 76.3 75.9 76.9 75.4 75.3 77.0 76.4 77.2 75.9
Department B: 56.4 70.3 93.2 45.6 83.9 78.9 88.6 99.6 51.8 95.5

If we want to compare the two groups it’s reasonable to calculate the mean for each group to see if there is any difference between them – we do this by adding each weight and then dividing by the number of subjects:

Department A: Mean = $763.8 \div 10 = 76.38$ kg
Department B: Mean = $763.8 \div 10 = 76.38$ kg

So the means are the same but it seems obvious that to say the two groups are the same is incorrect and this is highlighted when we list the weights in order of size:

Department A:	75.3	75.4	75.9	75.9	76.3	76.4	76.9	77.0	77.2	77.5
Department B:	45.6	51.8	56.4	70.3	78.9	83.9	88.6	93.2	98.5	99.6

Clearly, the results from Department A are less spread out than those from Department B. We can give a measure of this by calculating the standard deviation of each sample. Technically we first calculate what is called the variance and from that we calculate the standard deviation.

Calculating Variance

1. Calculate the difference between each data point and the mean:
e.g. $(x_1 - \bar{x}) + (x_2 - \bar{x}) + (x_3 - \bar{x}) \dots$
2. Sort out the problem that if we add those all together the answer will be zero as half the results will be positive and half negative!
e.g. $(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + (x_3 - \bar{x})^2 \dots$
3. Divide by the number in the sample – 1:

$$[\sum (x_1 - \bar{x})^2] / (n - 1)$$

Calculating the Standard Deviation

$$\text{Standard Deviation} = \sqrt{\text{variance}} = \sqrt{\frac{\sum (x_1 - \bar{x})^2}{n - 1}}$$

Key – what do all these symbols mean?

The symbols in the equations are mathematical shorthand and are supposed to make things easier! For clarity:

x = any bit of data

x_1 = first bit of data, x_2 = second bit of data, and so on

Σ = ‘sum of’ (i.e. add everything together)

\bar{x} = mean of all the data

n = total number of pieces of data

$\sqrt{}$ = square root

EXAMPLE

A WORKED EXAMPLE:

So let's first calculate the standard deviation of the weights of Department A.

1. Calculate the variance:

$$\begin{aligned} & (75.3 - 76.38)^2 + (75.4 - 76.38)^2 + (75.9 - 76.38)^2 + (75.9 - 76.38)^2 + \\ & (76.3 - 76.38)^2 + (76.4 - 76.38)^2 + (76.9 - 76.38)^2 + (77.0 - 76.38)^2 + \\ & (77.2 - 76.38)^2 + (77.5 - 76.38)^2 \end{aligned}$$

$$(10 - 1)$$

$$\text{Variance} = 0.58$$

2. Calculate the standard deviation:

$$\text{Standard Deviation} = \sqrt{\text{variance}}$$

$$= \sqrt{0.58}$$

$$= 0.76\text{kg}$$

If we were to perform the same calculation to find the standard deviation of the weights of Department B, we'd find it is 19.40 kg. Therefore although the means of the two groups are the same, the standard deviations of the two groups are very different thus highlighting that there is an important difference between the two groups.

INFERENCES ABOUT THE WHOLE POPULATION FROM SAMPLE DATA

It is important to understand that from a single sample (a set of observations) you can make inferences about the whole population. The population values are fixed. All we are trying to do by taking a sample is to make an estimate of the population.

If we look at the blood pressure drop achieved in a sample of patients given a new antihypertensive, we are able to measure the mean drop in blood pressure and the spread of the results in our sample but what we are really interested in is what is likely to happen in the population as a whole if we give them the new drug. This is where statistics can become our friend (really!).

MORE REASONS WHY A NORMAL DISTRIBUTION IS SO VALUABLE

The first important point to understand is that normally distributed data has one further mathematical quality which we have not yet discussed, and that is that we can

guarantee that certain proportions of the results will always fall within certain numbers of standard deviations (SD) either side of the mean (Figure 2.5). Thus it is always true that:

- 68% of the results will lie within ± 1 SD of the mean (i.e. one SD above and one below the mean)
- 95% of the results will lie within ± 1.96 SD of the mean (many people will round this up to 2)
- 99% of the results will lie within ± 2.5 SD of the mean

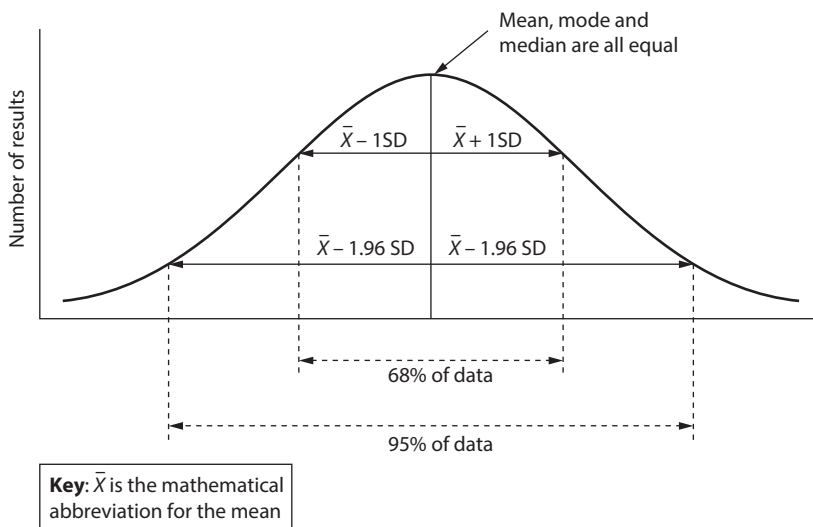


FIGURE 2.5 The proportion of data within various standard deviations of the mean of a normal distribution.

A second and equally important feature of normal distributions concerns how we use our sample data set to draw conclusions about the whole population. Let's think about a study looking at the mean height of all the emergency physicians (EPs) in the world. It is impractical to think we could actually measure every EP's height but we could look to measure the heights of a randomly selected subset of EPs and from this calculate the mean height of that sample. If there are enough people in our sample, it might reflect the whole EP world population, but given random variation it might not. However, if we were to repeat the study so that we had a number of sample groups we would be able to calculate a mean for each group, and the interesting part is that if we were to plot the means of enough sample groups these means would themselves fall into a normal distribution.

Even more interesting this is true even if the data from the whole population does not have a normal distribution (see ‘Spod’s Corner’).

HOW ACCURATE AN ESTIMATE OF THE POPULATION MEAN IS OUR SAMPLE MEAN?

Moving back to our example about blood pressure reductions – we’re not just interested in what happened in our sample group but also what is likely to occur in the whole population. Therefore we have to ask how accurately does our sample result reflect what we would see if we treated (and measured) the whole population?

Intuitively, our sample result is likely to be closer to the population value if:

- The sample is large
- The variation is small

We can also make a reliable estimate of the accuracy of our sample mean compared to the whole population by calculating something called the standard error of the mean (SEM) and using it in a similar way to how we use standard deviation. A fuller explanation of this is given in the following ‘Spod’s Corner’ but for those who want to keep it simple the following explanation should suffice: We can calculate SEM from the SD of our sample and the number of people within our sample. We can then estimate the accuracy of our sample result compared to the whole population. In a similar way that 95% of all our sample values fall within 1.96 SD of the mean, we infer that 95% of all the means in all the samples will lie within 1.96 SEM from our sample mean. From this, you can understand that we get to the idea of confidence intervals (CIs) – a range which is likely to include the true population mean and that we can calculate this from a single sample.

The flow chart in Figure 2.6 demonstrates a simplification of how we move from the result we obtain from a single sample of patients to draw conclusions about an entire population.

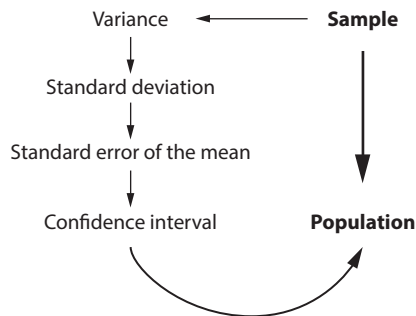


FIGURE 2.6 Making predictions about populations from a single sample.

SPOD’S CORNER

- 1. If we repeat the study a number of times and the studies are large enough ($N > 30$), we can calculate the means from each of the studies and they will themselves fall into a normal distribution (Figure 2.7).

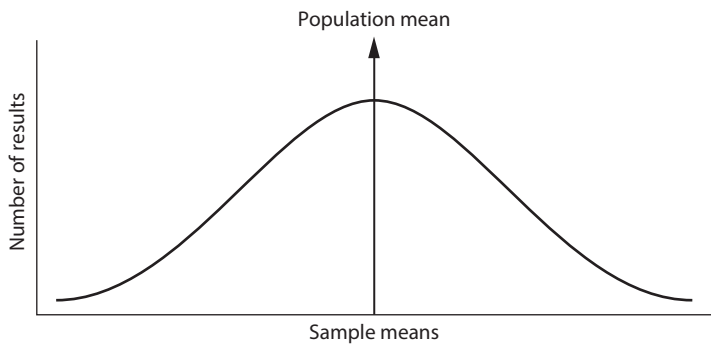


FIGURE 2.7

- 2. The same rule applies as with any normal distribution: 95% of values will fall within 1.96 SD of the mean. The difference is that when we’re moving from a sample mean to estimate a population mean we don’t call it ‘standard deviation’, we call it ‘standard error of the mean’ (Figure 2.8).

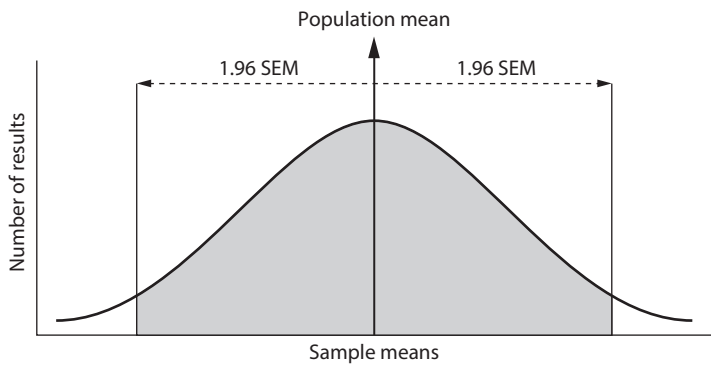


FIGURE 2.8

3. Calculating the standard error of the mean (SEM).

$$\text{Standard error of the mean} = \frac{\text{Standard deviation}}{\sqrt{n}}$$

It's worth noting that if you want to halve the size of your SEM (and hence your CI) you need to quadruple the number of people in your sample.

4. Formal definitions of 95% CIs can often be confusing – we find the following the easiest way to explain them. Imagine taking 100 equally sized samples from a population and calculating means and SEM for each. Then we can calculate for each sample the 95% CI (sample mean $\pm 1.96 \times \text{SEM}$). 95% of the CIs will contain the true population mean (Figure 2.9).

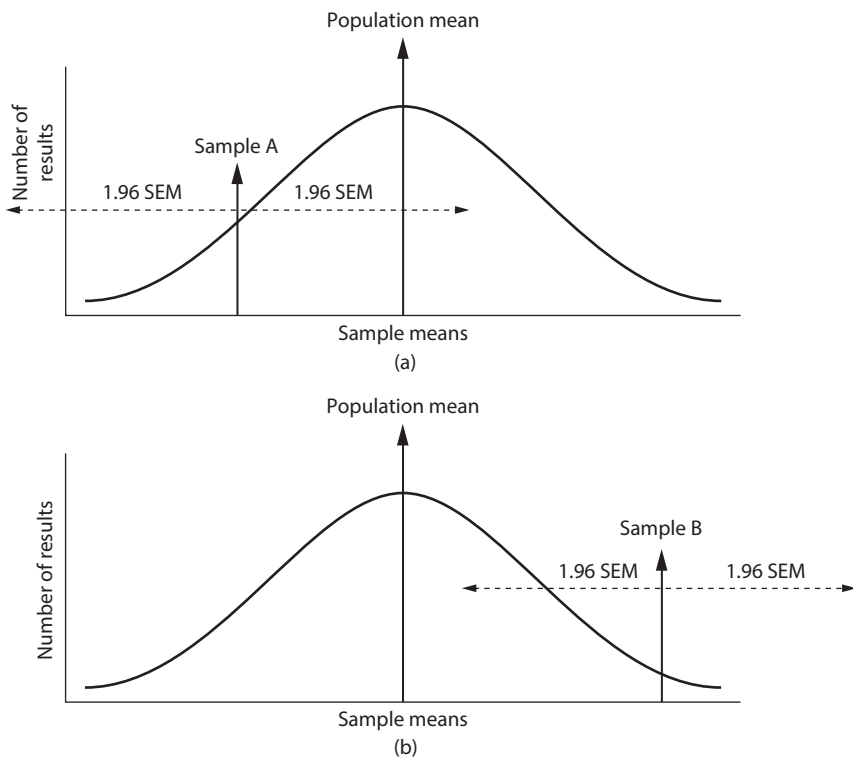


FIGURE 2.9 95% Confidence interval = sample mean $\pm 1.96 \text{ SEM}$. (a) 95% of CIs will contain the population mean but (b) 5% of CIs will not.

TAKE HOME MESSAGE

From a sample of data we can calculate the mean (a measure of location, or centrality) and the standard deviation (a measure of spread). From this standard deviation we can calculate a standard error of the mean and from this estimate a range which is likely to include the true population mean. This is called the 95% confidence interval for the population mean.

3 Statistics in papers – p-values, confidence intervals and error

Although when you read papers you might be struck by an overabundance of statistical terms, the truth is that what you actually need to know is thankfully a fairly small amount. In this chapter we cover:

1. Significance
2. Hypothesis testing, p-values and confidence intervals (CIs)
3. Error – or what if the researchers are wrong?
4. Statistical power and how many people are needed in a study
5. Looking out for error in studies
6. Statistical tests

WHAT IS SIGNIFICANCE?

We talk a lot in scientific studies about ‘significance’ but it’s important to realize there are two types of significance: statistical significance and clinical significance. Imagine a study looking at a new antihypertensive that showed a drop in blood pressure of 2 mmHg in the treatment group. We could perform a statistical analysis and might find that the result is *statistically significant*. By this we mean that we think the result demonstrates a difference between the two groups that is unlikely to have occurred by chance. However, it would be reasonable for you to ask whether, in the real world, a drop of 2 mmHg in blood pressure was of any real value, and in doing so you would be questioning the *clinical significance* of the result.

HYPOTHESIS TESTING, P-VALUES AND CONFIDENCE INTERVALS

HYPOTHESIS TESTING

When authors embark on a study they will have begun with a research hypothesis. For example we may have observed that people given aspirin at the time of their acute myocardial infarction (AMI) seem to do better than those who weren’t and so we start with a research hypothesis of ‘aspirin is better than placebo in AMI’. In superiority trials, we always set up the study so that the default position is we expect that there is no difference between two groups (even if this isn’t what we actually expect!). This is called the *null hypothesis* (H_0). The alternative hypothesis, or study hypothesis (H_1), is the statement that there is a difference between the two groups.

In our example, the null hypothesis is there is no difference in mortality between patients given aspirin or placebo for AMI. The study hypothesis is there is difference in mortality between patients taking aspirin or placebo.

We then set about running an experiment and using statistics to help us either reject or not reject the null hypothesis. The slightly awkward ‘reject or not reject’ is important because even if we reject the null hypothesis we can say we support the original hypothesis (that aspirin is good for you in AMI). If we do not reject the hypothesis we are not saying that the null hypothesis is correct, we are just saying that there is not enough evidence to say it is false. This is similar to being found ‘not guilty’ rather than ‘innocent’ in a court of law. The answer to whether we reject the null hypothesis can really only be yes or no. This is called *hypothesis testing*.

In hypothesis testing we are asking a yes or no question: Is there a difference, yes or no?

With hypothesis testing, we use statistics to generate a p-value and if it is above or below a predetermined value (by convention 5%, or 0.05) we say we can either reject or do not reject the null hypothesis.

P-VALUES

We calculate a p-value to enable us to decide if we will reject or not reject the null hypothesis. If we see a difference in mortality rates between our two AMI groups (aspirin and placebo), we need to be aware that while the difference might have occurred because of the treatment, it might also have occurred by chance. In calculating a p-value we are trying to answer the question: How likely are we to have obtained our results if we had sampled from a population where there was no difference between the groups (i.e. if the null hypothesis were true)? The higher the p-value, the more likely that any difference we saw between the groups occurred by chance.

A formal definition of a p-value can be given as:

A p-value is the probability of obtaining a result at least as extreme as the one observed, assuming the null hypothesis is true. The lower the p-value, the less likely the result is if the null hypothesis is true (and the less likely the result is to have occurred by chance).

For many people, that’s a bit of a mouthful, and so we prefer:

$p < 0.05$ means that the probability of obtaining the result by chance was less than 1 in 20. By convention, a p-value of less than 0.05 is the accepted threshold for statistical significance and is the level at which the null hypothesis can be rejected.

EXAMPLE

We want to know if a new painkiller is effective in minor injuries, and we run a study to compare the pain scores in 100 people given the new painkiller versus 100 given just paracetamol. Each patient has a visual analogue pain score on

arrival and one hour post-analgesia. We can calculate the drop in pain score for each patient, average this out for each group and arrive at a mean pain reduction score for each group: 5 in the treatment group and 2 in the paracetamol group. We can then use a statistical test to calculate a p-value which in this example comes out as $p = 0.01$.

Because the p-value is less than 0.05 we can say that the result is statistically significant.

CONFIDENCE INTERVALS

Although a simple yes or no answer is sometimes attractive to busy clinicians, in most situations we prefer a bit more information. Many studies are set up so that the information they provide is a measurement. In our example of aspirin in AMI, we might be looking at a mortality rate that could be quoted as a percentage. In that example it's naturally of interest to us what the actual mortality rate in each group was. The result of one sample can only give us an estimation of the result we would get if we were to look at the whole population. However we can use the results from a single sample, and with statistical testing estimate how accurate our result is in comparison to the whole population.

A CI provides the reader with a range of plausible values for the true population value. A 95% confidence is the standard size CI used in the medical literature when reporting results. In the same way, a value of less than 0.05 is taken as the threshold of significance when using p-values in hypothesis testing.

SPOD'S CORNER

Although 95% CIs are commonly thought of as representing the range that has a 95% chance of containing the true (overall population) value, this is sadly not statistically strictly true (sadly, because this is easy to understand!). In fact a 95% CI means that if we were to repeat our experiment ad infinitum and were to construct 95% CI from each of our sample values, 95% of those intervals would contain the true population value.

It is self-evident and logical that the narrower the CI the more precise our sample result is in comparison to the overall population result and thus the more accurate our result.

Confidence intervals can be used in a number of ways:

- Assessing the clinical significance of a result.
 - If the mean drop in blood pressure seen with a new antihypertensive is 10 mmHg, we might feel this is clinically significant. If the 95% CI

is 8–12 mmHg we would agree that even if the true population result was at the lower end (8 mmHg) this would be a clinically significant drop. However, if the 95% CI was 2–18 mmHg, we'd probably agree that if the true population result was only 2 mmHg that this would not be clinically significant. Therefore we must always interpret the mean value alongside the 95% CI to assess clinical significance.

- If we are simply comparing two groups and want to see if there is a real difference between them (i.e. did not occur by chance), we can look to see if the CIs of each group overlap. If they overlap, the result is likely not to be statistically significant.
- If we're looking at a ratio between two groups (e.g. an odds ratio or risk ratio), if the CI crosses 1 (the point of no difference) the result is not statistically significant.

P-VALUES VERSUS CONFIDENCE INTERVALS

Some studies publish both p-values and CIs, some just one and some just the other, but both p-values and CIs tell us about statistical significance; so is one better than the other? Confidence intervals give us more information; whether the result is statistically significant or not while allowing us an idea of where the true value might lie (and hence helping us to make a decision about clinical significance). We think CIs are of more use than p-values but in an ideal world, all studies would give both.

EXAMPLE

In the same analgesia study described earlier we calculate the mean drop in pain score for each group but we also calculate the 95% CI for each group. In the treatment group the reduction was 5 with a 95% CI from 4.8 to 5.2 (written as '95% CI 4.8–5.2') and in the paracetamol group 2 (95% CI 1.8–2.8). As the CIs in each group do not overlap, we can be confident that the difference between the two groups is real i.e. statistically significant.

ERROR – OR WHAT IF THE RESEARCHERS ARE WRONG?

By convention we say a result is statistically significant if the probability it occurred by chance is less than 5%, and we write this as ' $p < 0.05$ '. (Ideally the exact value should be given.) It's possible, although unlikely, that a result we take to be statistically significant has in fact occurred by chance. Equally, and for a number of reasons we will later discuss, the results from our study might *not* be statistically significant

when in fact there is a difference between the two groups. If any of these happen we would describe it as an ‘error’ which we divide into what we call type 1 and type 2 errors (Table 3.1).

TABLE 3.1
Type 1 and type 2 error

	Alternative hypothesis is true (there is a difference between the groups) +ve	Null hypothesis is true (there is no difference between the groups) -ve
Experiment shows a significant result +ve	True positive NO ERROR	False positive TYPE 1 ERROR
Experiment shows no significant result -ve	False negative TYPE 2 ERROR	True negative NO ERROR

TYPE 1 ERROR

Making a type 1 error is saying there is a difference between two groups when in fact there isn’t one. Making this error might lead you to recommend a treatment to patients that is of no benefit to them.

In statistics, we say the probability of making a type 1 error is called alpha (α). We don’t actually know the probability of making a type 1 error before we start our experiment but we can state what we consider to be an acceptable probability of making a type 1 error. This is usually taken to be 5% (or in mathematical terms, 0.05). Alpha is also called the ‘threshold for statistical significance’.

It is important we set α before we start our experiment (usually at 0.05) but calculate our p-value once we have our results (and say the result is statistically significant if it is below 0.05). The lower we set α to start with, the lower our p-value needs to be to be statistically significant and the lower the chance our result is a false positive.

TYPE 2 ERROR

Making a type 2 error is saying there is no difference between two groups when in fact one is present. This might lead us to not providing a treatment to our patients that in fact would be of benefit to them. Usually this is due to either the variance in the sample being too great or the sample size being too small (see Chapter 2).

In statistics, we label the probability of making a type 2 error beta (β). For a fixed sample size, we can consider α and β to be balanced on either end of a set of scales – as we set α (the chance of a false positive result) lower, we make β (the chance of a false negative result) larger. The reason that α is usually set at 0.05 is that conventionally this is taken as a reasonable balancing point. All clinical trials are a balance of the risk of a type 1 error, a type 2 error and sample size.

STATISTICAL POWER AND HOW MANY PEOPLE DO I NEED IN MY STUDY?

STATISTICAL POWER

We can describe how likely a study is to detect a difference if one truly exists by subtracting the probability of a type 2 error (a false negative) from 1. We call this ‘statistical power’ and as a formula we would write it as $1 - \beta$.

By convention the power of a study should be 80–90%, i.e. we allow ourselves a 10–20% chance of not demonstrating a difference if one exists.

Determinants of statistical power (probability of detecting an effect when it truly exists)

1. Alpha – the threshold of significance. Most commonly this is 0.05 and corresponds with accepting a p-value of 0.05 as the cutoff for statistical significance.
2. The variability of the data (measured by the standard deviation).
3. The minimum difference considered to be clinically significant.
4. The sample size.

SAMPLE SIZE

A major part of designing any study is working out how many people you need to enter into the trial. If we set the statistical power we’re aiming for, which by convention is 80–90%, we can calculate the sample size we need by knowing the other aforementioned three determinants. Alpha is usually also set by convention to be 0.05. The variability of the data is a fundamental feature of the data set, and although obviously you can’t know the variability of your data before you start you can use previous studies or pilot studies to get a good estimation of what it’s likely to be. The minimum clinically significant difference (or minimum clinically important difference, MCID) should be stated before starting a study and ideally should be taken from an independent source such as another research study. It is important that the MCID is a clinical threshold and is therefore determined by the clinician setting up the study and should be explicitly stated.

EXAMPLE

In our study looking at the effect of our new painkiller, we need to know how many people we are required to recruit *before* we start the study. We start by asking: What chance of a type 2 error am I willing to accept? Let’s say 10% (by convention we’d usually accept somewhere between 10% and 20%). That means we’re aiming for a power of 90% ($1 - \beta = 1 - 0.1 = 0.9$ [90%]). We then set our type 1 error risk, most commonly set at 5% (i.e. $\alpha = 0.05$; if we find a p-value less than 0.05 we’ll say the result is statistically significant). Additionally we decide (perhaps from

other studies) what is considered to be the minimum reduction in pain score that is clinically significant, let's say a drop of 1.0 on the visual analogue scale we're using. We may also have obtained an estimation of data variability from previous studies.

Finally we perform a *power calculation* to determine how many people we need to enrol into our study. Exactly how the power calculation is performed is thankfully beyond the scope of the FCEM examination but you should understand that the smaller the sample size, the lower your power will be and the higher the chance of a type 2 error (false negative).

Good quality studies should describe their sample size/power calculation together with all parameters and constants so an independent reviewer can replicate the results.

LOOKING OUT FOR ERROR IN STUDIES

SPOTTING TYPE 1 ERRORS

We've described that the chance of a type 1 error (a false positive) corresponds to the p-value level we set for accepting statistical significance. Basing your practice on a positive study result that is in fact a false positive (type 1 error) may result in you giving a treatment to a patient which is of no benefit and may even cause harm. Often type 1 errors occur due to a flaw with the study, either bias or poor methodology.

If we set α at 0.05, we're saying that we're willing to accept a 5% chance of a false positive and saying that if the p-value we calculate is under 0.05 we'll accept that the result is statistically significant. It's really important to note that this is only true if we're talking about looking at one single result. The more results we look at within one trial, and hence the more p-values we calculate, the greater the chance that one of them will be statistically significant purely by chance. That is to say, the more statistical tests you perform on one sample, the more likely you are to get a statistically significant result; or the wider the net you cast on the data, the bigger chance you get of catching a significant result; or if you overanalyse your data set, you increase your chance of getting a false positive result.

EXAMPLE

Consider running our analgesia study again. Let's keep α at 0.05 and hence the p-value we'll call statistically significant at $p < 0.05$, but this time instead of only quoting a result for reduction in pain we'll also give results for the following: total number of doses taken, days off work, days until pain free, days requiring analgesia, time spent in the ED, re-attendance rates at the ED, attendance rates at the GP, length of time needing a splint and overall satisfaction.

When we had only one result, the probability that we might find a statistically significant result by chance was 5%, but now because we have 10 results, the probability that one of them will be statistically significant by chance alone jumps to 40%.

SPOD'S CORNER

Probability of at least one false negative = $1 - (\text{probability of no false negatives})$

Probability of no false negatives = $1 - (\text{probability of a false negative}) = 1 - 0.05 = 0.95$

Therefore:

Probability of one false negative in 10 tests is $1 - (0.95)^{10} = 1 - 0.6 = 40\%$

(If multiple tests are done on a single data set, each result is not necessarily independent of the others and hence the calculation is a little more complicated than above but the general rule holds true.)

Looking for multiple results from one set of data has been termed *data dredging* and is surprisingly common. It is often seen in observational studies and database studies and may reflect a poorly planned study. Equally the risk of a type 1 error should be particularly watched for when authors report positive results from a subgroup analysis of a study. It is possible to power your study to allow for subgroup analyses and equally statistical correction can be done (e.g. the Bonferroni correction) to account for doing a number of tests on a data set. If neither of these has been done you should be very wary of authors looking for a positive subgroup result from a trial that overall did not show a statistically significant difference.

QUESTIONS TO ASK TO AVOID MISSING A TYPE 1 ERROR

- Are there multiple hypothesis tests (are the authors asking more than one question)?
- Is this a subgroup analysis? (In ISIS-2, subdividing patients by astrological sign suggested aspirin had a slightly adverse effect on mortality in Gemini patients but was beneficial in all other star signs!)
- Do the authors report and discuss the results they said they were going to look for in the design of the study?
- Does the result make sense in the context of other known results?

SPOTTING TYPE 2 ERRORS

Type 2 errors are what occur when we fail to find a statistically significant result when a real difference was actually present. When a type 2 error occurs it is usually because the sample size was too small.

In studies that reported a non-statistically significant result it is easy to dismiss the study and conclude that the intervention that was studied is of no value. This risks

missing the fact that valuable information can be gained from assessing whether a type 2 error may have caused a negative result and whether in fact you should still keep an open mind about the treatment or diagnostic test the authors looked at.

QUESTIONS TO ASK TO AVOID MISSING A TYPE 2 ERROR

- How wide are the CIs? (If they are very wide this suggests the sample size was too small and hence there might be a type 2 error.)
- Do the CIs encompass a potentially clinically important difference?
- Did the authors do a power calculation and were all the variables included appropriately calculated/estimated?
- Was the minimum clinically important difference the authors set actually too big?

TYPES OF STATISTICAL TESTS

Statistical tests are used mainly to make inferences about a population from which the sample was taken. We can use statistics to:

- Look at the differences between samples
- Look for equivalence and non-inferiority
- Look for correlation between groups

To decide which statistical test is the most appropriate you need the following information:

- How many groups do you need to compare?
- What type of data do you have?
- What type of distribution does the data produce?
- Do you need to adjust for confounders?

Once you have this information you can use Figure 3.1 to figure out the appropriate test. These tests can then be used to calculate a p-value. For the FCEM examination it is important to understand what a p-value is and where it is of use but you don't need to know how you calculate that value. Knowledge of the statistical test may, however, be of value when writing your CTR.

Statistical tests can be used to assess relationships. We call this correlation and regression. Correlation assesses the strength of the relationship between two variables. The correlation coefficient is the measurement of this relationship, although it does not measure causation. Tests for correlation include Pearson's, Spearman's rank and Kendall's.

R value	Degree of Correlation
-1	Negative correlation
0	No correlation
1	Perfect correlation

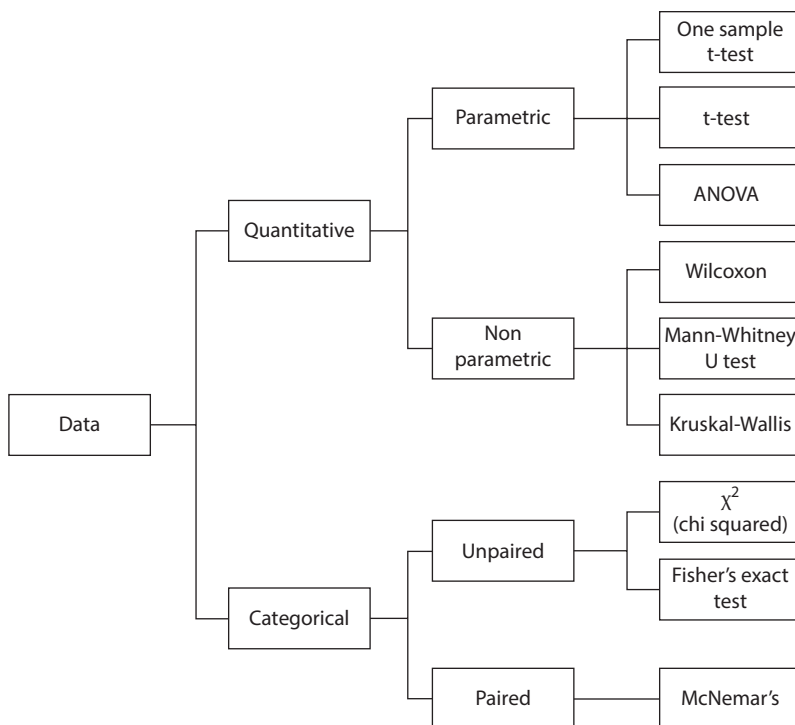


FIGURE 3.1 Statistical tests depending on data type.

When you have a relationship between two things, you can plot this relationship on a graph. Regression determines the nature of the relationship between the two variables. It is possible that multiple factors can affect the relationship and in that case multiple regression analysis is used.

TAKE HOME MESSAGE

1. The difference between p-values and CIs and the fact we think CIs are more valuable.

p-value

Needs a statement at the start of the study about what the authors think is *clinically significant*

Gives a yes or no answer to whether the result is significant

By convention a p-value of less than 0.05 is taken to be the cutoff at which we say a result is statistically significant

Confidence Interval

Allows the reader to interpret the size of the result and helps determine whether it is *clinically significant*

Gives a plausible range for the true population value

By convention a 95% confidence is the standard size CI used in the medical literature when reporting results

2. Type 1 error = false positive rate; type 2 error = false negative rate.
3. The necessary sample size of a study is calculated from: the power of the study (set before you start the study and usually 80–90%), the predicted variability of the data and the minimum difference that is deemed clinically significant. In a good quality scientific paper the power calculation should be described by the authors with sufficient details so that the reader can replicate it.

4 Probability and odds

It may seem strange to have a chapter on A-level mathematics in a medical textbook, but it is important to understand probability and odds when reading the results section of a paper. We understand that for some people this part of critical appraisal can be challenging and off-putting but please stick with it. We have kept it as brief as possible and with the clinical significance highlighted.

PROBABILITY

Unless you like betting on the horses or dogs probabilities are much more intuitive than odds.

- Probability is the likelihood of any event happening as a proportion of the total number of possibilities. Or put more simply:

Probability = Number of desired events \div Total possible number of events

- The probability of something which is certain to happen is 1.
- The probability of an impossible event is 0.
- Probability can easily be turned into a percentage by multiplying by 100, therefore a probability of 1 can also be described as a probability of 100%.
- Probability is a way of expressing risk or chance.

Let us give some simple examples.

Example 1: A coin

- When I toss a coin, there are two possibilities: heads or tails.
- The probability of a head is $1 \div 2$ which equals 0.5 or 50%.

Example 2: A dice

- When I roll a dice, there are six potential outcomes: 1, 2, 3, 4, 5 or 6.
- I have a 1 in 6 chance of getting a 6: $1 \div 6 = 0.167$ or 16.7%.

RISK VERSUS PROBABILITY

Risk is calculated in the same way as probability and to all intents and purposes the definition of risk is identical. So when we look at therapeutic papers we can work out the risk in the control group and the experimental group.

Experimental event rate (EER) is just another way of saying risk in the experimental group, i.e. the number of patients who had the experimental treatment and died (if that's the outcome of interest we're looking at) divided by the total number who had the experimental treatment.

Likewise, control event rate (CER) is just another way of saying risk in the control group, i.e. the number of patients who had the standard treatment and died divided by the total number who had the standard treatment.

One way of assessing the value of the experimental treatment is to compare the two ratios by dividing the EER by the CER; we call this relative risk.

$$\text{Relative risk} = \frac{\text{Experimental event rate}}{\text{Control event rate}}$$

All the other calculations required in the results section of therapeutic papers (e.g. absolute risk reduction, number needed to treat, etc.) are based on these simple calculations. This is discussed in greater detail in Chapter 6.

SPOD'S CORNER: A BIT MORE ABOUT PROBABILITIES

Probabilities can be used in three ways:

1. The probability of something not happening is 1 minus the probability that it will happen:
i.e. the probability of anything but a 6 when I roll a dice is $1 - 1/6$.

$$= 5/6 = 0.83 = 83\%$$

2. The probability of either one event or another event happening is the probability of each event happening added to each other:
i.e. probability of a 6 or a 5 is the probability of a 5 plus the probability of a 6.

$$= 1/6 + 1/6 = 1/3 = 0.33 = 33\%$$

3. The probability of an event happening and then another independent event happening = probability of each event happening multiplied together:
i.e. probability of a 6 and then another 6 is $1/6$ times $1/6$.

$$= 1/36 = 0.028 = 2.8\%$$

ODDS

So if probabilities are so easy to understand, why do we also use odds?

Before looking at why we use odds, let's get an understanding of them first. They are a lot less intuitive than probabilities (and perhaps this is why they are used by bookmakers).

Odds are calculated as:

$$\text{Odds} = \text{desired events} \div \text{undesired events}$$

- The odds of something certain to happen is infinity.
- The odds of an impossible event is 0.

Let us give some simple examples.

Example 1: A coin

- When I toss a coin, there are two possibilities: heads or tails.
- The odds of a head is 1:1.

Example 2: A dice

- When I throw a dice there are six potential outcomes: 1, 2, 3, 4, 5 or 6.
- The odds of a 6 are 1:5.

(See 'Spod's Corner: Converting probabilities to odds and vice versa' later in this chapter.)

ODDS IN THERAPEUTIC TRIALS

Odds can be used to look at outcomes in therapeutic trials. For example, the odds in the treatment group of a trial are calculated by the number of patients who received the new intervention and died (if that was the outcome of interest) divided by the number who received the new intervention and didn't die.

Similarly, the odds in the control group of a trial are calculated by the number of patients who received the control treatment (standard treatment) and died divided by the number who received the control treatment (standard treatment) and didn't die.

In exactly the same way as risk, we can compare the odds in the two groups by calculating a ratio.

$$\text{Odds ratio} = \text{Odds in treated group} \div \text{Odds in control group}$$

ODDS VERSUS PROBABILITY (RISK)

We asked earlier why we need to worry about odds when probabilities seem easier to understand. We'll start answering that by showing how odds and probabilities compare in a couple of examples.

Trial A

	Dead	Alive
New treatment	100	900
Control treatment	500	500

$EER = 100/1000 = 0.1$

Experimental odds = $100/900 = 0.11$

$CER = 500/1000 = 0.5$

Control odds = $500/500 = 1$

$RR = 0.1/0.5 = 0.2$

Odds ratio = $0.11/1 = 0.11$

Trial B

	Dead	Alive
New treatment	1	999
Control treatment	5	995

$EER = 1/1000 = 0.001$

Experimental odds = $1/999 = 0.001001$

$CER = 5/1000 = 0.005$

Control odds = $5/995 = 0.005025$

$RR = 0.001/0.005 = 0.2$

Odds ratio = $0.001001/0.005025 = 0.199$

It can be seen that when events are rare, risk and odds are similar. But as the prevalence increases, then the odds ratios and risk ratios become very different. It's easy to see that you might choose one over the other to present a particular slant to the results from your trial – this is shown in the example in the following ‘Spod’s Corner’.

SPOD’S CORNER: ODDS RATIOS APPEARING TO EXAGGERATE CLAIMS

Consider the following (fictional!) data on FCEM pass rates:

	Fail FCEM	Pass FCEM	Total
Male	33	14	47
Female	16	72	88
Total	49	86	135

It is obvious that a female candidate was more likely to pass than a male candidate but the question is, *how much more likely?* You could work out either the relative risk or the odds ratio to answer this question.

The odds ratio compares the relative odds of passing in each group. The odds of men passing were $14/33 = 0.42$ but for women $72/16 = 4.5$. The odds ratio is therefore calculated to be 10.7 ($4.5/0.42$), i.e. there is a more than ten-fold greater odds of passing for females than for males.

The relative risk (also called the risk ratio) compares the probability (risk) of passing in each group rather than the odds. For males, the probability of passing is 30% ($14/47 = 0.3$), while for females, the probability (risk) is 82% ($72/88 = 0.82$). The relative risk of passing is 2.7 ($0.82/0.3$), i.e. there is a 2.7 greater probability of passing for females than for males.

We can see there is a massive difference between the odds ratio and the relative risk and this can confuse people. Both measurements show that men were more likely to fail, but the odds ratio implies that men are much worse off than the relative risk. Which number is a fairer comparison? We would argue that as the risk ratio is the simplest to understand it should usually be used unless there is a specific reason to use the odds ratio. As we have demonstrated, be aware that the odds ratio might sometimes be used to inflate the difference.

Given that risk (probability) is intuitively more easily understood we would suggest these should be given most often in the results of a study. So why do we use (and need to understand) both? Odds essentially have ‘superior mathematical properties’, which makes them very useful.

ODDS HAVING ‘SUPERIOR MATHEMATICAL PROPERTIES’

There are three main areas where odds are superior to risk (probability) but in normal life and for the exam, you only need to be concerned with one. The other two are explained in ‘Spod’s Corner’ later in the chapter.

Calculating post-test probabilities using odds

The rest of this chapter covers the mathematics behind how we use investigations to make diagnoses in the ED: the ‘Bayesian theory of medicine’ in practice. This is an important concept to grasp in preparation for the diagnostic chapter later in the book.

You may not know exactly what Bayesian theory means – but you do it every day. When we see a patient we make an estimate of how likely they are to have a specific condition. We do this using our clinical ‘gestalt’ (or in other words the combination of our clinical skills and experience, the patient’s demographics, their history and examination combined with the population prevalence). This is the pre-test probability of a patient having the condition.

We may then do a test to help us decide if the patient has the condition in question. Unfortunately tests are not perfect and a yes-or-no test result will invariably not give us a 100% definitive answer. What the test actually does is help us decide the probability of the patient having the condition: the post-test probability.

Each test has a likelihood ratio attached to it, i.e. how much more likely is a positive test to be found in a person with the disease as opposed to a positive test in a patient without the disease. The post-test probability is the pre-test probability combined with the likelihood ratio.

There are likelihood ratios for positive tests (which range from 1 upwards – the higher it is, the better the test) and there are likelihood ratios for negative test results (which range from 0 to 1 – with the closer to zero, the better the test).

We may not realize it, but we do estimations of this calculation every single time we see a patient. While we are taking our history, we are making an estimate of the patient's pre-test probability. When we examine them we refine our probability. When we have a test result back, we use our knowledge of that test – our personal estimate of the likelihood ratio – to see how likely they are to have the condition. For example in someone who we are concerned may have a pulmonary embolism, the patient's management will depend more often than not on the CT pulmonary angiogram (CTPA) and not necessarily on the D-dimer. This is because the likelihood ratio of a CTPA is greater than that of a D-dimer for the diagnosis of a pulmonary embolism.

If we are trying to put actual numbers onto these post-test probabilities, as opposed to just estimates, then we must use odds and not probabilities. We can only calculate post-test probabilities by going via odds.

To get from pre-test probability to post-test probability, we go through a number of steps:

1. Convert pre-test probability to pre-test odds
2. Multiply pre-test odds by the likelihood ratio to calculate post-test odds
3. Convert post-test odds to post-test probability

This seems like all a bit of a faff but let us look at why we need to do it. Remember probabilities can only go from 0 to 1, but odds go from 0 to infinity. If we multiply the (positive) likelihood ratio with a probability, we may often get a number greater than 1, which in terms of probabilities makes no sense at all.

Let's look at an example:

- We're looking after a trauma patient.
- We suspect he has an 80% chance of a pneumothorax.
- We perform an ultrasound, which has a likelihood ratio of 9, i.e. we are nine times more likely to diagnose a pneumothorax by ultrasound in someone who has one as opposed to someone who doesn't. (Remember each operator will have a different likelihood ratio for subjective tests such as ultrasound, based on their skills and competencies.)
- Once you have the ultrasound result, now what is the chance of having a pneumothorax?
- Using probabilities:
 - Post-test probability = Pre-test probability combined with the likelihood ratio
 - Pre-test probability: 80%, mathematically written as 0.8
 - Likelihood ratio = 9
 - Post-test probability would be $0.8 \times 9 = 7.2$
 - In probability terms this makes no sense: remember that a probability of 1.0 (100%) means something is absolutely certain to happen. In this example the result works out as 7.2 or 720% and clearly this does not make sense.

- Using odds:
 - (Pre-test) odds = 80:20 = 4:1 (based on probability of 80%)
 - Likelihood ratio = 9
 - Post-test odds = Pre-test odds \times LR = 4:1 \times 9 = 36:1
- Odds of 36:1 clearly makes sense (suggesting a very high chance of having a pneumothorax given a positive ultrasound).

CONVERTING POST-TEST ODDS TO PROBABILITIES

We can see that if we are doing tests, we need to use odds. But the 36:1 post-test odds of having a pneumothorax are not intuitive. We need to convert it to a probability (which can be expressed as a percentage). There are two ways of doing this.

The first is by using mathematics. In the following ‘Spod’s Corner’ we show how you go from probability to odds and back again. But you must remember that it is not necessary to know these formulae for FCEM and they are there just for those who have an excessive interest in odds and probabilities.

SPOD’S CORNER: CONVERTING PROBABILITIES TO ODDS AND VICE VERSA

For converting probability to odds:

$$\text{Odds} = \text{Probability} \div (1 - \text{Probability})$$

Let’s take the example of the coin:

The probability of a head is 0.5 and the odds = $0.5/(1 - 0.5) = 0.5/0.5 = 1$, which is written as 1:1.

Let’s take the example of a dice:

The probability of a 6 is 0.1667.

The odds = $0.167/(1 - 0.167) = 0.167/0.83 = 1/5$, which is written as 1:5.

Let’s take the example of the pneumothorax pre-test probability:

The pre-test probability of a pneumothorax is 0.8. The pre-test odds of a pneumothorax is thus $0.8:0.2 = 4:1$.

For converting odds to probability:

$$\text{Probability} = \text{Odds} \div (\text{Odds} + 1)$$

Let’s take the example of the coin:

The odds of a head is 1:1 and the probability = $1/(1 + 1) = 0.5$ (50%).

Let’s take the example of a dice:

The odds of a 6 are 1:5. The probability = $1/(5 + 1) = 0.167$ (16.7%).

Let’s take the example of the pneumothorax post-test odds:

The odds are 36:1. The probability is thus $36/(36 + 1) = 0.973$ (97.3%).

CONVERTING POST-TEST ODDS TO PROBABILITY THE EASY WAY (BY JUST KNOWING THE PRE-TEST PROBABILITY AND LIKELIHOOD RATIO)

The second way of doing this – with much less effort – uses a nomogram. It very accurately calculates the post-test probability using the pre-test probability and the likelihood ratio by drawing one straight line. All the complicated sums of converting probability to odds, multiplying by likelihood ratio and then converting post-test odds to post-test probability have already been done for us. All we need to do is plot the pre-test probability and the likelihood ratio and then draw the line (Figure 4.1).

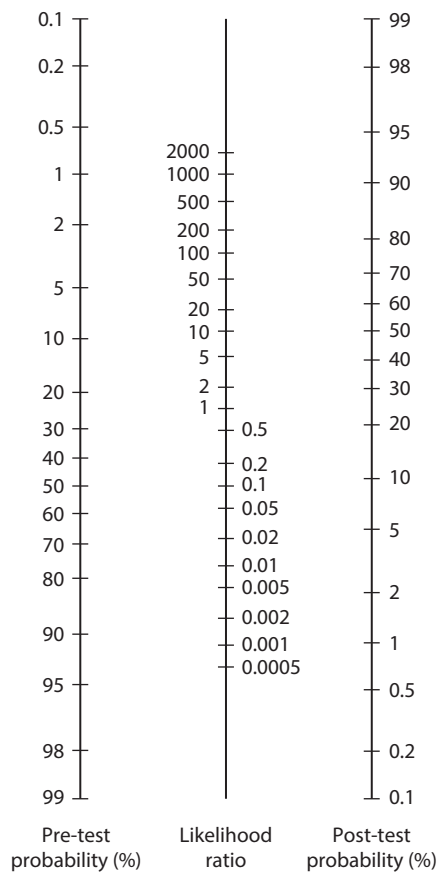


FIGURE 4.1 The nomogram for converting pre-test probabilities and likelihood ratios to post-test probabilities.

We draw a line from a point on the left-hand line, from where our clinical gestalt tells us the pre-test probability lies, through the likelihood ratio for the test and extend it on to the post-test probability. Where it cuts this line is the patient’s post-test probability.

We have to remember that the post-test probability is really an estimate, based on a best guess pre-test probability and the likelihood ratio, which if it is a subjective test such as an ultrasound, is user dependent. The beauty of the nomogram is that the same one can be used for all diseases and all tests.

Let us look at how we would use the nomogram for the pneumothorax example we previously used. As long as the line and intercepts are drawn accurately, then the post-test probability is accurate (Figure 4.2).

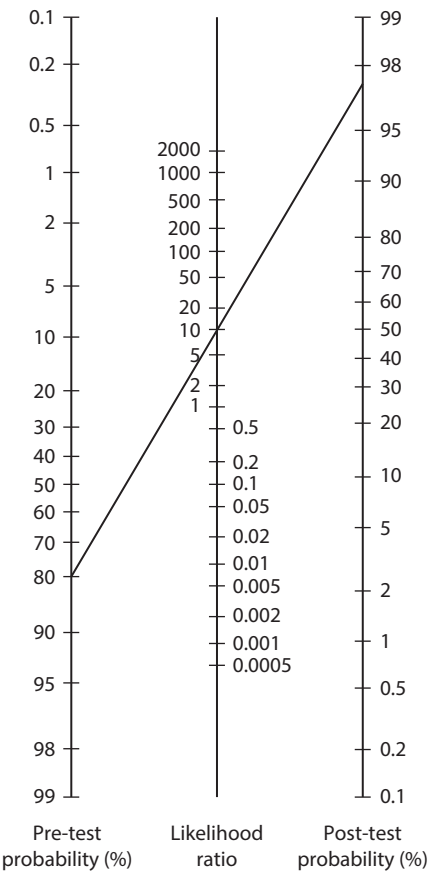


FIGURE 4.2 The nomogram for converting pre-test probabilities for post-test probabilities for pneumothorax. Pre-test probability = 0.8; likelihood ratio = 9. Post-test probability = 0.973 (97.3%), which is the same as the calculated post-test probability shown in ‘Spod’s Corner’.

SPOD'S CORNER: OTHER USES OF ODDS**Odds in meta-analysis**

Relative risk depends on the prevalence of the disease being studied but odds ratios do not. So if you are comparing multiple trials with differing prevalence in each trial, then combining relative risks is not mathematically sound but combining odds ratios is.

Odds in case-control studies

In case-control studies, the researchers define how many patients there are with and without the disease, i.e. they define their own prevalence. It therefore is not appropriate to use a calculation (risk ratio or relative risk) that is influenced by prevalence.

TAKE HOME MESSAGE

1. Odds and probability are two ways of expressing the chance of something happening but they are not interchangeable. Probabilities are easier to understand than odds and in appropriate situations should be used in preference to odds.
2. Bayesian theory is the process of using the likelihood ratio of a test to move from pre-test probability to post-test probability.

Section 2

Critically appraising papers

5 How to dissect a paper

Papers are not single units; they consist of multiple parts. In this chapter, we look at what these parts are. This chapter refers predominantly to interventional trials, whether diagnostic or therapeutic, as this is what is needed for FCEM. We also explore the different facets of study design, and, over the next few chapters, how it impacts on critically appraising a paper.

The purpose of research is to answer a clinical question, which we all meet in our day-to-day clinical practice. For example in a patient with a wrist fracture, we might ask ourselves a number of questions: which is better to make the diagnosis, x-ray or ultrasound (diagnostic study) or whether haematoma block or a Bier's block is a better anaesthetic (therapeutic study). To answer this we would need to perform a literature review and if we find there is insufficient evidence to choose one over the other, it may prompt us to do a study.

In a study, we enrol a sufficient quantity of patients and once we have the raw data we use statistics to present the evidence and decide whether one is better than the other. We can then answer our question in light of our results plus any previous research. We're now in a position to write our paper and submit it for publication.

Prior to 1993, the way in which studies were written for publications varied tremendously. The CONSolidated Standards Of Reporting Trials (CONSORT) statement has meant that randomized controlled trials (RCTs) should now be reported in a specific way. A full description of the CONSORT statement is available online at <http://www.consort-statement.org>. This systematic way of presenting studies is aimed at therapeutic studies; however it is equally valid for diagnostic papers.

Putting it into the simplest terms, papers can be divided into the following sections.

- Title: a quick to read explanation of what is being studied
- Abstract: a short summary of the whole paper
- Introduction: states the clinical question and a summary of the findings of the authors' literature review
- Methods: explains how the study has been designed
- Results: shows the number of patients enrolled, those who dropped out and what the researchers' findings were
- Discussion and conclusions: explains the results in the setting of previous research and answers the clinical question

Now, let's discuss these sections in more detail.

TITLE

The first thing we notice on picking up any paper is the title. The title is often used to sieve papers when performing a literature review and as such is important. The title should be clear and give us an idea of what clinical question is being studied. It may also give an indication of the study design.

ABSTRACT

An abstract is a summary of the whole paper, usually set out in the same way as the main text. It allows the reader to get a feel for what was done, what results were found and what conclusions were drawn by the authors but will not allow the reader sufficient detail to properly appraise the paper. It does however allow a second level of sieving of papers when performing a literature review. In the FCEM critical appraisal examination you will not be provided with the abstract of the study.

INTRODUCTION

The introduction should state the research question the authors have designed the paper to answer. It should describe the importance of the clinical question, what is already known about the clinical question and why further study is necessary. *The research question can usually be found in the last paragraph of the introduction – this is an important shortcut in the context of the FCEM exam.* Once we know this information, we can decide if we think that the clinical question is important and relevant to our practice.

At this point in reading the paper, it's valuable to think about how *you* would have designed a study to answer the question the authors have set themselves. As you read through the methods compare the design that the authors have chosen to the one you would have chosen.

METHODS

The methods section describes the design of the study and how it was subsequently run. There may be an ideal study design, but this may not be practical for many reasons such as finances, ethics or availability of patients. The study design takes these practicalities into account but in reality, trials do not always run to the protocol they were designed to follow.

The elements that should appear in the methods section of a trial can be remembered by the acronym 'PICOT':

- P – Population
- I – Intervention
- C – Comparison
- O – Outcome Measures
- T – Time

POPULATION

The first thing to check in assessing the methodology of a study is the population being studied.

- How did the authors pick their study population?
- Where was it set: in the ED, in an outpatient clinic, or somewhere else?
- Is the study population representative of the population that we are interested in, i.e. emergency medicine patients?

Now consider who was included and who was excluded:

- What were the inclusion criteria?
- Did they include all the patients we would have included?
- What were the exclusion criteria?
- Were any groups of patients excluded that should not have been?

Having too strict inclusion and exclusion criteria may lead to an experimental group that is no longer representative of the population as a whole. This is known as diagnostic purity bias.

Next, think about the sample size in the study:

- How did the authors come to this sample size?
- Did they calculate it using a sample size calculation and if so did they quote all the factors required for this calculation?

In many diagnostic studies the sample size calculation is not mentioned in the paper; however it should be.

Finally we consider how we draw our sample from the general population.

- What type of sampling did the authors use?

When taking a sample, we want it to represent the population in which we are interested (Figure 5.1). There are many different methods of taking a sample.

- *Random sampling* is where patients are randomly taken from the population in which you are interested.
- *Systematic sampling* is where every n th patient is chosen (n th can be any number, e.g. every third or every fourth, etc.).
- *Stratified sampling* involves taking a random sample from different strata of a population, e.g. a random sample from smokers and non-smokers. Generally the best method is for the numbers in each sample group to reflect the size of each of these groups in the general population. For example if 30% of the population were smokers and we were to sample 100 people

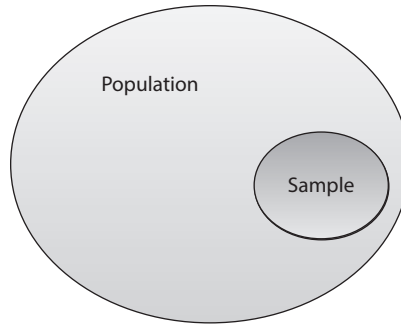


FIGURE 5.1 Sample size in respect to population.

then 30 people would be enrolled who are smokers. Stratified sampling may allow for a smaller sample size to be used and hence can be more cost-effective. It provides for a more representative sample and can allow better inferences to be made about specific subgroups.

- *Cluster sampling* is where instead of enrolling individual patients, we enrol groups of patients. Imagine we were looking at the impact of *bare below the elbows*. We could look at 10 hospitals where we compare bare below the elbow and 10 where we don't. We would have enrolled the hospitals, not the individual patients. Cluster sampling can be used when sampling of individuals is difficult or impossible. A regular example of this is in papers examining changes in service delivery. However, cluster sampling may not reflect the diversity of the population as a whole and sampling errors are more common. If you take a random sample of clusters it is a one-stage cluster sampling; if you also take a random sample within each cluster as well as a random sample of clusters, this is known as a two-stage cluster sampling.
- *Convenience sampling* is a sample taken when it is convenient for the researchers of the study, for example enrolling patients during office hours. This method of sampling is quick and easy but is prone to bias as the sample is likely not to be representative of the whole population.
- *Consecutive sampling* is similar to convenience sampling but with this form of sampling we try to enrol every patient who is available to us.

Fundamentally, we need to know whether the sample is representative of the population the authors wanted to study and of the population we are interested in.

INTERVENTION AND COMPARISON

In therapeutic studies, those enrolled undergo some form of intervention, which is compared with current practice or another treatment. In diagnostic studies we

compare a test with a gold standard. We discuss these types of studies in more detail in Chapters 6 and 7, but when reading a paper it is important to ask the following questions:

- What is the intervention in the study? What treatment or diagnostic test is being looked at?
- What is the comparison? In a therapeutic trial, what is the control and in a diagnostic study, what is the gold standard?
- What are the baseline characteristics of the patients and in therapeutic trials are they the same between these two groups?
 - Generally, 'Table 1' in the paper gives a list of patient baseline characteristics and it is always worth looking at this to make sure that there are no obvious differences between the groups.
- Have the authors identified any confounding factors? If so, how have they been taken into account?
- If none have been identified, are there any confounders that should have been taken into account?

There are specific questions to ask for therapeutic and diagnostic studies that are covered later. Do not worry if you are unsure of some of the concepts mentioned here, as we cover them in the next chapter.

OUTCOME MEASURES

The type of outcome we look at depends on whether it is a diagnostic or a therapeutic trial. In therapeutic papers we need to ask:

- What are the outcome measures?
- What is the primary outcome measure?
- Is this appropriate?
 - Outcome measures should be reliable (they should give the same result each time a measurement is performed if there has been no change), valid (how well they measure what we are expecting it to measure) and responsive (able to detect a change).
- Is the outcome measure clinically relevant, such as mortality or is a surrogate end point used?
 - A surrogate end point is an outcome measure that is used instead of a clinically relevant one, usually because it allows the trial to be smaller and to take less time. It also allows studies to be performed that would otherwise not be possible due to cost or an inability to enrol enough patients, e.g. using gastric pH instead of rate of GI bleeding when looking at proton pump inhibitors.

In diagnostic studies the outcome measure is a measure of the ability of the test under investigation to diagnose the disease.

TIME

It is important to assess if adequate time has been allotted to enrol and follow up patients.

- How long did it take to enrol the required number of subjects?
 - Taking a long time to enrol might suggest that the inclusion and exclusion criteria were too restrictive.
- Were patients followed up and if so, for how long and was this length of time appropriate?
 - Ideally patients should be followed up for an adequate period to determine whether they return to full function or have any complications of the treatment.
- Were there many dropouts and if so why did they drop out?
 - It is not unusual for patients to drop out of a study or be lost to follow-up. However, if this number is large it may be because of treatment side effects or a difficult to tolerate regimen, which may limit the treatment's usefulness. The authors should explain why patients dropped out so we can assess this.

RESULTS

The results section of the paper tells the findings of the study. It should tell how many patients were enrolled and how many patients dropped out or were lost to follow-up. We must consider:

- Did the authors enrol the required sample size?
- How many patients were excluded?
- How many patients dropped out or were lost to follow-up?
- Were all the patients compliant with the study protocol?
- Is the number of patients who dropped out of the study similar to the number estimated by the authors when they performed their sample size calculation?
- What was the number of protocol violations?
- Did the authors take into account dropouts and protocol violations in their analysis?
- Is there a CONSORT flow diagram or similar graph? (See Chapter 6.)
- Have the authors reported the appropriate statistics?
 - For therapeutic trials the authors should mention relative risk, absolute risk reduction, relative risk reduction and number needed to treat. For diagnostic studies they should quote sensitivity, specificity and ideally likelihood ratios and predictive values as well.

The authors should present their results in a way that can be easily understood by the reader. You should be able to draw a contingency (2×2) table from the data and ensure that the authors correctly quoted statistics. You might think that the quoted statistics are never wrong, but it can happen.

DISCUSSION AND CONCLUSIONS

The discussion should critically look at the study, outline how the results of the study fit into the current literature and explain how the authors have come to the conclusion they have. It should give us the limitations of the study and explain any ways in which the authors tried to take these into account. When reading the discussion and conclusion section consider:

- Are the results a surprise or expected?
- Why do the authors think this is?
- What limitations do the authors give?
- Are there any other limitations the authors haven't mentioned?
- Are their conclusions valid?
- Do they match your own conclusions?

At this point it is important to consider the validity of the study. There are two types of validity: *internal validity* and *external validity*. Internal validity is our assessment of the methodology of the study and whether this permits us to believe the results from this sample of patients. External validity is the extent to which the study can be extrapolated to the wider population beyond this sample of patients and in particular to our own local population.

After having decided about internal and external validity, the next thing to decide is whether the result of the study is important for your patients. If it is, then consider implementing it. This is not as easy as it sounds; there are many barriers which may prevent implementation of the intervention. We remember this as the 4 Bs:

- Badness: What harms are caused by the intervention?
- Burden of disease: Is it so rare that no one cares?
- Budget: Can we afford to make the changes?
- Beliefs: How will we educate people about this new intervention and how will we deal with intransigent colleagues?

OTHER FACTORS TO CONSIDER

Bias

Bias is a systematic error in the way observations were made, or the experiment is carried out, which leads to inaccuracy. Many types of bias have been described. Figure 5.2 and Table 5.1 give a list of the most common types of bias. Although we have listed a number of types of bias for information, it is more important to understand that bias leads to inaccuracy and how bias occurs, rather than knowing the name of each type of bias.

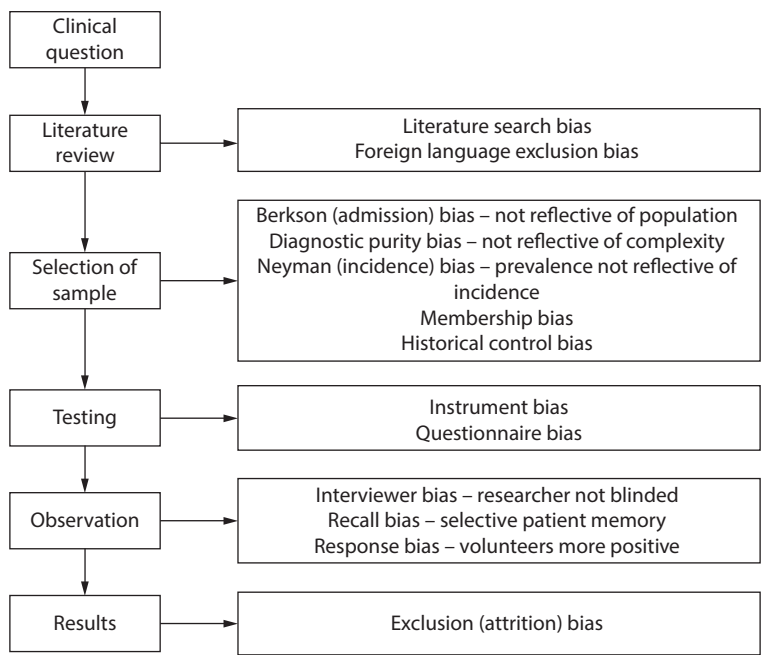


FIGURE 5.2 Types of bias from each stage of a study.

Bias can be combated by means of:

- A comprehensive literature review
- Consecutive patient enrolment
- Prospective study design – designing the study before you collect data
- Randomization
- Concealed allocation
- Blinding
- Use of intention to treat analysis

AUTHORS AND INSTITUTION

The list of authors of the paper is given a highly prominent and therefore important place in the writeup of any study. We need to be aware of the potential for our opinion about a paper to be influenced by the person or people who have written it. There are certain authors whose papers we are more likely to read and whose conclusions we are more likely to believe. This is particularly true of those authors who have a research pedigree. It is reasonable for us to consider that studies performed by these people are likely to be of good quality but we must never lose sight of the fact that each paper must be judged on its individual merits.

TABLE 5.1
Types of bias

Type of bias	Type of study	Explanation
Attrition	Therapeutic	Loss of subjects in a study leading to inequality between the two groups in respect to exposure and/or outcome.
Berkson	Case control	Patients who present to hospital may not be representative of community as a whole.
Diagnostic purity	Therapeutic	Exclusion criteria are too tight and so the study group may not be representative of the complexity of disease in the community.
Foreign language	Systematic review and meta-analysis	Less likely to include studies not published in English in a literature review.
Historical control	Observational	Apart from the intervention being studied, there may be a difference in management and prognosis between the research group and the controls simply because medical care improves (or should improve) over time.
Interviewer	Observational	The interviewer's opinions, prejudices and even non-verbal cues may influence the responses of subjects.
Language	Meta-analysis	English language papers are more likely to be published.
Literature search	Meta-analysis	Studies are more likely to be used if from a journal with a better reputation in a literature review.
Measurement	Diagnostic	Inaccurate measurement tool used.
Membership	Therapeutic	People who choose to be members of a group may differ in important respects to the general community.
Multiple publication	Meta-analysis	Variations of the same research or even the same study published in different journals.
Neyman	Observational	A bias seen in observational studies where we look at the links between a risk factor and the chance of getting a certain disease. If that disease causes the death of patients and our study only looks at survivors at any one point, we may underestimate the effect of the risk factor.
Non-respondent	Observational	Responders may differ significantly from non-responders.
Publication	Meta-analysis	Positive results are more likely to be published.
Questionnaire	Observational	Questions may be written to lead the respondent towards a particular reply.
Recall	Observational	Subjects with the disease are more likely to remember a potential causative agent than controls.
Response	Observational	Subjects may shape their response in order to please the interviewer.

It is also worth taking note of the institution or institutions in which the study was performed. Some institutions, like some authors, will be renowned for high quality research but remember that the environment in which the study is run impacts heavily on its external validity.

FINANCIAL INTERESTS

We have mentioned that studies cost money to run. Ask yourself:

- Where did the money come from for the study?
- Did the authors have any financial interest in the results of the study?

REFERENCES

The last section of the paper is the reference section. You might question the use of a list of studies. In fact the reference section is often a goldmine of information. When reading the paper look through the reference list to see what foundation the authors have built their study upon and to check whether you think they have missed any significant papers from the literature. Assess whether the number of references given, and the papers that the authors have drawn them from, represent a reasonable evidence base for the study performed.

PUBLISHING JOURNAL

A final point to note is in which journal the paper has been published. In the same way as considering who the authors are, we need to remember to not be too influenced by what journal a paper is published in, but to be aware that the larger journals have a bigger reviewing staff and in-house statisticians to check the internal validity of a trial before it is published. Because of this, most people feel instinctively more confident in the results of papers published in the bigger journals.

Many journals have an impact factor associated with them. The impact factor of a journal reflects how widely the papers it publishes are cited in the medical literature and as such can be seen as a measure of how important or 'big' a journal is.

TAKE HOME MESSAGE

1. A study consists of many parts and these can be divided into title, abstract, introduction, methods, results, discussion and conclusions.
2. The methods section can be divided into PICOT (population, intervention, comparison, outcome measures and time).
3. When reading a paper, consider the factors that are affecting internal validity and the factors that are affecting external validity.
4. There are many forms of bias – knowledge of each individual type is not necessary but understanding the concept is.

6 Therapeutic studies

Over the next two chapters we discuss therapeutic and diagnostic studies. We have given an overview of some of the important aspects of all trials in Chapter 5, and so if you haven't read Chapter 5 yet it would be worth reading before embarking on this chapter.

Therapeutic and diagnostic studies are by no means the only types of trial but they *are* the only types you are likely to get in the FCEM examination. We think this represents a sensible approach by the College, as a thorough understanding of both types will allow you to be proficient in day-to-day critical appraisal.

In much of clinical medicine we are interested in whether an intervention (a drug, a practical procedure, an operation, a change in process etc.) is better than what we already do. Of course what we already do may be nothing at all or may be an older, established treatment. Trials that compare the outcomes of two or more groups that are given different interventions are termed therapeutic trials.

In this chapter we work our way through various aspects of the methodology of therapeutic trials, initially with specific reference to randomized control trials (RCTs). We then discuss therapeutic studies of service delivery changes and the specific features that relate to them. We follow this with an explanation of how to understand and use the results that the authors of the trial produce. We have distilled much of this information into the checklist in Appendix A.

Let's start by looking at the different aspects of design and methodology of the trial that you should consider when reading an RCT.

DESIGN AND METHODOLOGY FEATURES OF RANDOMIZED CONTROL TRIALS

WHAT IS THE INTERVENTION BEING COMPARED AGAINST?

In most trials we're interested in how good an intervention is and this usually involves comparing it against another intervention, which we term the 'control' – hence the randomized *control* trial.

Our choice of control may vary depending on what we're looking at: if there is no standard treatment for us to use as a control, we may compare our intervention to no intervention at all. Approximately 1 in 3 people will report a difference in their symptoms when put on a treatment even if that treatment has no clinical effect, and this phenomenon is known as the *placebo effect*. Some would argue that this is why homeopathy purports to work. In order to remove this potential bias, in trials where the control group is going to receive 'no intervention', this should be disguised by using a placebo treatment for the control group.

In an ideal world the control group should run alongside the intervention group and they should receive what would be considered the ‘standard’ current treatment. When you read a paper and this is not the case, you need to factor this into your assessment of the result of the study.

In some disease processes the use of a contemporaneous control group is not possible. An example of this is in rarer illnesses where it may take many years to enrol sufficient people. In these cases it might be appropriate to use historical controls – i.e. compare your intervention against patients who have received the control intervention in the past. Although this is an understandable decision, we must be aware that medicine is (we hope!) an improving discipline and therefore the historical control group is usually more likely to have a poorer outcome simply due to the fact that medical care will have improved over time.

In disease processes where there is a guaranteed 100% mortality rate, it might be appropriate to not use a control group and give everyone you enrol the intervention under investigation. In this case the moral argument for this will not necessarily hold true if the intervention itself has significant side effects such as pain.

SELECTION OF TRIAL PARTICIPANTS

Selection of participants is covered in detail in Chapter 5. The important points to remember are:

- Inclusion and exclusion criteria should be listed within the trial and these determine which patients can be selected for the trial and which cannot.
- Too-tight inclusion and too-broad exclusion criteria run the risk of making the trial result hard to generalize to the rest of the population.
- Selection should be determined only by the stated inclusion or exclusion criteria; otherwise the result will be biased.
- Selection happens first, then requesting the patients’ consent to enter the trial, then (in an RCT) randomization to one arm of the trial.

ALLOCATION OF TRIAL PARTICIPANTS AND RANDOMIZATION

Once we have selected and gained consent from the patients we are going to enrol in the trial, we need to decide which group they will go into: intervention or control. If we are allowed to pick which patients go into which group, we will be able to influence the trial unfairly – perhaps consciously or unconsciously – by allowing certain types of patients access to one trial arm but not the other. In order to be confident that the result of the trial is a fair reflection of only the difference between the two interventions, we need to be certain that patients who are enrolled have an equal chance of being put in either of the groups. In RCTs we do this by means of two processes: concealed allocation and randomization.

Concealed allocation is the process whereby patients and researchers are unaware of what group trial subjects will end up in once they have been enrolled in the trial.

When patients agree to enter a trial they need to be fully informed about what might happen to them. When there is more than one potential arm to the trial what

we need to avoid, however, is letting the patient know which of the treatments they are going to get. This is because the decision on whether to consent to taking part may be influenced by whether the patient feels they are likely to get a ‘new and exciting’ (or indeed ‘experimental and scary’) drug or an ‘old and outdated’ (or ‘tried and tested’) one.

This potential for bias may also influence the researchers enrolling patients into the trial and so they too must be prevented from being aware of which group patients will be put into after enrolment into the study. This prevents the researcher from either encouraging or dissuading an individual into a particular trial arm. A good example where allocation concealment fails is when allocating people by days of the week (e.g. those that attend on a Monday go into group A, those that attend on Tuesday go into group B etc.). Although it gives the impression of random selection of patients it is totally predictable and therefore does not achieve allocation concealment. Researchers enrolling patients will know which group patients will go into and this may well affect whether they enrol certain patients.

Randomization is important as it distributes known and unknown confounders (see later) between the experimental and control groups as long as there are sufficient numbers in the study. We must not just accept that a study is randomized but should look at how randomization was performed. There are various types of randomization, some of which are described here:

- *Simple randomization* allows patients to be randomly divided between both groups. This can be performed with a coin toss, random number tables or ideally with a computer program.
- *Block randomization* occurs when patients are randomized in blocks. Keeping things simple, imagine you want to divide patients equally between two interventions; for each block of six patients, three would go to one arm and three would go to the other arm but the order of the patients within the block would be random. The main aim of doing this is that it will result in more balanced numbers in each group at each time period of the trial and allow for easier interim analysis. The smaller the blocks the higher the risk that the allocation process may be predicted and hence allow for the possibility of bias. This system can be used for any number of arms within a trial.
- *Stratified block randomization* is block randomization with stratification to take account of a confounder. For example, there is a new antibiotic to treat pneumonia but we are worried that smoking is a confounder. We group people as smokers and non-smokers and then block randomize them within the strata of smokers and non-smokers (Figure 6.1). This will ensure an appropriate spread of smokers and non-smokers in the treatment and control groups. In trials with smaller numbers this may not occur with simple randomization.
- *Quasi-randomization* is when there appears to be randomization but it is not truly random, e.g. randomization by day of the week.
- *Cluster randomization* is when the unit of randomization is not an individual study participant but rather a group of individuals treated in the same place or by the same person (Figure 6.2). Cluster randomization allows us to study interventions that are applied to groups rather than individuals,

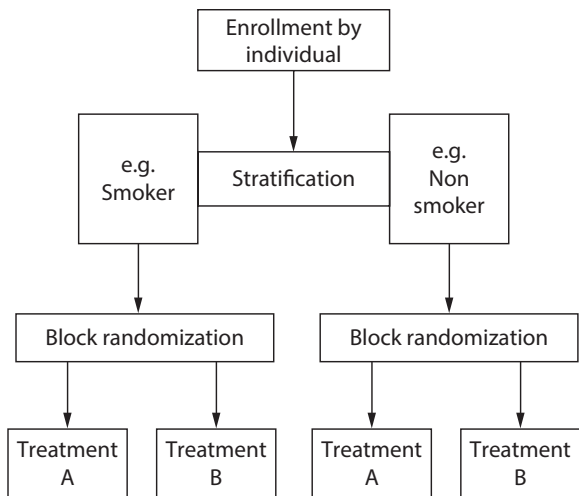


FIGURE 6.1 Stratified block randomization.

such as process changes like the introduction of a Clinical Decision Unit (CDU). If we were to try having two different processes operating in the same clinical area (e.g. some patients having access to a CDU and others not), there is the potential for confusion, protocol violation and bias; cluster randomization works to avoid this. Unfortunately, cluster randomization tends to require a greater complexity of design and analysis. How this is done is beyond FCEM but it is important that the authors tell you they have taken this into account.

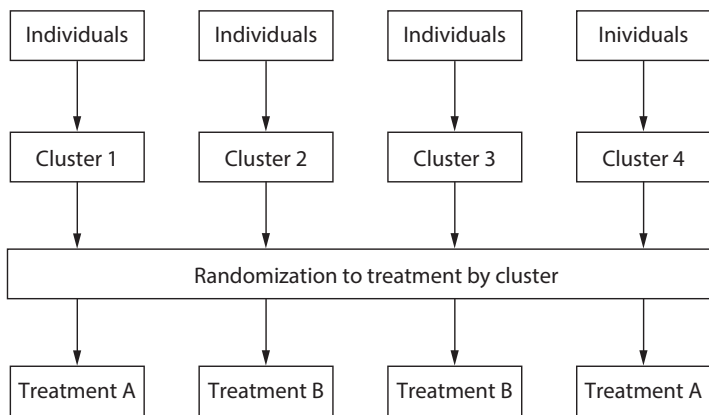


FIGURE 6.2 Cluster randomization. The number of clusters will vary depending on the study.

- *Adaptive randomization* is a mechanism by which the probability of being randomized to a group varies dependent on previous patient allocations.

When discussing randomization there are several important points to be aware of:

- The more that outside influence (researchers, clinicians, patients etc.) can alter trial allocation, the greater the chance of biasing the result of the trial.
- Randomization and trial allocation happen once the participant has been enrolled irreversibly in the trial – although the patient can decline to be involved at any point, once enrolled the researchers should count them as having been part of the trial.
- There are various forms of randomization process, each with its pros and cons but as a general rule, the further the researchers are kept from the randomization process the smaller the likelihood that it will be able to be subverted (intentionally or otherwise).

BLINDING

How we act is influenced by our knowledge and beliefs and this can consciously and subconsciously introduce bias into a trial. To prevent this from happening a process known as blinding is used. When we use blinding we prevent the people involved in the trial (e.g. researchers, patients, clinicians or statisticians) from knowing whether the patient is in the control or intervention arm of the trial. The more people who are blinded in a trial the lower the chance of bias. There are many forms of blinding depending on who is unaware of the patient designation.

- No blinding: In some trials everybody knows who received the intervention and the control; this is referred to as an open trial.
- Single blind: Patient or researcher are blinded to whether the patient has received intervention or control.
- Double blind: Both patient and researcher are blinded.
- Triple blind: Patient, researcher and analyst are blinded.

As part of blinding, it is often necessary to camouflage the treatment the patient is receiving. Unless the intervention looks exactly the same as the control, both the patient and clinician may be able to figure out which is which and this may introduce bias. To camouflage treatment we may have to make the control and treatment look similar (e.g. same colour and shape tablet, same packaging, etc.). If this is not possible because the treatments are too different (e.g. an oral and an intravenous medication), a double dummy trial where the patient appears to receive both treatments can be used.

If blinding is unsuccessful, the following biases can occur:

- Expectation bias: Patient, carers, clinicians or researchers expect a certain outcome from a particular treatment (or in the case of the patient, based on the fact they are getting a ‘new’ treatment) and this influences how they experience or report their symptoms.
- Attention bias: Patients recognize that they are receiving a new and therefore ‘special’ treatment and hence are receiving greater attention than normal, which makes them more likely to interpret their situation more positively.

The difference between allocation concealment and blinding is an important one and it is important that you understand that allocation concealment occurs in order to prevent people from knowing what trial group participants *will* go into whereas blinding prevents people knowing which group they are in *once they are in the trial* (Figure 6.3).

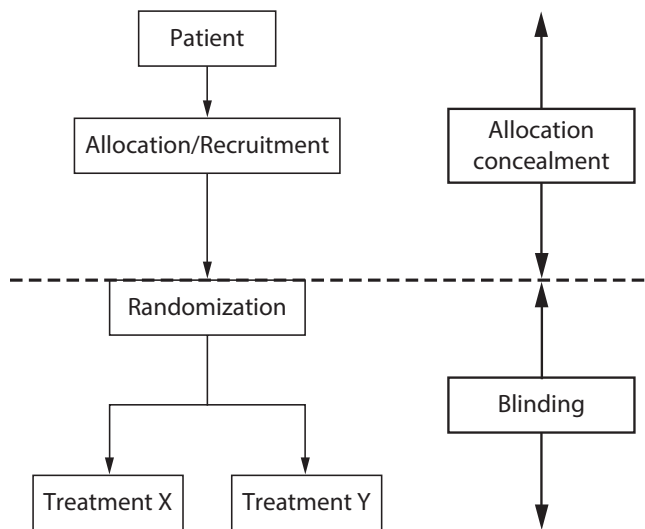


FIGURE 6.3 Many FCEM candidates struggle with the difference between allocation concealment and blinding. Allocation concealment is used to prevent bias up to the point of randomization whereas blinding is used after randomization.

Although in the ideal world all people involved in a trial should be blind to the intervention each patient has received, this is not always possible. For example, in the cooling post-cardiac arrest papers, given that those in the intervention group needed to have their temperature maintained at a particular level, those providing the intervention needed to know who was in each group. This does not, however, mean that no blinding can happen at all. For example, the people interpreting the results from the patients, including analysing clinical outcome, can still be blinded to what treatment group each patient was in.

As a final point, it’s worth mentioning that in some cases a lack of blinding may be more likely to cause problems than others. When measuring outcomes such as ‘died or didn’t die’ the interpretation is less prone (although not immune!) to bias than when recording outcomes such as patient satisfaction, where a degree of interpretation of the answer from the patient may be possible.

WERE ALL THE PATIENTS ACCOUNTED FOR?

Once enrolled in the study (which includes gaining the patient’s consent) the patient becomes irreversibly part of the study. Of course, the patient may or may not complete the study as planned. This may be for a variety of reasons including protocol violation (for instance due to physician choice or intolerability of the intervention) or being lost to follow-up. This is to be expected in a proportion of the trial participants

and should be planned for when the sample size is being calculated. However, a good RCT should record and publish how the patients flowed through the trial and what happened to the patients who for whatever reason either failed to get enrolled in the first place or didn't reach the end of the trial.

International guidance, the Consolidated Standards of Reporting Trials (CONSORT) statement 2010, tells researchers the recommended format for reporting their trial. As part of this, RCTs should produce a flow chart showing what happened to all the patients in the study and how many had their results analyzed. It is usually the first diagram in a study and is labelled the CONSORT flow diagram, an example of which is given in Figure 6.4.

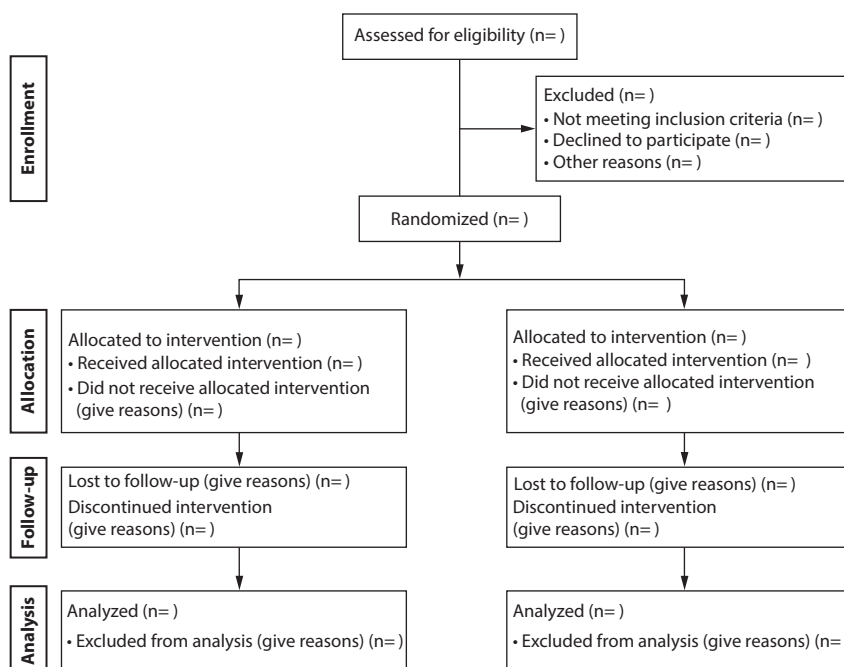


FIGURE 6.4 The CONSORT flow diagram.

Once you've looked at the CONSORT flow diagram consider whether the amount of patients successfully followed up seems appropriate. It's reasonable to expect some patients to be lost to follow-up but the acceptable rate will vary depending on the outcome that is being measured. For example, we would expect a low rate if the study is looking at in-hospital mortality as this data should be relatively easy to track down, whereas an outcome measured by postal follow-up is likely to result in far greater loss to follow-up.

WHAT SORT OF OUTCOME DID THEY LOOK AT?

The aim of most trials is to look at an intervention and see what it does. As clinicians we're most interested in whether an intervention makes people better, but how do we measure that? If we think about a trial comparing two antihypertensives, what is

‘better’? Perhaps it’s a lower rate of cardiac disease or of stroke, perhaps it’s the one with the least side effects or perhaps it’s something else. We divide outcomes into two main groups: clinical end points and surrogate end points.

Clinical end points

Clinical end points are features that are important to the patient and as such are the outcomes that measure how the patients feel or cope or how likely they are to survive or develop a certain disease.

Outcomes such as death or the rate at which people develop a disease are conceptually easy to understand and their measurement, if done properly, should be reliable. They are usually termed ‘hard clinical end points’. It may take a long time for these outcomes to happen and therefore it may cost a lot to run the trial, meaning that not infrequently they are not reported by trials.

Measures of how the patient feels or assessments of their quality of life are often referred to as ‘soft clinical end points’ and because they may be prone to interpreter and reporter bias, many people are instinctively wary of them. When measuring outcomes that may be seen to be open to these sorts of bias, it’s important that the researchers use a scoring system that has been properly validated which has the effect of making them slightly less ‘soft’.

The *primary end point* is the main parameter of interest. It is the outcome that is used to calculate the sample size and for which a study is powered.

Secondary end points are parameters that are measured as part of the study but that the study is not designed to look at specifically. As such they may be interesting, but the results of a secondary outcome should not be used as definitive evidence since the study is not powered to look at them.

Composite end points are combinations of outcome measures, e.g. all cardiovascular deaths rather than strokes or myocardial infarctions individually.

Surrogate end points

Surrogate end points are usually factors such as laboratory test results that can often be obtained relatively quickly and easily. With these sorts of outcomes, the researchers are suggesting that a change in the surrogate end point (e.g. blood pressure) equates to a clinical end point (e.g. a reduction in the risk of stroke). Surrogate end points, given that they are by definition not the actual things that patients are interested in, should be viewed with a degree of caution. Certainly if we are to accept that a reduction in a surrogate outcome is important, we need to be convinced of the clinical argument that changes in the surrogate lead consistently to changes in the clinical outcomes that we are interested in.

LOOK FOR CONFOUNDERS

Confounding is different to bias, which we discussed in Chapter 5. It is an error in interpretation of the meaning of the study result and as such the results may be accurate and precise but false conclusions are drawn due to assuming the different groups are the same in every way except the intervention. Confounders are

extraneous factors which can affect the outcome and will lead to invalid inferences if not distributed equally or taken into account. They may be positive (inferring an association when not really present) or negative (hiding an association which is present). Common confounders in most medical trials are age, gender, smoking and socioeconomic group.

Before we look at the result we need to ask ourselves whether there are any differences between the two groups that may have influenced the result. Any random method of sampling may simply, due to chance, result in there being differences between the groups that are beyond the control of the researchers. Good quality randomization of sufficiently large numbers of patients will reduce this but there is always a chance that it will happen and if it does we need to take that into account when assessing the result of the trial. These differences could be ‘confounding factors’.

Interventional trials should provide a table showing the baseline characteristics of the different groups and this is usually labelled as Table 1. When reading a trial, we need to check this table to see if we pick up any differences that we think might be important. In some trials, the authors will help by providing results of tests of hypothesis (p-value) to try to tell us whether any difference between the groups is statistically significant.

Although the table we’ve described will report the baseline characteristics that the authors have collected, we cannot rely on the fact that they will have considered or collected every confounding factor. Imagine our trial comparing one antihypertensive against another looking at the rate of stroke at 6 months; we collect data about the age, sex and body mass index and present that in the paper that we publish. If, however, we hadn’t thought to collect data about whether patients smoke or not, it’s possible that one group might have significantly more smokers than the other and it’s obvious that this may well have a large impact on the outcome (stroke) that we’re interested in. As such smoking would be classed as an ‘unknown confounder’. Although the confounder is unknown in this case, it’s obvious that it should have occurred to us and we would list this as one of our critiques of the paper. With other less obvious confounders, we might not notice this however. The purpose of randomization is that it acts to spread the confounders that we know, the ones we can guess about, and the ones we don’t know, evenly between the groups.

Confounders can be controlled by means of:

- Good study design.
- Randomization if there are sufficient patients in the study.
- If confounders are known then they can be adjusted for in an analysis.

If the confounder is known prior to study design, then the study design can:

- Restrict it (for example avoiding respiratory arrests that are more likely to be PEA [pulseless electrical activity] or asystolic in cooling studies).
- Allocate by minimization. This is a process whereby instead of patients being allocated to treatment groups purely by a random process, we control the allocation to ensure equal numbers of confounders are put in each of the treatment arms. This can be pure minimization where there is no

randomization or by randomization that is weighted by a confounding factor – the mechanics of this process do not need to be known for the FCEM examination.

- Match confounders so that there are equal numbers of confounders in each group – stratified randomization.

Confounders can be controlled during the analysis:

- Standardisation: risk in the exposed group is calculated as if the confounder was the same, e.g. adjusted hospital mortality by age.
- Stratified by separating the data into groups (or strata) based on the confounder.
- Multivariate analysis, which takes into consideration more than one confounder.

HOW DID THEY ANALYZE THEIR RESULTS?

Below we explain the mathematics that you need to know for the FCEM examination. Before we get to that we need to look at the two major ways in which the results might be analyzed.

Consider having a trial with two groups of patients: one randomized to the intervention, one to the control. Once patients have been enrolled into the trial, we've talked about the fact that they must be accounted for in the CONSORT flow diagram. We've also mentioned that not all the patients allocated to a group will complete the treatment that they were allocated to.

There are two main ways of dealing with the fact that all the patients entered into a trial may not make it through to the end in the way that was planned: intention-to-treat analysis and per-protocol analysis.

Intention-to-treat analysis

In an intention-to-treat analysis the patients are analyzed by the group they were originally enrolled to regardless of what treatment they actually received. For example, in our antihypertensive trial, it may be that a number of people who were assigned to the new antihypertensive chose to stop taking it partway through the trial due to feeling nauseated. Although they didn't manage to complete the course of antihypertensive we still use their outcomes as part of the intervention group. This has the effect of lessening any difference that might occur between the two groups. Analysing by intention to treat means that it is *less* likely we will show a difference if one exists (i.e. a greater chance of a type 2 error). This sort of analysis takes account of what happens in real life – that patients do not always comply with the treatments they are prescribed by clinicians – and as such is considered the best way to calculate the results from a trial.

Per-protocol analysis

Per-protocol analysis is considered the opposite of an intention-to-treat analysis. Patients are only included in the results if they managed to complete the treatment that they were allocated to. If they drop out or for any other reason do not adhere

to the trial protocol their outcomes do not get included in the statistical analysis of the results. By doing this it might be argued that we get a ‘purer’ look at the actual effect of the treatment protocol versus the control intervention – after all, failure to have received the full protocolled treatment in the intervention group may not just be due to side effects or intolerability of the treatment but also because of errors of the clinicians running the trial. This form of analysis makes it more likely that the trial will demonstrate a difference if one exists (i.e. a lower chance of a type 2 error) but it may hide potential problems with the treatment that make it less likely to be tolerable by patients.

SERVICE DELIVERY THERAPEUTIC TRIALS

Most therapeutic trials that you will come across (both in the FCEM examination and in real life) will be variations on the theme of RCT, with the variation being in the amount of blinding and the degree and type of randomization. They will compare one or more interventions – often a new drug – against a control and look at defined outcomes.

Before we move on to talk about diagnostic studies in Chapter 7, it’s worth mentioning one other type of therapeutic trial and some of the features that commonly come up regarding them – service delivery therapeutic trials.

It is very much appropriate for us to want to look at how we might improve patient care by changing *the way we deliver* a service rather than simply by trying to deliver a new treatment. A relatively recent intervention is the use of an early senior assessment team in the ED with the aim of improving the speed and quality of decision making that patients receive. It is right to wonder whether the efficacy of this is something we could test in a clinical trial. Any sort of service delivery change (e.g. triage processes, educational interventions, staffing templates, facilities such as short-stay wards, etc.) like this can be assessed in a therapeutic trial that is very similar to the RCT we previously described.

BENEFITS OF CLUSTER RANDOMIZATION

We randomize some patients to receive early senior assessment (let’s call it Patient Assessment Team or PAT) and a control group to receive standard care. Our next question is how we go about this randomization. The introduction of PAT is not something that lends itself easily to randomizing one patient to PAT and the next to standard care within a single ED. However, cluster randomization, discussed earlier, would be ideal for this. We could enrol a number of hospitals to our trial and randomize each hospital rather than each patient to either PAT or standard care. Service delivery papers lend themselves very well to cluster randomization.

NON-RANDOMIZED SAMPLING

Cluster randomization does, however, have drawbacks: allocation concealment is not usually possible resulting in the potential that patients and carers will choose which

hospital to visit depending on the service they are offering. Also, clustering adds challenges to the statistical analysis of the data, and the statistical power of the trial may be substantially reduced. Because of this, service delivery papers may instead use non-randomized methods of sampling, mostly by comparing the intervention group, when the service delivery change has been made, with a historical group when it hasn't. There are two main issues with this: first, as we hope that medical care improves over the course of time, the result may be biased to show an improvement in outcomes even if this is unrelated to the service change we are looking at. Second, if the historical control group were unaware that they were taking part in a trial but the intervention group are, the Hawthorne effect (see below) may apply unfairly to the two groups.

HAWTHORNE EFFECT

Discussion of service delivery papers would not be complete without mentioning the Hawthorne effect. In the 1920s and 1930s a series of experiments were performed in the Western Electricity Company in Chicago, the most often cited of which was an experiment where the researchers varied the intensity of light in the working environment. The results showed that productivity increased both when light intensity was increased *and* reduced. The generally accepted interpretation of this is that the act of observing a situation (including a clinical environment) improves the performance (and hence medical outcomes) of those *in all* groups within the trial. One conclusion of this is that the results obtained in research studies may be expected to be better than those that would be obtained in real life.

GENERALIZABILITY AND SUSTAINABILITY

When assessing service delivery change studies, perhaps even more so than RCTs on pharmacological interventions, we need to be mindful of whether the change that is being suggested is something that is both generalizable to our clinical situation and something that we would be able to sustain. Many service delivery changes are possible within the confines of a funded and well-staffed trial but may not be so successful when exposed to the realities of the staffing template and financial environment of everyday clinical practice.

HOW TO INTERPRET THE RESULTS SECTION OF A THERAPEUTIC TRIAL

The key to answering any questions in the therapeutic results section is drawing the 2×2 table (often called a contingency table). Based on the 2×2 table, we can then make a series of calculations based on this table. But remember the exam is about understanding what you are calculating and not just rote learning.

Prior to doing this section make sure you have an understanding of probability (discussed in Chapter 4) and remember that in clinical research, risk equals probability.

Risk = Number of times an event occurs/Total possible events

To draw the 2×2 table, you need to know the following:

- The number of patients in the control group who got the disease
- The number of patients in the control group who didn't get the disease
- The number of patients in the experimental group who got the disease
- The number of patients in the experimental group who didn't get the disease

These numbers go into the 2×2 table, which is traditionally written in the following way: The disease outcome is on the horizontal axis with 'positive' outcome (meaning they got the disease) on the left-hand column and 'negative' outcome (meaning they didn't get the condition) on the right-hand column. The treatment or exposure to risk is on the vertical axis with being treated written above the control. An example is shown here:

		Disease Status/Outcome	
		+	–
Treatment/Exposure	+	A	B
	–	C	D
	Experimental		
	Control		

In whatever format the table is drawn in the journal article and the exam, we would suggest that you redraw it in the above way, as otherwise there is potential to confuse yourself, especially when stressed.

Whenever you are asked to perform a calculation from this table, always start from the definition and show your workings. It shows you understand what you are doing and will prevent you from making an error.

CALCULATIONS BASED ON A 2×2 THERAPEUTIC TRIAL TABLE

Risk in the control group = Control event rate (CER)

$$= \text{Number who had a 'positive' outcome while having the control treatment} \div \text{the total number in the control group}$$

$$= c/(c + d)$$

Risk in the experimental group = Experimental event rate (EER)

$$= \text{Number who had a 'positive' outcome while having the experimental treatment} \div \text{the total number in the experimental group}$$

$$= a/(a + b)$$

From these two risks, all the other important results can be calculated.

Absolute risk reduction (ARR)

- = Simple difference in risk (outcome rate) between the two groups
- = Control event rate (CER) – Experimental event rate (EER)

Advantages of ARR

The ARR takes baseline event rate into account – we explain this more fully in the examples later. ARR is easy to understand. It is very straightforward to comprehend the idea of the difference in percentages between two groups. For those of us interested in critical appraisal (and those trying to get through FCEM) it allows us to calculate ‘number needed to treat’ (NNT).

A note about baseline rate: Before we start to interpret the result of the trial, one of the first things we need to assess is how important this disease process is. Look at the rate of disease in the control group; if it’s a very serious disease resulting in death then even low rates will be of interest (and hence small drop important). If it is a less serious disease, low rates and small absolute differences between groups may not be so important.

Number needed to treat/harm (NNT/NNH)

- = Number of patients needed to be treated to achieve one additional good (or bad) outcome
- = $1/ARR$

Advantages of NNT

Like the ARR, the NNT takes baseline event rate into account. In our opinion, it is the most useful of all the results that can be derived from a therapeutic trial to help doctors (and other health care professionals) to decide if we should use a treatment. Given that, if the paper does not quote the absolute risk reduction and the number needed to treat (or harm) – we have to ask why.

The other calculations are less easy to comprehend and more open to make treatments appear more important than they really are. Although they compare the two groups, they are not influenced by baseline rate. Drug reps often use them to promote their treatment and so we must really understand what they are telling us.

Relative Risk (RR)

- = Ratio of the risk (or outcome rate) in the experimental group ÷ the control group
- = EER/CER

Relative Risk Reduction (RRR)

- = Difference in risk (or outcome rate) between treatment and control group divided by risk (or outcome rate) in control group
- = $(CER - EER)/CER$
- = ARR/CER

We must be very cautious with the relative risk reduction. Even if there is a small baseline figure for getting a bad outcome, you can show your drug is effective. See the following example.

WORKED EXAMPLE 6.1

THE NEW STENT THAT BEATS ALL OTHERS – SUPER STENT

The cardiology team have been seduced by the rep of Super Stent. The leaflet says ‘with a relative risk reduction of over 75% and costing only £10,000 more, how can you afford not to use Super Stent?’ The small print says ‘of 1000 people in the group given the standard stent, 10 people died while in the Super Stent group of 1000 people, 2 died’. The chief executive says ‘ok’ and says ‘it can be paid for by the ED surplus’! How much will it cost to save a life? Go through the following steps.

Questions

1. Draw a 2×2 table.
2. Calculate risk in the control group.
3. Calculate risk in the experimental group.
4. Calculate ARR (absolute risk reduction).
5. Calculate RR (relative risk).
6. Calculate RRR (relative risk reduction).
7. Calculate NNT (number needed to treat).
8. Calculate the increased cost per extra life saved.

WORKED EXAMPLE 6.1 ANSWERS

1. Draw a 2×2 table.

	Positive outcome, Death	Negative outcome, Alive
Super stent (experimental group)	2 (a)	998 (b)
Normal stent (control group)	10 (c)	990 (d)

2. Calculate risk in the control group.
 - = Control event rate (CER)
 - = the number who had died while having the control ÷ the total number in the control group
 - = $c/(c+d)$
 - = $10/1000 = 0.01$ (1%)
3. Calculate risk in the experimental group.
 - = the number who had died while having the experimental treatment ÷ the total number in the experimental group
 - = $a/a+b$
 - = $2/1000$
 - = 0.002 (0.2%)

4. Calculate the ARR.
 - = the simple difference in risk (outcome rate) between the two groups
 - = Control event rate (CER) – experimental event rate (EER)
 - = $0.01 - 0.002$
 - = 0.008 (0.8%)
5. Calculate the RR.
 - = a ratio of the risk in the experimental group ÷ the risk in the control group
 - = EER/CER
 - = $0.002/0.01$
 - = 0.2
6. Calculate the RRR.
 - = the difference in outcome rate between treatment and control group ÷ outcome rate in control group
 - = $(CER - EER)/CER$
 - = ARR/CER
 - = $(0.01 - 0.002)/0.01$
 - = 0.8 (80%)
7. Calculate the NNT.
 - = $1/ARR$
 - = $1/0.008$
 - = 125
8. Calculate the increased cost per extra life saved.
 - = $NNT \times \text{increased cost per treatment}$
 - = $125 \times 10,000$
 - = 1.25 million
 - Tell the cardiologists ‘no way!’ (especially not in the current economic climate).

Odds ratios

For many people odds are difficult to understand. As such, in therapeutic trials (though not meta-analysis) there is no reason for the results sections to contain odds of getting conditions in the control and experimental groups. There is no reason why ratios or statistical tests should be based on odds, either. The authors may quote odds ratios, but beware they may exaggerate the effectiveness of a treatment. This is shown in detail in Chapter 4.

If there are odds ratios quoted, this is how they are calculated. Remember to always refer back to the 2×2 table.

		Disease status/Outcome	
		+	–
Treatment/Exposure	+	A	B
	–	C	D
	Experimental		
	Control		

(Remember that a ‘positive outcome’ in a 2×2 table is actually a bad outcome, e.g. death.)
The odds in the treatment group of a trial

= Number of patients who received the new intervention and had a positive outcome \div Number who received the new intervention and had a negative outcome.
= a/b

The odds in the control group of a trial

= Number of patients who received control treatment (standard treatment) and had a positive outcome \div Number who received control treatment (standard treatment) and had a negative outcome
= c/d

$$\text{Odds ratio} = \frac{\text{Odds in treated group}}{\text{Odds in control group}}$$

WORKED EXAMPLE 6.2

ODDS FOR SUPER STENT TRIAL

Questions

Using the 2×2 table, work out the following:

1. Odds in the treatment group
2. Odds in the control group
3. Odds ratio

Answers

Before you start, redraw your 2×2 table.

	Positive outcome, Death	Negative outcome, Alive
Super stent (experimental group)	2 (a)	998 (b)
Normal stent (control group)	10 (c)	990 (d)

1. Odds in the treatment group.
= the number of patients who received Super Stent and died \div the number who received Super Stent and survived
= $2:998$
= $1:499$

2. Odds in the control group.
= the number of patients who received Super Stent and died \div the number who received Super Stent and survived
= 10:990
= 1:99
3. The odds ratio.
= odds in the treatment group \div odds in the control group
= 1:499/1:99
= 0.198

(Note: As the outcome (death) is rare in this study, the odds ratio is similar to the relative risk, but with higher prevalence this may not always be true.)

TAKE HOME MESSAGE

1. Therapeutic trials compare two or more interventions.
2. Randomization is the process whereby patients have an equal chance of going into either the control or the treatment group.
3. Randomization reduces the impact of known and unknown confounders.
4. Blinding is the process whereby those involved in the trial (patients, researchers etc.) do not know which arm of the trial the patient is in. This reduces bias.
5. The most valuable result in a therapeutic paper is the number needed to treat. This is calculated as $1 \div$ absolute risk reduction.

7 Diagnostic studies

Studies that look at the ability of a test to predict the presence or absence of a disease are called diagnostic studies. In the FCEM critical appraisal examination you are likely to be given either a randomized controlled trial (RCT) or a diagnostic study, and your first job in the exam is to recognize and write which it is.

Diagnostic studies are a type of prospective experimental study. Often in this type of study, the researchers look to compare a new test against an established test to see how the new test performs. The test being investigated is termed the ‘index’ test and the test it is compared to is called the ‘reference’ or ‘gold’ standard. However, it’s worth highlighting that although a diagnostic study often uses one individual index test versus one reference test (e.g. ED ultrasound versus CT aortogram), the ‘tests’ used may be a combination of clinical features (e.g. hypotension, pain radiating to back, presence of pulsatile mass, over 65 years old) or an individual test with addition of clinical features (hypotension, pain radiating to back, presence of pulsatile mass, over 65 years old *plus* ED ultrasound). This concept is equally applicable to the gold standard as it is to the index test.

In this chapter we cover the aspects of the methodology necessary to assess the quality of these sorts of trials. We then move on to what should appear in the results section. We describe how the data generated by the trial is used to allow different groups of people with differing priorities to assess the value of the test being studied. The checklist in Appendix A can be considered a simplified version of this chapter.

DESIGN AND METHODOLOGY FEATURES

When reading a diagnostic paper, closely consider the following five areas of the trial design and methodology: patient selection, the reference (or gold) standard, blinding, biases, and qualities of the test.

PATIENT SELECTION

Who were the subjects in the study?

As with all studies, we need to look at the environment in which the study was done and what sort of patients it was done on. If the test being investigated has been done on only children in a general practice environment it is unlikely to be applicable to our adult ED patients.

How were the patients selected?

Next, we must ask how the researchers have selected the patients for their investigation. Is it likely that the selected sample is truly reflective of the general population

or is it possible that the selection technique has resulted in a biased sample? For example, if the researchers have only selected people who presented on Monday through Friday perhaps these patients have particular characteristics that make them unrepresentative of patients who present at weekends.

What was the prevalence of disease in the study sample?

If the prevalence of disease in the study sample is not similar to that in either the wider population or the population that you treat, you should be suspicious that this study may not be applicable to your practice. If the prevalence of disease in the sample is particularly high it may suggest that it was a highly selected sample.

REFERENCE (OR GOLD) STANDARD

What was the reference standard?

The researchers need to know whether the patient actually has the diagnosis – otherwise they can't comment on the value of the test under investigation (the index test). They should decide whether the subject has the disease by using the best test available – termed the reference or gold standard. Although we might argue that the best reference standard is a post-mortem, this is not always appropriate! As an example, if we were looking at a new test for pulmonary embolism (PE) we would want to design a trial using CT pulmonary angiogram (CTPA) as our reference standard rather than using d-dimer alone.

We might use a combination of tests as our reference standard; it might be considered that CTPA plus telephone follow-up at 6 months to ensure the patient wasn't subsequently diagnosed with a PE would be better than CTPA alone. We must be convinced that the gold standard truly is the best mechanism available to make the diagnosis.

Was the reference standard applied to all the subjects?

Every subject in the study needs to have the reference standard test so that the index test can be compared to it. If everyone does not get the reference standard test, then bias will be introduced to the results. You might ask why everyone *wouldn't* be given the reference standard in the study. This might occur due to failures in processing the subjects during the trial or because some of the researchers came to believe that the index test was so accurate that they didn't want to expose the subjects to a reference standard that has risks (e.g. the radiation of a CTPA or indeed the cost of the test).

Was there only one reference standard?

There are often side effects to an investigation, for example the radiation dose associated with CTPA. In our d-dimer example, researchers might design a study where the reference standard for those with a positive d-dimer is CTPA but those with a negative d-dimer receive just telephone follow-up at 6 months. Although the reasoning can be understood it still needs to be recognized as a flaw and having introduced bias into the study.

Was the reference standard independent of the test under investigation?

In some studies there is a blurring of the lines between the index test and the reference standard. Knowing the result of the index test may affect the interpretation of the gold standard and vice versa.

BLINDING

Who knew what? To avoid bias, all the participants involved in the trial (people conducting and interpreting the index test, people conducting and interpreting the reference standard, and the patients) ought to be blinded to the results of the tests. At the least, those conducting the index test should be unaware of the results of the reference test and vice versa.

BIASES SPECIFIC TO DIAGNOSTIC STUDIES

Bias has been discussed in Chapter 5 and so we will not cover it in detail here other than to mention the types of bias that you should specifically be aware of and look out for when reading diagnostic studies. Remember that bias occurs when the trial design, method or interpretation of results is done in such a way as to build in a flaw that might influence the result. Table 7.1 reviews the various ways bias can be introduced.

TABLE 7.1
Examples of bias in diagnostic studies

Type of bias	Explanation
Work-up bias	Patients with a positive index test result get one reference standard test whereas those with a negative index test result get another, for instance our example of a study looking at d-dimer where those with a positive d-dimer get a CTPA but those with a negative d-dimer get telephone follow-up at 6 months.
Incorporation bias	The test under investigation is incorporated into the reference test. In some situations the reference standard may be an ‘expert’s’ opinion. Consider deciding whether someone is medically fit for psychiatric assessment having attended the ED drunk and suicidal. We want to know whether an alcohol breath test can be used to determine medical fitness in these patients and we look to compare it against our gold standard. However, our gold standard is our current practice: clinical assessment combined with an alcohol breath test. Clearly the alcohol breath test is both the index test and part of the gold standard.
Diagnostic review bias	This is when the results of the index test (test under investigation) are known to the person interpreting the gold standard test. This knowledge may influence the person in their interpretation of the gold standard test, e.g. the radiologist knowing the d-dimer result when interpreting the CTPA.
Test review bias	The opposite of the above – the person interpreting the index test result knows the result of the gold standard.
Clinical review bias	Where those interpreting either the gold standard or index test result are allowed to know clinical information about the patient. Remember that when interpreting test results people are naturally influenced by the clinical information about the patient including such things as age, sex, BMI and symptoms.

QUALITIES OF A TEST

There are three important qualities that must be assessed when looking at a test: reliability, validity and practicality.

Reliability

How does changing the operator of the test affect the result?

We’re interested in how good a test is at picking up disease, but one very important aspect of this is whether if two people use the test they are likely to get the same result. If we are doing a blood test with automated machines we would expect the result of two tests on the same sample to be the same within the given limits of the machine (which the manufacturer should state). However, where we are using a test that depends on the user’s technical ability (ultrasound being a very good EM example) the result may not be the same between users. For example, when one person performs a FAST scan they might pick up blood in the abdomen; if it is repeated by a colleague they may get the same result or they may not find blood. The fact that they disagree with the first person may be due to their colleague’s technical abilities with the ultrasound scanner. The ability of getting consistent results with a test when different people use it is called *inter-observer reliability*.

How reliable is the test done by the same operator?

Not only is it important to know that two operators would get the same result using a test, it is also important that if the same person repeats the test, it will consistently deliver similar (or the same) results. This is called *intra-observer reliability*.

What is the Kappa?

Inter-observer and intra-observer reliability are measured using the Kappa value; we do not need to know how to calculate it for the FCEM exam but should have an understanding of what different Kappa values mean. See Table 7.2.

TABLE 7.2
Interpretation of Kappa value

Kappa value	Agreement
0	Chance agreement
0–0.2	Poor
0.2–0.4	Fair agreement
0.4–0.6	Moderate agreement
0.6–0.8	Good agreement
0.8–0.99	Very good agreement
1	Perfect

Validity

Validity is the ability of a test to accurately measure what it is supposed to measure. A valid test has to be able to show a difference if there is a difference present.

Practicality

Practicality is the ability of the test to be used in clinical practice. We must consider the economic costs, side effects and difficulty of performing the test when assessing this.

HOW TO USE AND INTERPRET RESULTS

As with the results of a therapeutic study, you could be forgiven for feeling anxious about having to calculate the results yourself. However, here we show how you can do this simply, guiding you through step by step, and by the end you will have a thorough understanding of how to interpret the results of a diagnostic study.

WHICH RESULT IS FOR WHOM?

Let's start by considering from whose perspective we are looking at the test (Figure 7.1).

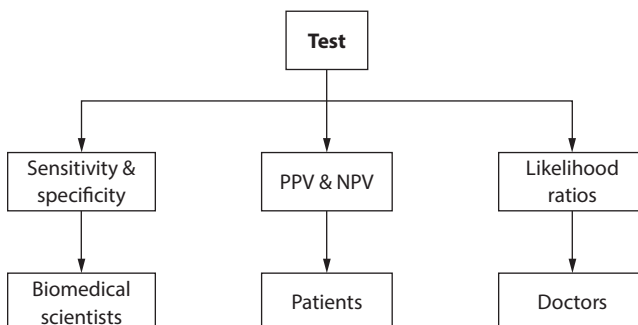


FIGURE 7.1 Different viewpoints demanding different results.

The clinician

The clinician will already have an idea about how likely a diagnosis is, based on their clinical evaluation of the patient (the history and examination). What they want the test to tell them is how much more or less likely the diagnosis is now that they can combine their clinical findings with the test result.

The patient

The patient is only interested in the effect that the test result has on them. They will ask 'If the test comes back as positive, how likely is it that I have the disease?' or 'If the test is negative, how confident can I be that I don't have the disease?'

The biomedical scientist

The biomedical scientist, the person responsible for buying the test, wants an overall assessment of how good a test is at ruling out or ruling in pathology. They need an overall answer rather than one that is applicable to one subset of the population.

Although sensitivity and specificity are often quoted in diagnostic papers, it's clear that we also need to know the positive and negative predictive values and the likelihood ratios. These are surprisingly frequently not given. Therefore we need to be able to calculate them for ourselves and this in turn makes them ideal for examiners to ask about – fortunately this isn't too difficult to do.

CALCULATING RESULTS FROM THE DATA

As with therapeutic trials, drawing a 2 × 2 table will be the key to answering any questions on the results section of diagnostic papers. Patients are classified into some simple definitions that need to be committed to memory:

- Case positive: person with disease (positive reference standard)
- Case negative: person without disease (negative reference standard)
- Test positive: person with positive test under investigation
- Test negative: person with negative test under investigation

Each patient who gets the test under investigation and the gold standard will fall into one of four categories:

- True positive: positive test result and has disease (A)
- True negative: negative test result and does not have disease (D)
- False positive: positive test result but does not have disease (B)
- False negative: negative test result but has disease (C)

The number of patients falling into each of these four categories can then be put into the 2 × 2 table. We need to have a thorough understanding of how to draw one, and remember, however it is drawn in the paper, we recommend drawing it in the way shown here.

		Gold Standard/Have the disease	
		Disease positive	Disease negative
Index test	Positive	A	B
	Negative	C	D

		Gold Standard/Have the disease	
		Disease positive	Disease negative
Index test	Positive	True positive	False positive
	Negative	False negative	True negative

Calculations based on the 2 × 2 table

Before doing any calculations we need to ensure we understand what we're calculating rather than relying upon rote learning of formulae. The FCEM

examination questions will be phrased so that you will get maximal points if you can demonstrate your understanding. When answering the FCEM questions, we feel it's important to show all the workings, starting with the definition.

Prevalence

From the 2×2 table, we can calculate the prevalence of the disease on the population examined.

$$\begin{aligned}\text{Prevalence} &= \text{Proportion of patients with condition (as defined by gold standard)} \\ &= (a + c) \div (a + b + c + d)\end{aligned}$$

This is not traditionally quoted but is a good way for FCEM examiners to see if you really understand the 2×2 table.

Sensitivity

$$\begin{aligned}\text{Sensitivity} &= \text{Proportion of people with the disease who test positive} \\ &= a \div (a + c)\end{aligned}$$

A good way to think of this is as the true positive rate, i.e. proportion of positive tests in people with the disease.

If you have a highly sensitive test and get negative results, it effectively rules out the condition. Many people remember this by the phrase SnOut (Sensitivity rules Out).

Specificity

$$\begin{aligned}\text{Specificity} &= \text{Proportion of people without the disease who test negative} \\ &= d \div (b + d)\end{aligned}$$

A good way to think of this is as the true negative rate, i.e. proportion of negative tests in people without the disease.

If you have a highly specific test and get a positive result, it effectively rules in the diagnosis. Many people remember this by the phrase SpIn (Specificity rules In).

Remember, as a biochemical scientist, the sensitivity and specificity are what are crucial in telling you how good the test is. These results are independent of prevalence. This is key as in whatever population the test is done, the specificity and sensitivity do not change.

Positive predictive value (PPV)

$$\begin{aligned}\text{Positive predictive value (PPV)} &= \text{Proportion of people with a positive test who} \\ &\quad \text{have the disease} \\ &= a \div (a + b)\end{aligned}$$

Negative predictive value (NPV)

Negative predictive value (NPV) = Proportion of people with a negative test who don't have the disease

$$= d \div (c + d)$$

Sensitivity and specificity are independent of prevalence but PPV and NPV are not. As such the PPV and NPV may significantly change if we change the population in which the test is done. For example if a study of a new troponin test was performed on patients with chest pains presenting in primary care, it is unlikely to have the same PPV and NPV as if the study were performed on patients presenting to the ED. However, the sensitivity and specificity (and likelihood ratios) would still be the same.

WORKED EXAMPLE 7.1

NEW TROPONIN TEST

Question

A new ‘trop now’ test has come on the market. Of 136 people who were in the study all had trop now and all had an angiogram (gold standard). Of 33 people who had a positive angiogram, the trop now test was positive in 32 of those people. Of the 103 who had a negative angiogram, the trop now test was positive in 2 people.

Draw a 2 × 2 table and then work out prevalence, sensitivity and specificity.

Answer

		Disease (by Gold Standard)		
		Present	Absent	
Test	Positive	A (32)	B (2)	34
	Negative	C (1)	D (101)	102
	Total	33	103	136

Prevalence = Proportion of patients with condition

$$= (a + c) \div (a + b + c + d)$$
$$= 33 \div 136$$
$$= 0.24 \text{ (24\%)}$$

Sensitivity = Proportion of people with the disease who test positive

$$= a \div (a + c)$$
$$= 32 \div 33$$
$$= 0.97$$

$$\begin{aligned}
 \text{Specificity} &= \text{Proportion of people without the disease who test negative} \\
 &= d \div (b + d) \\
 &= 101 \div 103 \\
 &= 0.98
 \end{aligned}$$

$$\begin{aligned}
 \text{Positive predictive value (PPV)} &= \text{Proportion of people with a positive test who have the disease} \\
 &= a \div (a + b) \\
 &= 32 \div 34 \\
 &= 0.94
 \end{aligned}$$

$$\begin{aligned}
 \text{Negative predictive value (NPV)} &= \text{Proportion of people with a negative test who don't have the disease} \\
 &= d \div (c + d) \\
 &= 101 \div 102 \\
 &= 0.99
 \end{aligned}$$

Earlier in this chapter we introduced two terms: true positive rate (sensitivity) and true negative rate (specificity). Before going onto likelihood ratios we now need to introduce two new terms. These terms are not often thought about but make understanding of likelihood ratios much easier. They are false positive rate and false negative rate.

$$\begin{aligned}
 \text{False positive rate} &= \text{Proportion of positive tests in people without the disease} \\
 &= b \div (b + d) \\
 &= 1 - \text{Specificity}
 \end{aligned}$$

$$\begin{aligned}
 \text{False negative rate} &= \text{Proportion of negative tests in people with the disease} \\
 &= c \div (a + c) \\
 &= 1 - \text{Sensitivity}
 \end{aligned}$$

For those who are interested, we explain how false positive rate = 1 – specificity and how false negative rate = 1 – sensitivity in ‘Spod’s Corner’, but remember this is not needed for the exam.

SPOD’S CORNER

(This will not be tested in the FCEM examination.)

False positive rate = 1 – specificity.

To prove this we need to prove that $b/(b + d) = 1 - d/(b + d)$.

Step 1: Multiply both sides by $(b + d)$.

$$(b + d) \times [b/(b + d)] = (b + d) \times [1 - d/(b + d)]$$

Step 2: $b = b + d - d$

Step 3: $b = b$

Therefore $b/(d + b) = 1 - d/(b + d)$.

And therefore false positive rate = $1 - \text{specificity}$.

False negative rate = $1 - \text{sensitivity}$.

To prove this we need to prove that $c/(c + a) = 1 - a/(a + c)$.

Step 1: Multiply both sides by $(a + c)$.

$$(a + c) \times [c/(c + a)] = (a + c) \times [1 - a/(a + c)]$$

Step 2: $c = a + c - a$

Step 3: $c = c$

Therefore $c/(c + a) = 1 - a/(a + c)$.

And therefore false negative rate = $1 - \text{sensitivity}$.

Likelihood ratios

As doctors what we want to know is how the result of the test affects whether the patient in front of us really does have the condition. We can use likelihood ratios to help us with this. By combining a pre-test assessment (such as our clinical gestalt or a formalized scoring system such as a Well's score) with a likelihood ratio we can estimate how likely the patient in front of us is to have the condition.

Likelihood ratios can be either positive or negative. We use positive likelihood ratios when the test result is positive and negative likelihood ratios when the test result is negative.

Positive likelihood ratio: How much more likely is a positive test to be found in a person with the disease compared to a positive test in someone without the disease.

Another way to think of this is:

true positive rate \div false positive rate = sensitivity \div (1 – specificity).

Negative likelihood ratio: How much more likely is a negative test to be found in a person with the disease compared to a negative test in someone without the disease.

Again, a way to think of this is:

false negative rate \div true negative rate = (1 – sensitivity) \div specificity

Likelihood ratios are important tools in a diagnostic paper. They tell us how diagnostically useful a piece of information is (whether it is an examination finding, part of history or a test). Positive likelihood ratio values can be between 1 and infinity; the greater the number the better the test. Negative likelihood ratio values can be between 0 and 1; the closer to 0, the better the test (see Table 7.3). Similar to sensitivity and specificity, likelihood ratios are independent of prevalence.

TABLE 7.3
Values for likelihood ratios and their significance

Likelihood ratio	Significance
1	None at all
0.5 to 2	Little clinical significance
2 to 5	Moderately increases likelihood of disease, but does not rule in
0.2 to 0.5	Moderately decreases likelihood of disease, but does not rule out
5 to 10	Markedly increases likelihood of disease, may rule in
0.1 to 0.2	Markedly decreases likelihood of disease, may rule out
>10	Near enough rules in
<0.1	Near enough rules out

To put this in a clinical context, if we were trying to diagnose a patient with a *deep vein thrombosis* (DVT) we might use clinical assessment findings such as those in Table 7.4 to help us make our assessment. For example having a previous DVT makes you more likely to have a DVT by a likelihood ratio of 2.5. A common misconception is that this means that those who have had a DVT are 2.5 times more likely to now have a DVT than those that didn't have one previously. Actually, what this means is that the patient is more likely to have a DVT, but how much more likely is complex, and we explain why later in this chapter.

In Example 7.2, we show you how to calculate likelihood ratios, which could be a question in the FCEM examination.

TABLE 7.4
Likelihood ratios in a clinical context: diagnosis of a DVT based on clinical assessment

Feature	Likelihood ratio of positive finding
Past history VTE	2.5
Malignancy	2.6
Immobilization	1.9
Recent surgery	1.7
Difference in calf diameter	1.8
Homan's sign	1.4
Oedema	1.2

WORKED EXAMPLE 7.2

CALCULATING LIKELIHOOD RATIOS

Question

Using the example of the trop now test given in Example 7.1 calculate the likelihood ratio of a positive result and the likelihood ratio of a negative results. The 2 × 2 table is shown here.

		Disease (by Gold Standard)		
		Present	Absent	
Test	Positive	A (32)	B (2)	34
	Negative	C (1)	D (101)	102
	Total	33	103	136

Answer

Positive likelihood ratio

= How much more likely is a positive test to be found in a person with the disease compared to a positive test in someone without the disease?

Positive likelihood ratio = True positive ÷ False positive
= Sensitivity ÷ (1 – Specificity)
= 0.97/(1 – 0.98)

= 48.5

Negative likelihood ratio

= How much more likely is a negative test to be found in a person with the disease compared to a negative test in someone without the disease?

Negative likelihood ratio = False negative ÷ True negative
= (1 – Sensitivity) ÷ Specificity
= (1 – 0.97)/0.98
= 0.03

Pre-test probabilities to post-test probabilities

Without placing them in a clinical context the uses of likelihood ratios are limited; we generally don’t perform tests randomly on patients. We assess patients and decide on their pre-test probability of having a disease. However, what we really want to know is how likely is the patient in front of me, after my clinical assessment and the test, to have the disease; i.e. the post-test probability.

There are two ways of getting to the post-test probability (shown in Chapter 4). The ‘100%’ accurate way is by using a four-step process:

- 1. Estimate the pre-test probability using our clinical judgment
- 2 Convert pre-test probability to pre-test odds

3. Multiply pre-test odds by the likelihood ratio to calculate post-test odds
4. Convert post-test odds to post-test probability

However, this is complex and is not needed for the FCEM examination. It also is not really 100% accurate as it depends on the pre-test probability, which we as clinicians *estimate* anyway. For those interested, though, how this is actually done is shown in the following ‘Spod’s Corner’.

The second and simpler way is by using the likelihood ratio nomogram. The nomogram effectively creates a probability to odds table for us and we can use it to give us a good estimation of post-test probability.

SPOD’S CORNER: CALCULATING THE POST-TEST PROBABILITY USING ODDS

Using the example of the trop now test given in Example 7.1, calculate the post-test probabilities for a positive or negative test for trop now for the following two clinical examples.

The formulas for converting odds to probability and vice versa are as follows:

$$\text{Odds} = \text{Probability} \div (1 - \text{Probability})$$

$$\text{Probability} = \text{Odds} \div (\text{Odds} + 1)$$

$$\text{Post-test odds} = \text{Pre-test odds} \times \text{Likelihood ratio}$$

As calculated in Example 7.2: +ve LR = 48.5 and –ve LR = 0.03.

High-risk man

A 65-year-old overweight diabetic smoker presents with tightening chest pain. Your clinical gestalt says he has a 95% chance of having acute coronary syndrome.

$$\text{Pre-test probability} = 0.95$$

$$\text{Pre-test odds} = 0.95 \div (1 - 0.95) = 19 = (19:1)$$

$$\text{Post-test odds} = \text{Pre-test odds} \times \text{Likelihood ratio}$$

+ve results

$$\text{Post-test odds} = 19 \times 48.5 = 921.5 = (921.5:1)$$

$$\begin{aligned} \text{Post-test probability} &= \text{Post-test odds} \div (\text{Post-test odds} + 1) \\ &= 921.5 \div (921.5 + 1) \\ &= 0.999 = 99.9\% \end{aligned}$$

–ve results

$$\text{Post-test odds} = 19 \times 0.03 = 0.57 = (0.57:1)$$

$$\begin{aligned}\text{Post-test probability} &= \text{Post-test odds} \div (\text{Post-test odds} + 1) \\ &= 0.57 \div (0.57 + 1) \\ &= 0.36 = 36\%\end{aligned}$$

Low-risk man

A 20-year-old man just finished a marathon and was tired. He has a normal ECG. Your clinical gestalt says he has a 0.1% chance of having acute coronary syndrome.

$$\text{Pre-test probability} = 0.001$$

$$\text{Pre-test odds} = 0.001 \div (1 - 0.001) = 0.001001$$

$$\text{Post-test odds} = \text{Pre-test odds} \times \text{Likelihood ratio}$$

+ve results

$$\text{Post-test odds} = 0.001001 \times 48.5 = 0.0485485$$

$$\begin{aligned}\text{Post-test probability} &= \text{Post-test odds} \div (\text{Post-test odds} + 1) \\ &= 0.0485485 \div (0.0485485 + 1) \\ &= 0.046 \\ &= 4.6\%\end{aligned}$$

–ve results

$$\text{Post-test odds} = 0.001001 \times 0.03 = 0.00003$$

$$\begin{aligned}\text{Post-test probability} &= \text{Post-test odds} \div (\text{Post-test odds} + 1) \\ &= 0.00003 \div (0.00003 + 1) \\ &= 0.00003 \\ &= 0.003\%\end{aligned}$$

Calculating post-test probabilities using the nomogram

There are three lines in the nomogram. The left-hand line is the pre-test probability. The middle line is the likelihood ratio and the right-hand line is the post-test probability (Figure 7.2). Example 7.3 shows worked examples.

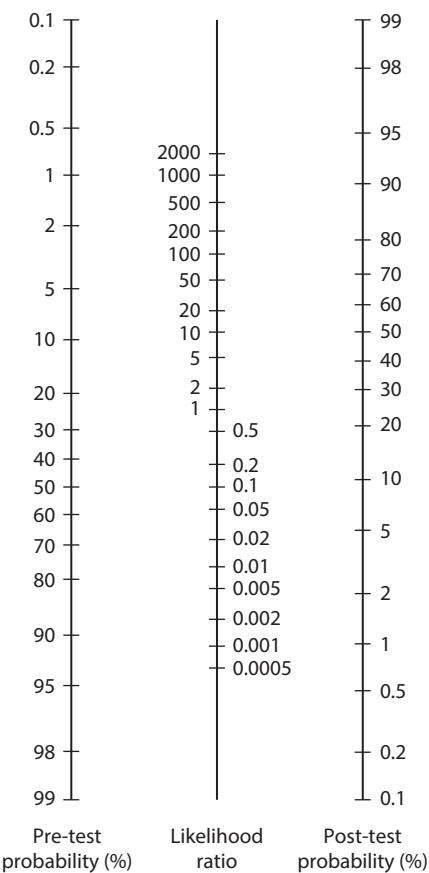


FIGURE 7.2 The nomogram for converting pre-test probabilities and likelihood ratios to post-test probabilities.

WORKED EXAMPLE 7.3

CALCULATING THE POST-TEST PROBABILITY USING THE NOMOGRAM

Question

Using the nomogram shown in Figure 7.2, calculate the post-test probabilities for the following patients using the trop now test, which we know to have a positive likelihood ratio of 48.5 and a negative likelihood ratio of 0.03.

- a. A high-risk patient (pre-test probability guestimate of 95%) with a positive test
- b. A high-risk patient (pre-test probability guestimate of 95%) with a negative test

- c. A low-risk patient (pre-test probability guestimate of 0.1%) with a positive test
- d. A low-risk patient (pre-test probability guestimate of 0.1%) with a negative test

Answer

The lines are shown in Figure 7.2. The answers are all approximations based on how accurate you drew the intercepts. But they show that the answers correlate with the exact calculations.

- a. >99% (Figure 7.3)
- b. 36% (Figure 7.4)
- c. 4.5% (Figure 7.5)
- d. <0.1% (Figure 7.6)

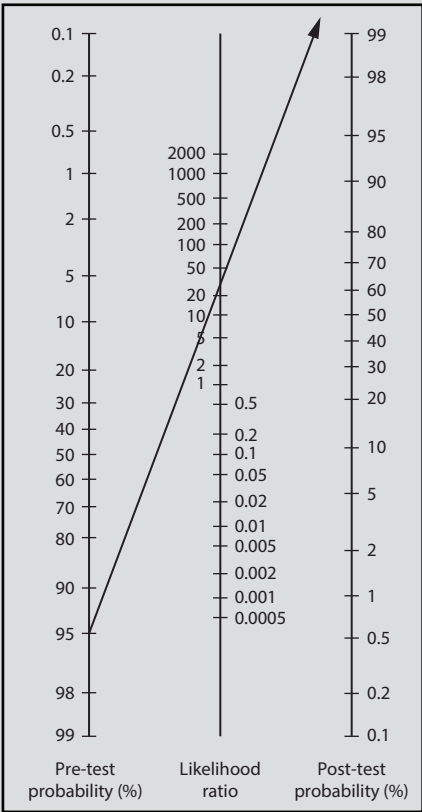


FIGURE 7.3

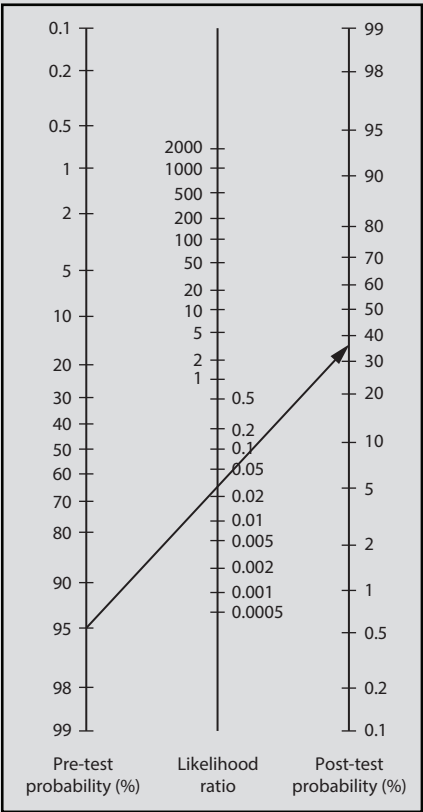


FIGURE 7.4

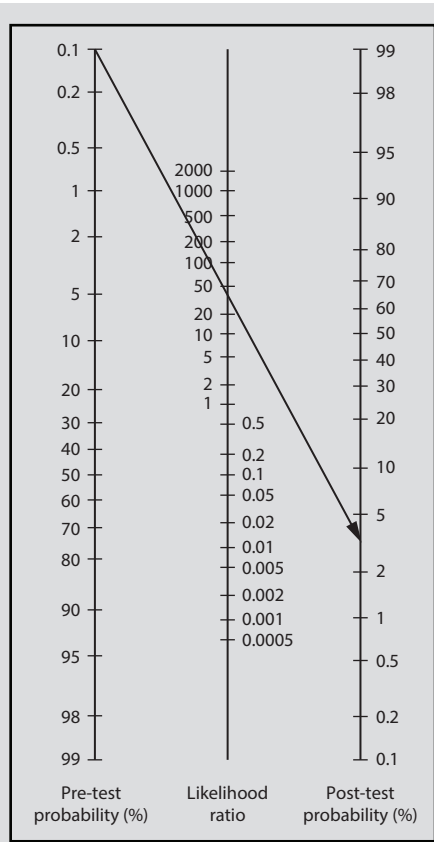


FIGURE 7.5

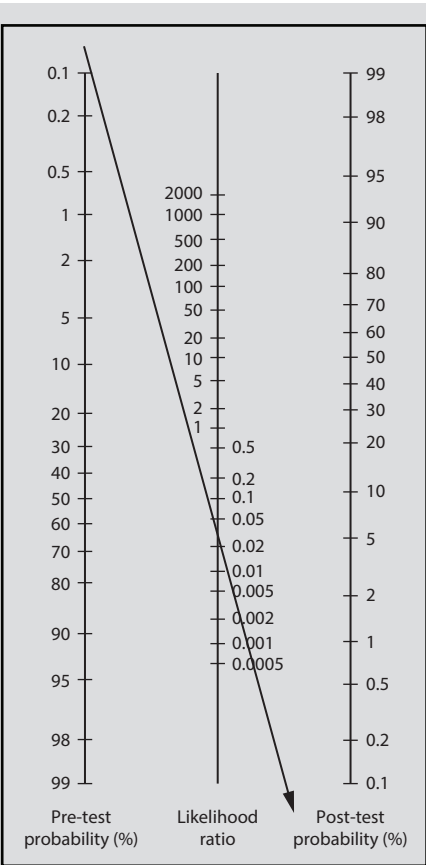


FIGURE 7.6

Before we move on, for those who want to stretch themselves it’s interesting to consider the effect of pre- and post-test probabilities if we were dealing with a truly random patient; see the following ‘Spod’s Corner’.

SPOD’S CORNER: PRE- AND POST-TEST PROBABILITIES ON TRULY RANDOM PATIENTS

In reality we evaluate patients and make our own pre-test probabilities, but what if we had a truly random patient? So using our trop now example:

		Disease (by Gold Standard)		
		Present	Absent	
Test	Positive	A (32)	B (2)	34
	Negative	C (1)	D (101)	102
	Total	33	103	136

Pre-test probability is the chance of the patient you see having the condition.

$$\begin{aligned}
 &= \text{Number of patients with the condition} \div \text{total number of patients} \\
 &= (a + c) \div (a + b + c + d) = \text{Prevalence} \\
 &= 33 \div 136 \\
 &= 0.243
 \end{aligned}$$

Pre-test odds are the odds of the patient you see having the condition.

$$\begin{aligned}
 &= \text{Number of patients with the condition} \div \text{number of patients without the condition} \\
 &= (a + c) \div (b + d) \\
 &= 33 \div 103 \\
 &= 0.32 \\
 &= 0.32:1
 \end{aligned}$$

From previous calculations LR positive is 48.5 and LR negative is 0.03.

+ve results

$$\begin{aligned}
 \text{Post-test odds} &= \text{Pre-test odds} \times \text{Positive LR} \\
 &= 0.32 \times 48.5 = 15.52
 \end{aligned}$$

$$\begin{aligned}
 \text{Post-test probability} &= \text{Post-test odds} \div (\text{Post-test odds} + 1) \\
 &= 15.52 \div (15.52 + 1) \\
 &= 0.94 = 94\% \\
 &= \text{Positive predictive value}
 \end{aligned}$$

-ve results

$$\text{Post-test odds} = \text{Pre-test odds} \times \text{Negative LR}$$

$$\text{Post-test odds} = 0.32 \times 0.03 = 0.0096$$

$$\begin{aligned}
 \text{Post-test probability} &= \text{Post-test odds} \div (\text{Post-test odds} + 1) \\
 &= 0.0096 \div (1 + 0.0096) \\
 &= 0.01 = 1\% \\
 &= 1 - \text{Negative predictive value (i.e. chance of a negative test if you do have the disease)}
 \end{aligned}$$

ROC (RECEIVER OPERATED CHARACTERISTIC) CURVES

Some tests have a yes or no answer (e.g. does the ultrasound show a pneumothorax?). Other things, such as clinical scores and assays, don't have a yes or no answer, so we need to choose a cutoff value to decide what is a positive and what is a negative result. For each different cutoff, the test will give you a different sensitivity, specificity, PPV, NPV and likelihood ratios.

To decide what the best cutoff value is, we draw a receiver-operated characteristic (ROC) curve. To draw the ROC curve we plot the sensitivity (true positive rate) against $1 - \text{Specificity}$ (false positive rate) for all potential cutoff values of the test. The point on the curve that is diagonally closest to the maximum sensitivity is the statistically best cutoff value for the test (Figure 7.7). However, this may not be clinically the most appropriate level for the cutoff.

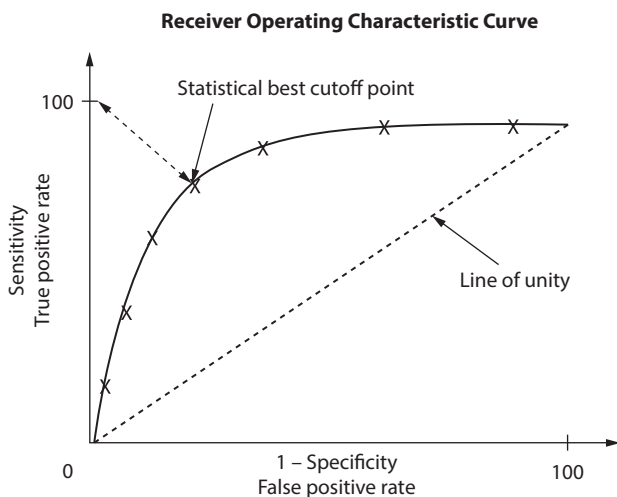


FIGURE 7.7 The receiver operated curve (ROC) with best statistical cutoff point.

For example, although trop now has an optimal statistical cutoff of 17 ng/L, we may decide to not use this as our cutoff value in clinical practice. We may choose a lower value (a cutoff of 10 ng/L for example) as this will give us a better sensitivity (i.e. you are less likely to miss any patients with acute coronary syndrome [ACS]), with the trade-off that this will also give us a worse specificity (i.e. more false positives and more unnecessary admissions). This will make the test a better rule-out test for ACS. Conversely, if we chose a higher value for trop now, it would make it a better rule-in test.

Area under the curve (accuracy of the test)

The area under the curve (AOC) is a measure of how accurate the test is.

1 = perfect

0.5 = useless test - no better than chance and the line would lie along the line of unity

TAKE HOME MESSAGE

1. Diagnostic trials should compare the test under investigation to a gold standard.
2. For there to be internal validity in the study, the researchers must ensure that each patient gets the test under investigation and the gold standard.
3. If the prevalence in the study is very dissimilar to the prevalence in your population the study will not be externally valid.
4. In any diagnostic study, the authors should state the sensitivity, specificity, PPV, NPV and likelihood ratios. You must know how to calculate these for the FCEM exam.
5. ROC curves help us to decide the optimal statistical cutoff value for a test. This may not be the same as the optimal clinical cutoff value.

8 Meta-analysis

Although reviews and meta-analyses have never come up in the critical appraisal portion of the FCEM examination, knowing about them will be helpful for your CTR and life in general. In this chapter we discuss different types of reviews and go through the stages of writing a meta-analysis.

DIFFERENT TYPES OF REVIEW ARTICLES

A systematic review is a scientific study using the IMRaD (introduction, methods, results and discussion) format discussed in Chapter 5. It addresses a specific question, using existing data, and as such is *secondary* research rather than *primary* research. A *systematic* review should be an unbiased synthesis of available evidence and every effort should be made to prevent bias during the process of writing it. A meta-analysis is a specific form of systematic review with a statistical combination of the results of all the appropriate trials in the review.

The main aim of meta-analysis is to increase the precision of the estimates and hence reduce type 2 error. They do not, however, overcome bias, and care must be taken to prevent introducing new forms of bias. A common criticism of meta-analyses is that in combining the results from a number of studies with flaws, we create one big but substantially flawed study. A meta-analysis is only as good as the papers that contribute to it.

A narrative review is not a systematic review; rather it is a broad overview of an issue. The preparation is generally not as rigorous as a systematic review; it may contain the author's opinion and only selected data may be presented. As such it may not necessarily be objective.

WRITING A META-ANALYSIS

There are four steps in creating a meta-analysis: (1) searching the literature, (2) selection of papers, (3) abstraction of data and (4) analysis.

SEARCHING THE LITERATURE

When writing a meta-analysis, the first step is to perform a literature review with respect to the clinical question. As with writing a CTR, a literature review should include all the common databases, relevant papers from the references of papers that have already been chosen and a review of the grey literature. Grey literature is data that is not easily available and could include unpublished studies, conference abstracts, drug company promotional material, etc. How and where the authors of the paper have searched must be explicitly stated.

We are always worried when reading a meta-analysis that relevant studies have not been included. This is generally not because the authors have not performed an appropriate literature review but rather because of publication bias. This occurs when the relevant studies are either not written up, submitted or accepted for publication or because they are published in an obscure journal or in a foreign language. To demonstrate the presence or absence of publication bias, a funnel plot should have been performed and presented. This is a plot showing a measure of study precision against effect size. For ‘Spod’s Corner’ aficionados it is worth noting that study precision is the inverse of the standard error of the result ($\text{Precision} = 1 \div \text{Standard error of the result}$).

Figure 8.1 shows an example of two funnel plots. The funnel plot should be the shape of an inverted cone (a funnel), and asymmetry suggests the presence of publication bias.

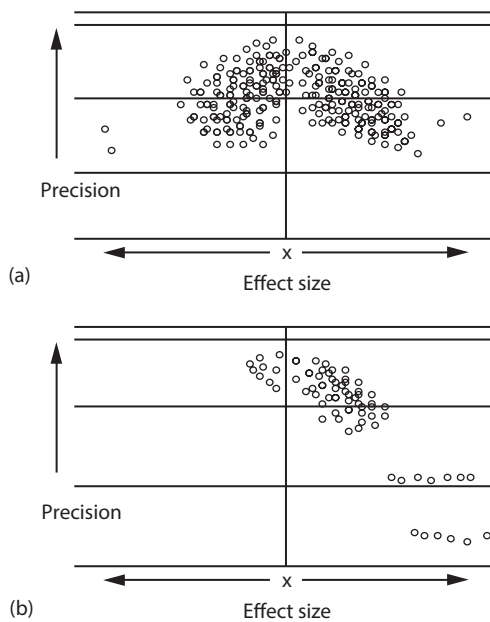


FIGURE 8.1 (a) A symmetrical funnel plot. (b) An asymmetrical funnel plot.

SELECTION OF PAPERS

Once all the potentially relevant papers have been gathered, then the right papers need to be selected. To do this we have inclusion and exclusion criteria similar to an interventional trial. These are known as eligibility criteria and examples include statistics quoted, study design, year of study, language of the article, etc. More than one author should perform an independent selection of papers, and discussion between the authors or use of a third author to resolve conflicts in selection is often necessary – this prevents bias. A Kappa value should be calculated to measure the agreement between the authors, and the closer this value is to 1 the better the agreement.

ABSTRACTION OF DATA

Once the appropriate group of studies has been identified, the author(s) have to abstract the relevant data from the studies. This can lead to misinterpretation and error, so mechanisms should be put in place to identify any errors. More than one person should independently extract data from every report. Forms for data extraction should be piloted to ensure that they do not miss data or collect unneeded data. There should also be a procedure for identifying and resolving differences in the data.

In addition to gathering the appropriate data from the appropriate studies, the studies should be assessed with respect to their quality. There is no set method for this, although a formalized grading scheme is usually used, e.g. the Jadad score (see Glossary, Appendix B).

ANALYSIS

Analysis occurs in two stages. First the statistic of interest is calculated for each study included. The pooled intervention effect is calculated as a weighted average of the intervention effects in each individual study. The standard error of this pooled intervention effect can be used to calculate CIs and p-values in exactly the same way we have described previously. Care must be taken to combine results that are the same, for example odds with odds and risks with risks. The study data is often presented in the form of a forest plot (see Figure 8.2).

Even though studies in a meta-analysis are similar, there will be differences in the variability among the studies and this is known as heterogeneity. It is measured using

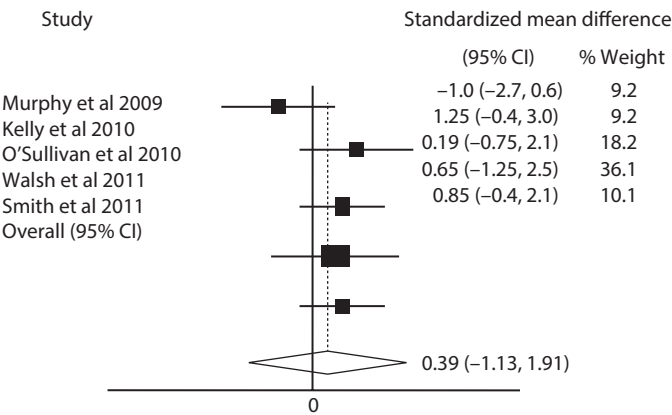


FIGURE 8.2 Forest plot. The left-hand column lists the name of the studies. The right-hand column plots the measure of effect (e.g. odds ratio). The area of each square is proportional to the study weight in the meta-analysis. The lines extending from the box are the CIs. The line in the centre is the line of zero effect. The diamond is the cumulated effect and this is extended upwards as a dotted line, which is the overall result of the meta-analysis. The lateral tips of the diamond are the CIs for the overall effect. In the example, because the CIs cross the line of zero effect, this result would be deemed non-significant.

an I^2 value and is given as a percentage. The lower the value the better; greater than 75% is considered significant. If the CIs of the individual studies have poor overlap this may be a sign of heterogeneity. If a meta-analysis shows heterogeneity then the authors must decide how best to proceed; examples of what they can do include:

- Recheck the data looking for errors
- Change the effects measure
- Exclude studies
- Explore the heterogeneity using subgroup analysis \pm performing a random effects meta-analysis
- Ignore the heterogeneity
- Do not do the meta-analysis

Although it is possible to combine heterogeneous studies, doing so is controversial and has a significant impact on the validity of the results.

TAKE HOME MESSAGE

1. A meta-analysis is a form of systematic review and as such should be evidence based using rigorous mechanisms to prevent the introduction of bias.
2. There are four steps to writing a meta-analysis: literature review, selection of papers, abstraction of data and analysis of data.
3. A funnel plot can be used to look for publication bias and we conclude bias is present when the funnel plot is not symmetrical.
4. The results of a meta-analysis are generally presented in the form of a forest plot.

Section 3

Passing FCEM

9 Passing the exam

Let's face it – passing the exam is why you've bought this book. The thought of missing out on a consultant post and not getting rid of your registrar rota because of a lack of knowledge of p-values is not a pleasant one.

The critical appraisal exam (and FCEM as a whole) is an incredibly good exam. It is challenging and requires a high standard of clinical acumen to pass and it covers the areas needed to help us function as senior emergency medicine doctors. Learning how to appraise papers is an essential requirement for a senior doctor – both to shape your practice as a 'best' evidence based clinician and also to help make decisions on departmental policies.

Historically, the critical appraisal exam has had a poor pass rate. Some people feel the exam is too hard – but we disagree; the exam is very fair. The exam is there to test your ability to critically appraise a paper and understand the process. It is not about turning you into a medical statistician. The questions are generally predictable and relevant, and require knowledge and understanding rather than rote formula learning; core concepts are examined and not 'far-out' minutiae.

But most importantly, the exam tests what you need to know – i.e. as consultants how will you assess the papers you read and how will that change your practice and departmental policy?

So why do people fail? There are generally three reasons: lack of preparation leading to lack of knowledge and understanding, lack of exam technique and lack of exam practice.

By this point in the book, we have addressed the knowledge and understanding aspects and over the next two chapters we move on to exam technique and practice.

Because the exam is often tight for time, those who have done multiple mock exams often do the best. There are lots of mock exams out there – practise them in exam conditions with colleagues and discuss the answers together. Even better, find recent relevant papers, set your own exams and go over the answers. By practising multiple papers, you will embed your skills and knowledge to make sure you pass. We have provided some mock examinations for you to use at the end of this section.

EXAM TECHNIQUE

Everybody who goes on an exam course or reads an exam preparation book thinks the exam technique section will just be common sense. There is no exception here, but (and a big but) there are specific tactics we used to pass FCEM and we advise others to use – and these have worked well for many candidates.

The first thing to remember is that the exam is often very tight for time – especially if there is a long journal article. You must use your time very carefully. The key is

reading the questions first and *then* reading the paper so you know where to concentrate your efforts. While reading the journal article, use a highlighter to make a note of key points, and write over the paper which question the key point is relevant to.

The questions are usually predictable and over the years have followed a standard format. (This is true up to the exams in 2014 when this book was written and we see no reason why the standard format would significantly change.)

There are usually six to eight questions. Each carries a certain number of marks; the number of marks is stated alongside the question. Trying to predict the pass mark is usually not helpful. Just try and get as many points for each question as possible.

The best way to pass is make sure you get the full marks for the questions that require the least amount of effort. These are often the first and last questions, and we suggest you consider answering these initially, as you don't want to run out of time and not have answered the 'easier-to-score' questions.

When answering questions, remember that you need to write in pencil (so you can rub out) and also you must stay within the boxes and (for question 1) word count. The boxes and associated marks for each question give an indication of how long the answer needs to be. It goes without saying that you should write in a sensible and readable size – sadly, people have actually failed because of their handwriting. The examiners have multiple papers to mark, and so you need to make it easy for them to give you the marks. Write clearly and legibly; when practising exams do it in pencil to get used to it. Writing in bullet points and underlining key points is acceptable, as is writing in capital letters.

Finally, just in case, bring your own pencil, rubber, sharpener and highlighter. Unfortunately, *you will not be allowed to take in calculators*.

Prior to answering any of the questions, there is one thing you must do. *It is imperative that you work out very quickly what type of paper you are reading*. Invariably it will either be a diagnostic paper or a therapeutic paper. (However there is no reason why the College cannot use a meta-analysis of these types of papers. Theoretically they could also use any other type of paper – but this has not yet happened and as therapeutic and diagnostic papers are the most important for EM doctors, we would doubt that this format would change.)

It should be very obvious what type of paper it is, but a large number of people who fail do so because they don't recognize what type of paper they are reading and thus how they should answer the questions. So, right at the start, think 'Is there a test being investigated (diagnostic paper) or is the paper looking at differences between two treatments (therapeutic paper)?'

We recommend that as soon as you have recognized what type of paper it is you write this at the top of your answer sheet.

The format of the questions is usually along the following lines:

- Question 1: Always a summary
- Questions 2–4: Usually about design and methodology of the trial
- Questions 4–7: Something about statistics and mathematics, which may be simple calculations, e.g. sensitivity
- Question 8: Your interpretation of the strength and limitations of the paper

QUESTION 1: A SUMMARY (USUALLY 7 MARKS)

These are the easiest marks in the paper to get – so make sure you get them. The question will read something like this:

Provide a no more than 200-word summary of this paper in the box provided. Only the first 200 words will be considered – short bullet points are acceptable. Maximum of seven marks.

Do as the instructions say. Provide a summary – not your interpretation, not the whole paper regurgitated, not if you would implement their findings. Nothing but the summary; a bit like an informal version of what you would expect in the abstract. The beauty of this question is that there are 7 marks for literally copying sentences from the paper – a bit like when the medical senior house officer (SHO) gets praised for copying your clerking.

In many ways this is a brilliant but simple question – it asks if you have understood the paper well enough that you can distil it down to the abstract. You should have no problem with this as long as you have practised doing this a number of times before the exam.

When writing the summary, think about the following points:

- It needs to stand alone, so don't refer to other specific papers.
- There is no need to detail the background information.
- You need to include results as well as statistics (don't just say that the treatment worked or the test was good). You must include the primary outcome, but may wish to include some of the secondary outcomes if they were important.
- Use the author's conclusions, not your own.
- Do not write your opinion of the paper – you will get a chance later in the exam.
- *The key point:* This is not a checklist assessing the quality of the paper. Many people make this simple mistake and lose out on very easy marks.

What goes into a summary depends on the type of paper: diagnostic or therapeutic article.

THERAPEUTIC PAPERS

You should include all the following sections to ensure you cover all the areas that could generate a mark. The College's marking scheme may be shorter than your summary but by including all the details you should cover all the bases.

It is often best if you start by writing out the headings and underlining them and then filling in the details as you go along. This makes it easier for the examiner to see where you have scored points. Bullet points are acceptable. The key sections are in italics:

- Objectives
- *Design*

- Setting
- *Population*
- Methods
- *Intervention/treatment*
- *Control*
- Outcome measurements
- *Results* – state the results that the authors give but it's especially important to give relative risk, number needed to treat and associated statistics.
- *Conclusions*

As with every list in medicine, there are various mnemonics to remember this by. We advise you to make up your own.

The key is to state that it is a therapeutic trial and to state how it was designed. You also must be sure to state the intervention and control. However, it can be hard to work out how much detail to put into this answer. An example is given here.

Austin et al. Effect of high flow oxygen on mortality in chronic obstructive pulmonary disease patients in pre-hospital setting: randomized controlled trial. *BMJ* 2010;341:c5462

Question 1: Provide a no more than 200 word summary of this paper in the box provided. Only the first 200 words will be considered – short bullet points are acceptable. Maximum of seven marks.

Objective	To compare pre-hospital high flow oxygen with titrated oxygen therapy for patients with acute exacerbations of chronic obstructive pulmonary disease
Design	A therapeutic trial – interventional, non-blinded, cluster randomized, controlled, parallel group trial
Setting	Pre-hospital patients attended by the Tasmanian Hospital Service between July 2006 – July 2007
Population	All patients were admitted to the Royal Hobart Hospital, Tasmania 405 patients aged ≥ 35 years with breathlessness and a history or risk of COPD Risk of COPD defined as greater than 10 pack years
Intervention and control	Titration oxygen in the pre-hospital setting High-flow oxygen in the pre-hospital setting
Primary outcome	Pre-hospital and in-hospital mortality
Results	An intention to treat analysis was performed. The mortality in the intervention group was 4% compared to 9% in the control group and this reduction was significant (RR 0.42, 95 CI 0.20–0.89, $p = 0.02$). This was also true in those with confirmed COPD (RR 0.22, 95% CI 0.05–0.91, $p = 0.04$). Titrated oxygen reduced the risk of death by 58% for all patients and 78% for those with COPD.
Conclusions	Compared with high flow oxygen therapy, pre-hospital titrated oxygen therapy significantly reduced mortality in patients with acute exacerbations of COPD.

DIAGNOSTIC PAPERS

For diagnostic papers, we suggest you include all the following sections. Again, as with therapeutic papers, it is often easiest if you start by writing out the headings and underlining them and then filling in the details as you go along. The key sections are in *italics*.

- Objective
- *Design*
- Setting
- *Population*
- Methods
- *Test under investigation*
- *The gold standard* (good idea to mention if the test was done to all patients and if it was done blindly)
- *Results* (PPV, NPV, sensitivity, specificity, LR, Kappa), also quote the relevant statistics
- *Conclusion*

Again, make sure you state that it is a diagnostic study and what is the test under investigation and the gold standard. Remember they may not be as simple as an assay/radiological diagnosis. For example in a test for appendicitis the gold standard could be operative findings plus follow-up for those who got discharged with a diagnosis of no appendicitis. In a similar way, the test may not be so easy to spot as a simple assay – it could be a clinical score, a set of rules or a combination of test and clinical score. This is explained in more detail in Chapter 7.

An example of a diagnostic summary is given here.

Diagnosis of intussusception by physician novice sonographers in the emergency department. Riera et al. *Ann Emerg Med*. 2012 September; 60(3):264–268.

Question 1: Provide a no more than 200 word summary of this paper in the box provided. Only the first 200 words will be considered – short bullet points are acceptable. Maximum of seven marks.

Objective	To investigate the performance characteristics of bedside ultrasound by paediatric emergency physicians who received limited and focused training in the diagnosis of ileocolic intussusception in children
Design	Prospective diagnostic study
Setting	Tertiary care children’s hospital emergency department
Population	Children with suspected ileocolic intussusception who were to undergo ultrasonography in the diagnostic radiology department
Methods	Patients underwent bedside ultrasound performed by one of six emergency physicians attending or fellows who had received a one hour focused training session
Test under investigation	Bedside ultrasound by an emergency physician
Gold Standard	Radiological department ultrasound

(Continued)

Results	Bedside ultrasound had a sensitivity of 85% (95%CI 54–97%), specificity of 97% (95%CI 89–99%), positive predictive value of 85% (95%CI 54–97%), negative predictive value of 97% (95%CI 89–99%), positive likelihood ratio of 29 (95%CI 7.3–11.7) and negative likelihood ratio of 0.16 (95% CI 0.04–0.57)
Conclusion	This prospective observational study demonstrated the good performance characteristics of paediatric emergency physician-performed bedside ultrasonography for the diagnosis of intussusception in children after a single, focused training session.

META-ANALYSIS

Although very unlikely, there is no reason why a meta-analysis could not be used in the exam. If there were a meta-analysis then it would be based on either a series of diagnostic studies or therapeutic studies. You would have to include the same headings as in the summary for the individual paper but be sure to mention some specific extra points. All these areas are described in detail in Chapter 8.

- Paper identification
- Paper selection
- Funnel plot to see missing papers
- Forest plot
- Statistical analysis (including test for heterogeneity)

QUESTIONS 2–4: USUALLY ABOUT DESIGN AND METHODOLOGY (6–10 MARKS)

Knowledge and understanding are needed for this section. There may be questions directly related to the paper and there may be free-standing questions, e.g. explain the importance of randomization.

Again, technique is required. The length of your answer should be proportional to the marks allocated and the size of the box should also give an indication. There are some specific ways of answering the questions to help maximize marks.

1. Read the question. If it asks for a specific number of points give that specific number of points. But if it says list the aspects of the design that were well done, then list all the factors you can think of starting with the most relevant as long as it fits within the box.
2. If it asks for strength of the design, do not write strengths of the paper but concentrate on the design.
3. You must explain ‘buzz’ words and not just put them in, in the hope that you will get some marks. For example instead of stating that the study was pragmatic so generalizable, state ‘the study was pragmatic in the fact that it was done with normal staff, using normal processes, with the type of patients

we see and no special resources. Therefore the conclusions could easily be applied to our everyday practice’.

4. You must give attention to detail to your answers. For example, don’t state that there was no loss to follow-up as an example of the strength of the design. This is careless as it is a *consequence* of the strength of the design. Instead you can say, ‘The design of the study was excellent in that all patients would have a telephone follow-up and if they didn’t answer the phone, a home visit was arranged. This ensured that there was minimal loss to follow-up’.

The key though is to remember that therapeutic and diagnostic papers will have different things to look for when asked to assess the strengths of the design. Always think about design in terms of internal and external validity. Meta-analysis will have additional things to think about. See Appendix A, which goes through a checklist of things to think about when asked to assess how good a study was.

QUESTIONS 4–7: USUALLY SOMETHING ABOUT RESULTS, STATISTICS AND 2×2 TABLES (6–10 MARKS)

These are easy marks to get if you have a good understanding of critical appraisal. You will not be expected to calculate p-values or power calculations, but you will be expected to understand how they are derived. You may be asked to calculate simple things based on the 2×2 tables, which are the key to the results section of both therapeutic and diagnostic papers. For questions based on 2×2 tables, you will get the marks if you understand what you are calculating and can do it however the table is presented. You should start by stating the definitions and then always show your workings for your calculations.

Two things that will help you score well in these questions are:

1. Know your definitions, e.g. number needed to treat, sensitivity, etc.
2. Be able to draw a standard 2×2 table no matter how the data is presented in the question or paper.

QUESTION 8: YOUR INTERPRETATION OF THE PAPER AND/OR LIMITATIONS OF THE PAPER (4–8 MARKS)

These should be ‘easy-ish’ marks to get. There is often no definitive right or wrong answer and marks can be awarded for coherent answers, whatever your interpretation of the paper.

The question will be put into a context of how they want you to answer the question. For example instead of saying ‘Is this paper any good?’, the question will be something like ‘Your new chief executive heard about this paper and has asked if the ED will be introducing the intervention. Explain what you would do if you were clinical lead and give your reasons.’

The best way to approach this question is to split your answer to cover different areas where marks can be gained:

1. What are the results? Are the differences in outcome significant or due to chance? Do not get confused by the difference between statistical significance ($p < 0.05$) and clinical significance (will it affect my patients?).
2. Is the methodology good enough for you to believe the results and conclusion – internal validity.
3. Decide if it could be generalized to your patients – external validity.
4. Decide if it is important for your patients. How common is the condition, what is the clinical significance of the result?
5. What could prevent implementation of the intervention in the ED and how easy would the identified barriers be to overcome? Think about the four barriers to implementation: they all begin with B:
 - a. Badness: What harm could be caused by the intervention?
 - b. Burden of disease: Is it so rare that no one cares?
 - c. Budget: Can we afford to make the changes?
 - d. Beliefs: How will we educate people about this new intervention and how will we deal with intransigent colleagues?

TAKE HOME MESSAGE

1. At the beginning work out what type of study it is.
2. Read the questions before you read the paper.
3. You will be asked to write a summary. Write nothing but a summary, use the author's words and write it as you go along.
4. Answer the easy questions first – first and last questions.
5. When considering your opinion of the paper, consider internal and external validity.

10 The Critical Topic Review (CTR)

The Critical Topic Review (CTR) part of the FCEM examination is designed to test your ability to ask a clinical question, search the evidence, critically appraise the papers you find and synthesize the answer to your question. It formalizes something we should be doing regularly as clinicians on the shop floor of the ED and assesses our ability to do this. For the betterment of our patients, being able to critically assess the literature and manage patients in an evidence-based manner is essential.

CLINICAL QUESTION

A clinical question can often be broken into three parts:

- Population (who, where, when, what's wrong with them)
- Intervention or test
- Outcome measure(s)

For your CTR, the question should be designed to interrogate the evidence-based literature and applicable to Emergency Medicine practice.

Let's look at an example:

In an adult patient with a sprained ankle is splinting equivalent to plaster in respect to pain control and fracture healing?

Another way to look at this is to divide the question into 'PICO':

- Population: Adult patients self presenting to ED
- Intervention: Splinting
- Control: Plaster
- Outcome measures: Pain control and fracture healing

If you start with a root question, for example 'In a patient with a sprained ankle is splinting equivalent to plaster?', you can then use qualifiers such as adult or paediatric patients, only studies performed in EDs and so on to narrow the number of papers your search will find. Equally you can refine the outcome measures you look for to select out a specific subset of the literature.

This process allows you to generate a manageable number of papers to critically appraise for your CTR – we suggest aiming for four to seven. Remember the CTR

has a word count of 3500 words and the result of your search needs to allow you to produce a CTR that sticks to this word limit. Conversely, if on your first literature review you find that you do not have enough papers then you can widen your search terms to try to include more papers.

If you find the clinical question gives you more papers than you would like, don't throw it away yet because you may be able to qualify the number of papers in your search by using exclusion criteria such as only including RCTs or not including non-English papers. Every time you use an exclusion criteria you must be able to justify why; remember that ready availability of online translation tools has made excluding non-English language papers harder to justify.

For most of us the CTR is an evolutionary process where we wrote multiple questions, often on different topics, before settling on one that gave us a topic and set of papers we were happy to appraise in 3500 words (not including appendices or texts).

It is important to recognize that marks are given in both the written and viva for personal work and this also will have a bearing on the topic you choose. Perhaps our strongest advice is to choose a topic that you are genuinely interested in because you should be spending a lot of time on it.

LITERATURE REVIEW

Once you have a clinical question you can then perform a literature review. This can be done using the National Health Service (NHS) portal website or through your local medical library. Ideally you should use more than one search engine. Examples include:

MEDLINE/PubMed: General medicine
 AMED: Allied health professionals and complementary medicine
 BNI: British Nursing index
 CINAHL: Nursing and allied health database
 EMBASE: General medicine and pharmacology
 PsycINFO: Psychiatry

Write your question using parentheses [()], quotation marks ("") and AND OR NOT (Table 10.1). Watch for spelling changes between American English and British English (for example z versus s); plural versus singular and spaces (soft-cast versus soft cast). Simple things like this can mean that you exclude potential papers and write a less impressive CTR.

Let's consider a literature search of all the relevant databases online in respect to immobilization of ankle sprains:

- Ankle sprain AND immobilization: Gives 41 results.
- Ankle AND (sprain OR strain) AND (immobilization OR immobilization): Gives 112 results.

TABLE 10.1
Terms to use to help us refine literature searches

Terms (Boolean operators)	Meaning
AND	Retrieves results that include all the search terms.
OR	Retrieves results that include at least one of the search terms.
NOT	Excludes the retrieval of terms from your search.
*	Includes all terms with the root before; e.g. fractur* will retrieve all the results for fracture, fractures, fractured etc.
()	Most search engines work through a phrase from left to right, i.e. when we search for acute MI AND smoking OR coffee drinking, the search will return all papers that contain acute MI and smoking, plus all papers that contain coffee drinking. This is not what we intended! We can force the order of search engines' process using brackets around 'smoking OR coffee drinking' to get results of papers about acute MI that also refer to smoking or coffee drinking.
"**"	Allows you to search for a specific phrase.

Now on review of some of the titles, you notice that there are a lot of abstracts that refer to fractures, which you do not want to include. So you can change your search strategy:

- Ankle AND (sprain OR strain) AND (immobilisation OR immobilization) NOT fracture*: Gives 83 results.

Removing duplicates leaves you with 50 abstracts to review, a much more manageable number. Reviewing the abstracts will reduce the number of papers you have to read and critically appraise.

Remember if there were too many papers for your CTR after going through the abstracts, you can qualify your clinical question, for example only look at grade III ankle sprains, use certain types of papers (RCTs) or a certain grade of evidence.

It is important that you don't miss any other easy to access evidence such as Cochrane reviews, BestBETS or NICE (National Institute for Health and Care Excellence) guidance. It is also worth checking for unpublished research – the national research register (<https://portal.nihr.ac.uk/>), the clinical trial website in the United States (<https://clinicaltrials.gov/>) and EudraCT for Europe (<https://eudract.ema.europa.eu/>). If you find any such research you can then contact the authors and see at what stage of the process they are; this will impress the examiners and gain you extra marks.

Trying to access grey literature may also give you more potential papers. Perform a Google/Google Scholar search for your topic and make sure all the results it finds are already in your search results.

After reviewing the abstracts and titles, you will end up with a shortlist of studies that you need to get the full papers for. Your local library can also help with this. Make sure to give yourself time to get and review the papers. Make sure to check the bibliography of each paper you are going to include for further papers that you may not have found from your literature review.

CRITICALLY APPRAISING THE PAPERS

Once you have gathered your papers, you need to decide if the paper pertains to your clinical question. Sometimes despite reading the title and the abstract, the paper may not actually answer the question that you are asking and as such should be excluded. Once you are happy that all the papers are relevant to your clinical question you can begin to appraise them.

Consider using a checklist (see Appendix A). In the critical appraisal part of the FCEM examination, this is obviously not possible, but for the written part of your CTR, this will allow reproducibility of assessment and easier identification of positive and negative factors of the study design. It is worth keeping a copy of the completed checklists to take with you if you have a CTR viva.

As part of your assessment use level of evidence or scoring systems, such as Jadad scores. This highlights that you have spent time appraising the papers. It also may allow you to exclude papers based on level of evidence if you are finding that you have too much information for 3500 words. Finally, it will make it easier to explain the strength of the evidence, when you are making conclusions about each of the papers.

Consider performing a basic meta-analysis if the papers are amenable to this. This would provide an excellent method of garnering marks for your personal work and can be done relatively easily using software from the Internet.

Once you have gathered and critically appraised your evidence, you need to come to some conclusions. Is the evidence robust enough to say yes or no to your clinical question? Is it applicable to practice and how would it change your practice? If it does, your personal work could be performing an audit on current practice within your department, introducing an evidence-based guideline and re-auditing to complete an audit loop. If there is not enough evidence to come off the fence, how do we go about designing further research to answer the question?

PERSONAL WORK

There are marks for personal work for both the viva and the written portions of the exam. These marks can be significant and in some cases can be the difference between passing and not passing. Most people write the CTR and then think about the personal work they are going to perform, but really it should be the other way around. You should be thinking of topics for a year to 18 months prior to the CTR exam so you can perform adequate personal work. A small survey is not counted as adequate and will not get you many, if any, marks. Think about introducing a new practice, writing a meta-analysis, doing some new research or completing an audit loop.

WRITING THE CTR

The College produces excellent guidance on how to write up the CTR. The College CTR guidance is found within the FCEM exam regulations at <https://secure.collemergencymed.ac.uk/Training-Exams/Exams/FCEM>.

The front sheet should include the title, candidate name, word count (not including tables or appendices) and a declaration that it is your own work. It must be written in size 12 font with double spacing and the margins should be 2 cm. References should be in Vancouver style. Despite the College of Emergency Medicine examination committee being explicit of their expectations for the CTR, people lose marks for not following their guidance. Despite the fact that your word count does not include appendices, don't be tempted to use appendices to provide all the evidence of your critical appraisal of the literature as this is frowned upon.

In your introduction to your CTR, you should explain why you have chosen the topic and why from an emergency medicine perspective it is important. You should clearly state your clinical question. You need to explain your literature search – what your search was, what databases you used, any other methods of gathering information used, how many results obtained from your search and how many papers you used. It is important to explain why you excluded papers from your finished CTR. A good way to explain much of this information is to use a flow diagram showing how papers were excluded. You need to explain your critical appraisal of the literature. You must also succinctly summarize the evidence and explain your conclusions. As part of this, explain your personal work and how it is relevant to the topic. References and appendices should follow your conclusions.

Make sure that at least one person reads your CTR prior to submission. Consider asking someone non-medicine based to read it as well; they often notice formatting abnormalities that we might otherwise miss. Having a librarian or a colleague check your literature search is also a good idea both prior to submission and prior to the viva.

MARKING

If you score less than 7/28 you have automatically failed and will have to repeat the examination. If you score more than 21/28, you have automatically passed. Scores in between necessitate attending the viva. You receive marks for your title, clinical question and the relevance of your clinical question to emergency medicine. You get marks for the formatting and your literature review, which can be maximized with a little work. There are also marks for how well you critically appraised your papers, your conclusions and recommendations. And of course there are marks for personal work.

The viva part of the examination is 15 minutes duration. The Consultants examining you will have repeated your literature search before your viva and so should you. Sometimes you will find a new paper has been released since you wrote your CTR and sometimes you will find a paper you missed. Don't despair if that is the case. Bring any new papers with you and think about how they alter the conclusions you have made. Make sure you read your CTR before the viva several times. You should know your CTR better than your examiner but as with anything, some of the detail will have been lost in the period since submission. Questions and marks will come in five areas: relevance, your literature search, critical appraisal, conclusions and your personal work. Think about all these areas before the examination and consider having a colleague read your CTR and ask you questions on it. It often hones your thoughts about your CTR and helps you recognize areas you could improve on, which might be asked about during the viva.

TAKE HOME MESSAGE

1. Start your CTR early. You need time to do personal work.
2. Choose a topic that interests you.
3. Refine your clinical question; broaden or narrow the search terms to get four to seven papers that you can critically appraise.
4. Repeat your literature search prior to the viva.

Appendix A: Checklists

TABLE A1.1
Therapeutic paper checklist: what to look for in the paper

General questions	Specific questions	Answer
What sort of paper is this?	Was this looking at the results of an intervention? Was there a control group and if so were patients randomized into each group? Was the clinical question and primary hypothesis clearly stated (usually at the end of the introduction)?	
What is the internal validity?	How were the patients randomized? Was it clearly described how this was done? Was adequate allocation concealment used? Were subjects and investigators blinded to the treatment allocation? Were treatment and control groups similar at the end of the trial (usually table 1 – baseline characteristics)? During the trial, is the treatment under investigation the only difference between groups? Was a power calculation performed? Are all the patients accounted for? If not is there an attempt to follow those lost to follow-up? Were there many protocol violations? If so was an intention to treat analysis performed?	
What were the results?	Did the paper have a CONSORT diagram? The paper should report (or the results given should allow you to draw a contingency table and calculate): <ul style="list-style-type: none">• Control and experimental control rate• Relative risk, relative risk reduction and absolute risk reduction• Number needed to treat (or harm)	
Applicability (external validity)	Did they quote CIs and p values? Were they significant? Are the patients in the study similar to your patients? Were all relevant outcomes including likely side effects quoted? Is the treatment likely to be acceptable to your patients? Do the number of dropouts from the study support this? If a difference has been demonstrated is this a clinically significant difference as well as a statistically significant difference? Is the cost of the treatment likely to be worth the clinical difference demonstrated?	

TABLE A1.2**Diagnostic paper checklist: what to look for in the paper**

General questions	Specific questions	Answer
What sort of paper is this?	Is the paper really a study of diagnostic accuracy? Does it compare a specific diagnostic test against a gold standard? Is it a validation study looking at a previous study?	
What is the internal validity?	<p>If a convenience sample was taken is this really representative of the population you are interested in?</p> <p>Do you believe the reference standard will make the diagnosis accurately?</p> <p>Was the time between testing with the reference and index (trial) test short enough that the target condition couldn't have changed?</p> <p>Did the whole sample receive both the index (trial) test and reference standard? Ensure that there weren't different reference standards in use for different parts of the trial.</p> <p>Does the test give a simple yes/no answer? If no, did they draw a ROC curve and explain how the cutoff was chosen? Is this clinically as well as statistically appropriate?</p> <p>Was it necessary to quote a Kappa value and if so was it greater than 0.8?</p> <p>Were those interpreting the reference test blinded to the index test result and vice versa?</p>	
What were the results?	<p>The paper should report (or the results given should allow you to draw a 2×2 table and calculate):</p> <ul style="list-style-type: none"> • Sensitivity and specificity • Positive and negative predictive values • Positive and negative likelihood ratios <p>Did they quote CIs and what were they?</p>	
Applicability (external validity)	<p>Is the population used similar to the people you would use the test on? Did they describe how they selected them?</p> <p>Did they come from the same population as yours – i.e. an ED with patients presenting 24 hours a day and with similar baseline characteristics?</p> <p>Is the index test acceptable to your patients? Is it affordable and available in your setting?</p> <p>Do the results suggest it is accurate enough to influence your practice?</p>	

TABLE A1.3
Meta-analysis paper checklists: what to look for in the paper

General questions	Specific questions	Answer
What sort of paper is this?	Is the paper a meta-analysis (includes statistical combination of results from a number of trials) or a systematic review (does not)?	
What is the internal validity?	Is there an appropriate and clearly focused question? Is there a description of the methodology used and have they produced a flow diagram to demonstrate paper selection? Does the review collect papers that you think are relevant to the stated question? Is the literature search well described? If needed, could you repeat the search based on the information given? Is the literature search rigorous enough to have identified all the relevant studies? Was the quality of the studies assessed and reported using a valid assessment tool, e.g. Jadad score? Are there enough similarities between the papers to make combining them reasonable? Was a formal assessment of the heterogeneity between the studies made?	
What were the results?	Was a funnel plot produced and was there any asymmetry? Have the authors produced a forest plot to demonstrate their results? Is the overall result statistically significant?	
Applicability (external validity)	Did the study populations in the trials meta-analyzed reflect the population you treat? Does the result represent a clinically important difference?	

Appendix B: Glossary

This list is not exhaustive but is representative of the things examined in the FCEM including terms that will be of use when doing the CTR.

Absolute risk reduction: See *Risk*.

Allocation concealment: The process whereby patients and researchers are unaware of which arm trial subjects will end up in once they have been enrolled in the trial. The important point is that if researchers or participants know what arm of a trial they are likely to go into, they may consciously or unconsciously change their decision about whether to enrol. Allocation concealment is the process where this is hidden until after enrolment and so reduces selection bias in a trial. See Chapter 6.

Alpha (α) value: α is the probability of making a type 1 error (a false positive – saying we reject the null hypothesis when in fact we shouldn't). α is used in power calculations and is called the 'threshold of significance'. When considering a single test, if we set α at 0.05, we're saying that we're willing to accept a 5% chance of a false positive and that if the p-value we calculate comes under 0.05 we'll accept the result as statistically significant.

Bayes' theorem: It's argued (mostly by Rob) that this is the most important concept in medicine. There are lots of complicated explanations of the theorem in other books and on the Internet but a simple way to explain it is to use an example. Essentially, the idea is that reaching a decision about how probable a diagnosis is relies on two stages. Assess the pre-test probability of someone having a diagnosis (e.g. your *clinical suspicion* whether a patient with chest pain is having an MI) and then do a test to try to make the diagnosis. Any given test (e.g. in this case a troponin) will provide a likelihood ratio for the diagnosis – but this is a feature of the test and not the patient. Bayes' theorem states that the probability of the diagnosis is proportional to the product of the pre-test probability (in this case a feature of your history and diagnosis) and the likelihood ratio of the test. If your pre-test probability was low (e.g. a 19-year-old patient with muscular sounding pain after lifting weights) even a raised troponin would make the diagnosis of MI unlikely; if your pre-test probability was high (a 75-year-old arteriopath who smokes 80 cigarettes a day and has central crushing chest pain) a low troponin would not be sufficient to convince you he hadn't had an MI.

In the case of troponin, it is such a good test in our population that we actually usually rely quite heavily on it. In the case of DVT however, our d-dimer test is not so good, and the process of working out the pre-test probability has been formalized in scores such as Wells criteria.

Bias: A systematic error in the way a trial is designed or run leading to inaccuracy in the result. There are many types of bias. When appraising a trial, one of the most important judgments to make is whether you think there was bias that might have affected the result. We have detailed various sorts of bias in Chapters 5 and 6 but as recommended by the Cochrane Collaboration, when considering bias it's worth breaking the areas that might have been affected by bias into the following steps:

1. Selection of participants: Was selection truly random and was allocation concealed before enrolment?
2. Performance bias: Were there any differences in the interventions (other than the treatment under investigation) given to the different trial groups that might have affected outcomes? Were all those involved in the trial adequately blinded to ensure that the different trial interventions were not known by those taking part in the trial or analysing the results?
3. Detection bias: Were those assessing participant outcome adequately blinded to the intervention that the participants received?
4. Attrition bias: Do the authors adequately explain the reasons behind patients dropping out or being excluded from the study? Do their explanations seem reasonable?
5. Reporting bias: Question how the authors decided upon what result they have chosen to report and whether it follows from the trial they originally designed.
6. Anything else: The above five areas are not exclusive, so also consider if there is anything else about the trial that would make you worried that the result had been influenced by bias.

Case control study: People with a disease are compared to those without the disease to find out what risk factors both groups have been exposed to in the past. It is therefore always retrospective. Pros: quick, cheap, need few participants, good for studying rare diseases or those where there is a long duration between exposure and outcome. Cons: rely on recall and records to assess exposure, not good for looking at rare exposures, temporal relationships between exposure and outcome may be difficult to establish.

Causation and correlation: Correlation assesses the relationship between two variables, but it is an old adage of science that correlation does not prove causation.

In 1965 Sir Austin Bradford Hill presented a list of criteria by which we might judge a relationship to have causation rather than simply correlation. These have very much borne the test of time and are still used today:

1. Strength: The larger the association the more likely causation exists.
2. Consistency: Repeated observations that confirm the link in different places, at different times, by different people.
3. Specificity: If the link is limited to specific situations, peoples and disease, then causation is more likely.

4. **Temporality:** Does the exposure to the variable always precede the development of the disease? This is particularly important in diseases where there may be a long time between development of the disease and diagnosis.
5. **Biological gradient:** Does a greater exposure to the variable (e.g. smoking more cigarettes) make one more likely to develop the disease?
6. **Plausibility:** Is there a biological rationale that is reasonable for why exposure to the variable might cause the disease?
7. **Coherence:** Does the evidence run in line with the other things known about the variable and the disease so far?
8. **Experiment:** If we intervene to reduce the exposure to the variable are we able to demonstrate a reduction in the disease or outcome?
9. **Analogy:** If the potential causation mirrors that seen in a similar exposure we might be willing to accept lesser other evidence of causality.

Chance: A random variation that leads to imprecision in a study result.

Cohort study: Groups with and without an exposure are followed up to see how the development of an outcome differs between the groups. Pros: good for resolving questions about aetiology, harm and prognosis, rare exposures, assessing temporal relationships and many outcomes, can give direct estimation of disease incidence rates. Cons: not good for rare outcomes, can be expensive (e.g. long time from exposure to outcome), open to bias if people drop out. Cohort studies can be performed prospectively or retrospectively.

Confidence interval (CI): Gives a range of values that is likely to cover the true but unknown population value. CIs are more valuable than the p-value (primarily because of the difference between statistical significance and clinical relevance). CIs allow the reader to assess the likely clinical significance of where the result might lie. A CI can be given for various parameters: means, difference between means, ratios (RR, OR etc.). See Chapter 3.

Confounding factor: A factor or variable within a study that is not part of the intervention and is not controlled for but which differs between the two groups being studied and which will affect the outcome of interest. Randomization is the tool we use to make sure that confounding factors (even those we don't know about) are equally spread between groups. See Chapter 6.

CONSORT diagram: A flowchart showing the flow of participants through each stage of a randomized trial.

CONSORT statement: The Consolidated Standards of Reporting Trials statement is a set of recommendations and tools intended to improve the way in which randomized controlled trials are reported and thus make them easier for readers to understand. It gives a suggested template depicting the passage of participants through a trial and shows what should appear in the various parts of the write-up (title, abstract, introduction, methods, results, discussion and other information).

Contingency table: A contingency table is a table used to show the results of diagnostic or therapeutic trials. In studies with two arms and two outcomes this is a 2×2 table and what we traditionally are used to working with.

In therapeutic trials, along the horizontal axis is the outcome and along the vertical axis, the intervention. In diagnostic studies the gold standard is along the horizontal axis and the study under investigation along the vertical axis.

Controlled trial: People in the study are given one of two (or more) treatments or interventions.

Convenience sample: A type of non-probability sampling drawn from the part of the population that is close at hand, i.e. a sample population is selected because it is readily available and convenient.

Crossover trial: All the people in the study receive one treatment and then halfway through the study cross over and receive the other treatment. Pros: useful in studying rare diseases as subjects are their own controls. Cons: difficulties with order effects, historical controls and carryover effects (e.g. if you give someone an analgesic then they cross over to the placebo, the painkilling effect of the first drug may still be present. This can be avoided by using a sufficiently long ‘wash out’ period between cross over).

Cross-sectional study: Disease status and exposure status are measured simultaneously in a given population at a set time point. Pros: good for establishing prevalence or establishing an association. Cons: not good for establishing causation, recall bias is a problem, you need large groups, groups may be unequal sizes with unequal distribution of confounders.

Data: Data can be described in various ways, as shown in Figure A2.1.

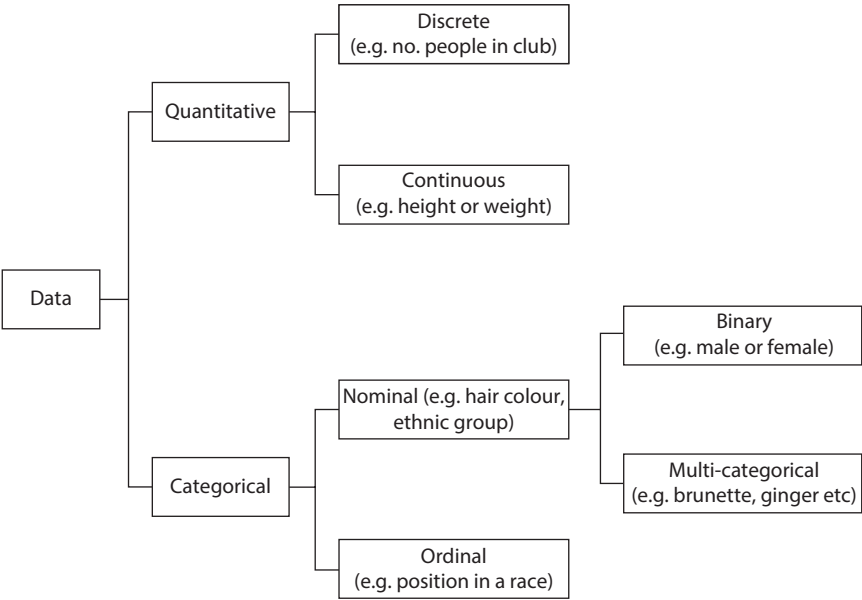


FIGURE A2.1 Types of data.

Effectiveness: Does the intervention have the desired effect under ordinary circumstances? If external validity is present and the trial demonstrates an effect of the intervention, we can accept that a treatment has effectiveness.

Efficacy: Describes the impact of the intervention under optimal (trial) conditions. If internal validity is present and the trial demonstrates an effect, we can accept that a treatment has efficacy.

End points: The end point is the result that you set out to look for. The main end point that the researchers are looking for (and the end point that the trial is powered to look at) is termed the primary end point. It should be clearly stated at the start of the trial and the result of the primary end point should be easy to find in the results section. Other end points that are collected but that are not the primary end point are termed secondary end points and it must be remembered that the trial has usually not been powered to look at them. A hefty degree of caution should be used when trying to draw conclusions from secondary end points. A further common way to divide end points is:

1. Clinical end point: An actual diagnosis of disease, symptom or sign. Death would perhaps be the best and most undeniable clinical end point!
2. Surrogate end point: A measure that may correlate with a real end point but is not necessarily guaranteed to. An example would be to measure gastric pH in an antacid trial rather than confirmed upper GI bleeding.

Error: (a) Type 1 error (α) = False positive. Falsely rejecting the null hypothesis when it is true. (b) Type 2 error (β) = False negative. Falsely accepting the null hypothesis when it is wrong.

External validity: The extent to which the results of a study can be extrapolated to a wider population. It is closely related to the term *generalizability*.

False negative: This is best explained using an example: when using a test, diagnosing someone as not having a disease when in fact they do have it.

False positive: This is best explained using an example: when using a test, diagnosing someone with a disease when in fact they do not have it.

Forest plot: A diagram produced in meta-analysis that shows the treatment effect from each trial and the overall treatment effect that has been calculated from the meta-analysis of the data. See Chapter 8.

Generalizability: Is the finding of the study applicable elsewhere?

Gold standard test: The best test available to diagnose a condition. It should be used as the comparison test in diagnostic studies looking at the ability of a novel test to detect the same disease.

Hazard and hazard ratio: Hazard is a measure of how quickly an outcome occurs. If we're looking at death from MI in an aspirin versus no-aspirin trial, we can calculate not only the overall death rate in each group but also the rate of death and this is termed the hazard in each group. Although we like to keep things easy and we strongly recommend using risk and risk ratios when interpreting the results of an interventional trial, clearly there is value in an assessment that includes a measure of time to the outcome.

Hazard ratio is simply the hazard in the intervention group divided by the hazard in the control group. It is not the same as risk ratio and tells us how many participants in the intervention group will have the outcome compared to the control group at the next time point. It assumes that the hazard ratio remains constant between the two groups over time and should always be quoted alongside a time frame.

Incidence: The rate of occurrence of new cases over a given time, in a defined population.

$$\text{Incidence} = \text{Number of new cases over a time period} / \text{Population size}$$

Intention to treat analysis: All the study participants are included in the analysis as part of the group to which they are allocated, regardless of whether they completed the study, or received the treatment to which they were allocated. Intention to treat generally signifies a higher quality study methodology and acts to reduce attrition bias. See Chapter 6.

Inter-quartile range (IQR): The IQR is a measure of spread or dispersion of a data set. It is the difference between the upper (75th) and lower (25th) quartiles and thus contains the middle 50% of the sample. It is smaller than the range and thus less affected by outliers. IQR (along with median) is used to describe data instead of standard deviation (and mean) when the data is not normally distributed, ordinal or from a small sample.

NB: Quartiles are values that divide sample data into four equal numbers of observations.

Internal validity: The extent to which the methodology of the paper permits a reflection of the truth, i.e. was the study methodologically good enough to demonstrate what was actually going on?

Jadad score: One example of a scoring system to assess the methodological quality of a trial – usually an RCT.

Three questions:

1. Was the study described as randomized? (1 point)
2. Was the study described as double blind? (1 point)
3. Was there a description of withdrawals and dropouts? (1 point)

Two additional points if:

- a. The method of randomization was described in the paper and was appropriate. (1 point)
- b. The method of blinding was described in the paper and was appropriate. (1 point)

Points are deducted if:

- a. The randomization was described but inappropriate.
- b. The blinding was described but was inappropriate.

Kaplan-Meier curves: A widely used form of graph used to compare survival over time between groups in an interventional trial. It uses a method that incorporates the facts that not all participants will reach the end of the trial (through drop out or death) or indeed have the outcome of interest during the time frame of the trial. Instead of simply calculating event rates

at various times, the event rate is adjusted for loss of participants (called ‘censoring’). The curves that are produced for the different groups in the trial can be compared to see if there is a difference in outcome over time. Although widely referred to as ‘survival curves’ they can in fact be used to look at any binary outcome measure (death, AMI, relapse, discharge from hospital, etc.).

Kappa value: The Kappa statistic indicates the level of agreement between measurements made by different raters and gives an indication as to whether this agreement is more than can be expected by chance. It is also used to assess the level of agreement when the same observer does repeated observations with a test. Kappa is used for measuring categorical data – for our purposes, generally a ‘disease present’ or ‘disease free’ answer. The Kappa value lies between 0 and 1.0 where 0 is chance agreement and 1.0 is perfect agreement.

Agreement	Chance	Poor	Moderate	Substantial	Perfect
Kappa	0	0.20	0.40–0.60	0.80	1.0

Likelihood ratio (LR): Useful when considering an individual’s test result in combination with your clinical findings. Also see *Two by two table*.

Mean: The sum of the values divided by the number of values. Normally distributed data is usually described using mean and standard deviation.

Median: The middle number of a group when they are ranked in order. The mean of the two middle values is taken if there is an even number of values. It is useful for skewed data and has the advantage of reducing the importance of outliers. Non-normally distributed (or non-Gaussian) data is usually described using median and inter-quartile range.

Minimization: A form of adaptive randomization (see Chapter 6). A running total is kept of all the levels of the prognostic factors of interest. The first patient in the trial is randomly allocated but all subsequent patients are randomized using randomization weighted towards which assignment would minimize the imbalance.

Mode: The most frequently occurring value in a list. It has the advantage in that it can be used with non-numerical data (e.g. colours of cars).

Negative predictive value (NPV): The proportion of people with a negative test who do not have the disease. See *Two-by-two table*.

Nominal data: A type of categorical data where the order of the categories is not important, e.g. country of birth.

Normal data: Data that is normally (i.e. bell shaped) distributed about the mean. Also called Gaussian data.

Null hypothesis: This is the statement that therapeutic trials start with. The investigators will have a hypothesis that, for example, a drug lowers blood pressure. However for the purpose of statistical inference, the researchers say that they are sampling from a population where the null hypothesis is true, i.e. we are sampling from a population in which the intervention has no

effect. The results of the trial are then used to either reject or not reject this hypothesis.

Number needed to treat (or harm): The number of people who need to be treated to achieve one outcome of interest. It is a powerful way to describe the value of an intervention. Mathematically it is $1 \div \text{Absolute risk reduction}$.

See *Risk*.

Odds: Odds are a way of describing risk. See Chapter 4.

Mathematical odds = Desired events \div Undesired events

Odds ratio = Odds in treated group \div Odds in control group

Open trial: Both the researchers and participants know what treatments are being given to individuals. It is the opposite of a blinded trial.

Ordinal data: Data that can be ranked (put in order) or have a rating scale attached to it. You can count and order, but not measure, ordinal data, e.g. first, second, and third in a race.

p-value: A formal definition of a p-value can be given as: A p-value is the probability of obtaining a result at least as extreme as the one observed, assuming the null hypothesis is true (i.e. assuming there is no effect in the population). The lower the p-value, the less likely the result is under the assumption of the null hypothesis (i.e. the less likely the result is to have occurred by chance).

For many people, that's a bit of a mouthful and so we prefer: $P < 0.05$ means that the probability of obtaining the result by chance was less than 1 in 20. By convention, a p-value of less than 0.05 is the accepted threshold for statistical significance and is the level at which the null hypothesis can be rejected.

Calculation of the p-value is beyond the scope of this book and the FCEM but a description of which statistical test to use in which scenario is given in Chapter 3. Also see *Statistical tests*.

Paired or unpaired data: Paired data is data for the same individuals at different time points; for our purposes usually before and after an intervention. Unpaired data is two groups having different members (also known as 'independent data') as when we compare an interventional and control group in an RCT.

Parametric statistics: Inferential methods that rely on the data following the normal distribution.

Per protocol analysis: Also called 'on-treatment analysis'. Only those who sufficiently complied with treatment (as per the trial protocol) are included in the statistical analysis. It is a good way to demonstrate the actual effect of an intervention on people who take it (or can tolerate it) as the researchers intended but it does not allow the reader to see how applicable the treatment is to the whole population.

Positive predictive value (PPV): The proportion of people with a positive test who have the disease. See *Two-by-two table*.

Power: The power of a study defines its ability to demonstrate an association or causal relationship between two variables, given that an association exists, i.e. when there is a relationship, the ability of the study to actually show it. A more formal definition would be: the power of a statistical test is the probability that the test will reject a false null hypothesis (i.e. it will not make a type 2 error).

If α has been set at 5% and the study powered to 80% (which is the lowest conventionally accepted level) then the study has an 80% chance of ending up with a $p\text{-value} \leq 0.05$ (i.e. a statistically significant result) if there really is a difference.

Power = $1 - \beta$, where β is probability of a false negative (type 2 error).

Four factors that affect the power of a study (and so should be mentioned in the methods of a paper!):

1. α value (set before starting – usually 0.05)
2. Sample size (therefore a power calculation should be done to calculate this before the trial starts)
3. Variability in data
4. Minimum clinically relevant difference/estimation of likely difference

Pragmatic trial: Pragmatic trials assess whether an intervention works under real-life conditions. It is simply concerned with whether an intervention works, not how or why.

All the people in a clinical setting (e.g. an outpatient clinic) are randomized to receive a particular treatment protocol. Pros: more reflective of clinical practice. Cons: more difficult to control, difficult to blind, difficulties with excessive drop out.

Prevalence: Point prevalence is the proportion of a defined population with a disease at any one time. Period prevalence is the proportion of a defined population with a disease during a given time period. Lifetime prevalence is the proportion of a population that has or has had a disease at a given time.

Randomization: The process whereby an individual has an equal chance of entering any group within a study. Randomization reduces bias by distributing characteristics of patients that may influence outcomes (confounders – known and unknown) randomly between the groups. There are various methods described in Chapter 6. See also *Minimization*.

Receiver operator characteristic (ROC) curve: Figure A2.2 is used to assess diagnostic tests whereby the cutoff point for a positive and negative test have been varied and then the specificity and sensitivity calculated for each cutoff. It is a tool that can be used to assess the best cutoff value for the test. The area under the curve (AOC) is a measure of how good the test is: 1.0 is perfect; 0.5 is no better than chance and the line would lie along the line of unity.

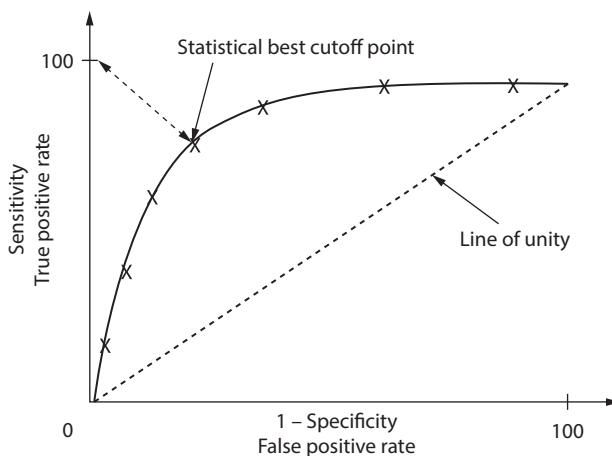


FIGURE A2.2 Receiver operation characteristic curve.

Statistically, the optimum cutoff point for a test (i.e. the balance between sensitivity and specificity) is where the line is closest to the top left corner of the graph however clinically we may choose a different point. See Chapter 7.

Regression analysis: Regression is a statistical method for looking at the relationship between one or more variables and/or outcomes. Consider a study measuring people's height and weight. If we were to plot the two on a graph, we might find that in general as height increases, so does weight. If the points lie approximately in a straight line, we can say there is a linear relationship between the two; if not there may be no relationship or there may be a non-linear relationship (e.g. the points lie in a curved shaped). The exact way this is done is beyond the scope of this book and FCEM but a general understanding is valuable.

Linear regression analysis is the statistical tool used to analyze the relationship between continuous (i.e. not categorical) data when there is a linear relationship between the two. If we are looking at one variable and one outcome we use *simple* linear regression; if we want to look at more than one variable we use *multivariate* linear regression.

Logistic regression analysis is used, as with multivariate linear regression, when we want to know the effect of multiple variables on categorical data (e.g. dead or alive). If necessary, it is possible to convert a continuous outcome (e.g. d-dimer result) into a categorical outcome (e.g. positive or negative) but this must be done in a clinically appropriate way. Unlike in linear regression, the relationship between the variables and the outcome does not need to be linear.

Risk (statistics from a therapeutic paper): This is covered in Chapter 6. In clinical research, Risk = Probability.

$$\text{Risk} = \frac{\text{Number of times an event occurs}}{\text{Total possible events}}$$

Using a 2 × 2 contingency table, we can calculate the following:

	Outcome +ve	Outcome -ve
Experimental group	a	b
Control group	c	d

Risk in the control group = Control event rate (CER)

$$= c/(c + d)$$

Risk in the experimental group = Experimental event rate (EER)

$$= a/(a + b)$$

From this you can then calculate:

Absolute risk reduction (ARR) is the difference in risk between the groups and is useful as it allows you to calculate number needed to treat.

$$\text{ARR} = \text{CER} - \text{EER}$$

Relative risk (RR) is a ratio of the risk in each group.

$$\text{RR} = \text{EER}/\text{CER}$$

Relative risk reduction = $(\text{CER} - \text{EER})/\text{CER} = \text{ARR}/\text{CER}$

Number needed to treat (harm) is the number of people who need to be given the treatment to generate one positive (or negative) outcome, i.e. = $1/\text{ARR}$.

Sensitivity: The ability of a test to rule out a diagnosis (SnOut). It is the proportion of people with the disease who test positive. Also see *Two-by-two table*.

Specificity: The ability of a test to rule in a diagnosis (SpIn). It is the proportion of people without the disease who test negative. Also see *Two-by-two table*.

Statistical tests: There are a number of different types of statistical tests, which can be used to analyze data and generate a p-value. See Chapter 3, Figure 3.1.

Study types: There are a number of different study types. These are shown in Figure A2.3.

Two-by-two table (statistics from a diagnostic study): This is used in diagnostic studies to compare the studied test to a gold standard.

		Gold Standard/Have the disease	
		Disease positive	Disease negative
Index	Positive	A	B
Test	Negative	C	D

		Gold Standard/Have the disease	
		Disease positive	Disease negative
Index	Positive	True positive	False positive
Test	Negative	False negative	True negative

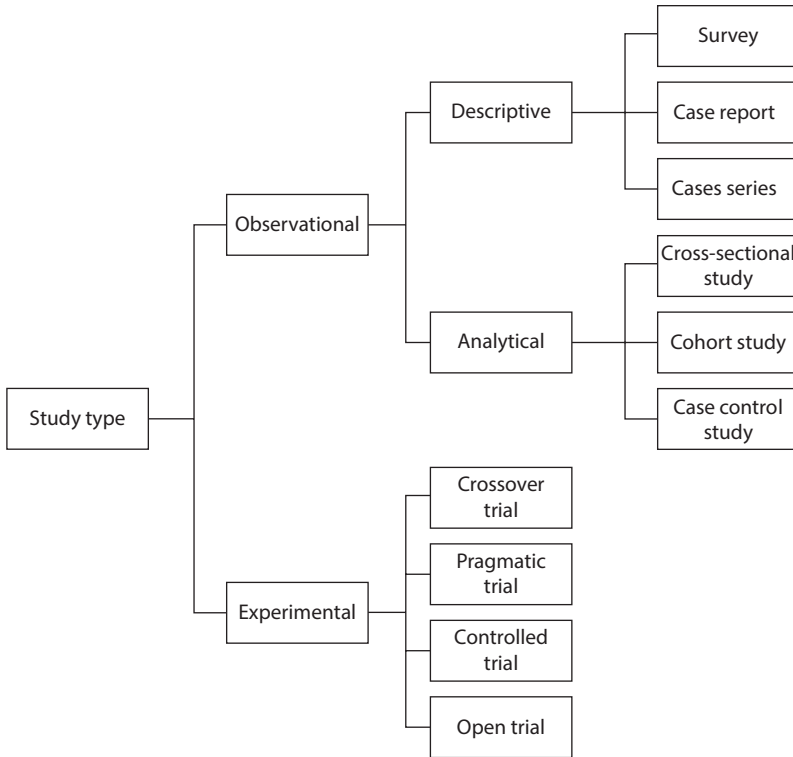


FIGURE A2.3 The different study types.

The table then can be used to calculate a number of values:

$$\text{Sensitivity} = a/(a + c) = \text{TP}/(\text{TP} + \text{FN})$$

= Proportion of people with the disease who test positive

= Sensitivity rules out (SnOut)

$$\text{Specificity} = d/(b + d) = \text{TN}/(\text{FP} + \text{TN})$$

= Proportion of people without the disease who test negative

= Specificity rules in (SpIn)

$$\text{Negative predictive value (NPV)} = d/(c + d) = \text{TN}/(\text{FN} + \text{TN})$$

= Proportion of people with a negative test who don't have the disease

$$\text{Positive predictive value (PPV)} = a/(a + b) = \text{TP}/(\text{TP} + \text{FP})$$

= Proportion of people with a positive test who have the disease

Positive likelihood ratio (LR+) = $\text{Sensitivity} / (1 - \text{Specificity})$
= Likelihood of someone with a positive result having the disease

Negative likelihood ratio (LR-) = $(1 - \text{Sensitivity}) / \text{Specificity}$
= Likelihood of someone with a negative result having the disease

Validity: Are the findings of the study true? See *Internal validity* and *External validity*.

Appendix C: Practice Papers

The only way to pass is by practising exams. We have provided two practice papers and we would suggest doing them in exam conditions. Please note that the papers we use here are fictional. The Critical Appraisal for FCEM course that we run has three mock exams during the two days and a minimum of a further four exams to take home (using real papers). For further information please go to: www.criticalappraisalforfcm.com

PRACTICE PAPER 1: JOURNAL ARTICLE

Provision of tea & coffee in the ED to reduce 'Left Without Being Seen' rates

Coughlan EP, Galloway RI, Bootland DJ

Royal Sussex County Hospital, Brighton

Introduction

In 2011, in response to concerns over the utility of the 4-hour rule, the UK Government introduced a new set of eight Quality Indicators (QIs) which it was claimed would act as a better measure of the care that National Health Service (NHS) emergency departments (EDs) provide. After further consultation, the QIs were altered such that EDs are now measured against five 'headline' QIs and these are divided into two main groups: Patient Impact (Unplanned re-attendance and Left without being seen) and Timeliness (Time to initial assessment, Time to treatment and Total time in the ED). Performance is currently judged against each group rather than each QI although Total time in the ED (essentially equivalent to the old 4-hour rule) remains as a standalone measure.

The measure of 'Left without being seen' (LWBS) 'reflects the satisfaction of patients with the initial management experience they receive'¹. Leaving before being seen by a clinical decision maker also represents clinical risk and it has been demonstrated that 49% of patients who are leaving before being seen subsequently require urgent treatment¹. The QIs suggest departments should not have LWBS rates at or above 5%. In our region, LWBS rates are currently running at a level of 6%².

Achievement of the standards against the QIs has proved challenging for many departments and has resulted in novel methods of improving quality. A recent telephone survey of patients from this department who had been designated as LWBS highlighted lack of comfort in the waiting area as a major cause of failure to wait to be seen³. We hypothesized that provision of free tea and coffee in the waiting room for both patients and relatives would reduce our LWBS rate and improve the service experience of patients attending the department.

Methods

Patients

The study population comprised all adult patients (age ≥ 18 years) presenting to one of four UK District General Hospital (DGH) EDs in the South of England over one week in August 2012. Patients admitted directly to a 'resus' cubicle, whether having been transported to hospital by ambulance or arriving by private transport, were excluded from the study.

Randomization

Due to the nature of the study, cluster randomization was deemed the most appropriate means of randomization. We invited four DGHs to be involved in the study and used computer generated random numbers to cluster randomize them to intervention or control group. Neither the patients nor the treating ED team were blinded to the group that their hospital had been randomized to. The local ethics committee waived the need for patients' consent.

Intervention

EDs randomized to the intervention group were provided with an automated tea and coffee making machine that was to be placed within the main seating area of the ED waiting room. Complimentary contracts with the owners of the vending machines were negotiated for the duration of the trial to ensure daily maintenance and stocking checks for each hospital in the intervention group. EDs in the control group were to leave their waiting rooms standard – none of the hospitals in the control group routinely provided free drinks for people in their waiting rooms.

Data collection

Patients' movement through the department and/or leaving before being seen was recorded as routinely done by each department. Demographic data, length of wait to be seen, pain score and final diagnosis on discharge from the ED were also recorded by each department's computerized patient system. On a weekly basis this data was retrospectively retrieved by four research nurses and entered into a computer database. A subset of the patients in each group was surveyed by a research nurse via telephone questionnaire to assess their satisfaction with their attendance in the ED. Patient satisfaction was measured using the Hartland questionnaire that has been validated in similar populations⁴; the Hartland score measures patient satisfaction on a scale of 0–10 with a change of 1 point considered to be significant.

Outcome measures

The primary outcome measure was rate of patients who left without being seen by a clinical decision maker – defined as a post-registration doctor or emergency nurse practitioner (ENP). Secondary outcome measures were patient satisfaction and waiting time.

Statistical analysis

Statistical analysis software that is routinely available commercially (SPSSTM) was used to analyze the data. The primary outcome measure was reported as a proportion of the total number of attendances and a test of independent proportions was performed with a chi-squared test. Secondary outcome measures were normally distributed and thus analyzed using Student's T-test.

Previous work by Parker et al.⁵ demonstrated a 3% absolute difference in LWBS after an intervention to improve comfort in the ED waiting area. We calculated that a total of 2400 patients would be required to demonstrate a reduction from 7% in the control arm to 4% in the intervention arm, with a power of 0.80 and a type 1 error rate of 0.05.

Results

During the study period, 2527 patients attended the four EDs enrolled in the study. Of these 91 patients were subsequently found to have incomplete data, predominantly due to having a discharge diagnosis and location recorded as 'unclassified'. As such 2436 patient attendances were analyzed; 1272 in the control group and 1164 in the intervention group (Figure A3.1). Characteristics of the study groups are given in Table A3.1 and it can be seen that the two groups did not differ in terms of gender or age.

The rates of patients classified as LWBS are shown in Table A3.2. The overall LWBS rate was statistically different between the two groups (Relative Risk 0.63, 95% CI 0.45–0.90, $p = 0.04$) with a LWBS rate in the control group of 6.37% and in the intervention group 4.04%. This was particularly true in the patients classified as 'minors' (RR 0.42, 95% CI 0.27–0.66, $p < 0.001$).

TABLE A3.1
Study group characteristics

Population characteristics	Total population	Interventional group	Control group	p-value
Sex, No. (%)				
Male	1368 (56)	660 (57)	708 (56)	0.227
Female	1068 (44)	504 (43)	564 (44)	
Age, y, mean (SD)	61.8 (1.4)	61.6 (1.3)	61.9 (1.4)	0.86
Mode of arrival, No. (%)				
Ambulance	504 (21)	192 (16)	312 (25)	<0.001
Private transport	1932 (79)	972 (84)	960 (75)	
Diagnosis group, No. (%)				
Medical	704 (29)	184 (16)	520 (41)	<0.001
Surgical	252 (10)	108 (9)	144 (11)	
Minor injury	1480 (61)	872 (75)	608 (48)	
Admission status, No. (%)				
Admitted	824 (34)	288 (25)	476 (37)	<0.001
Discharged	1612 (66)	876 (75)	796 (63)	

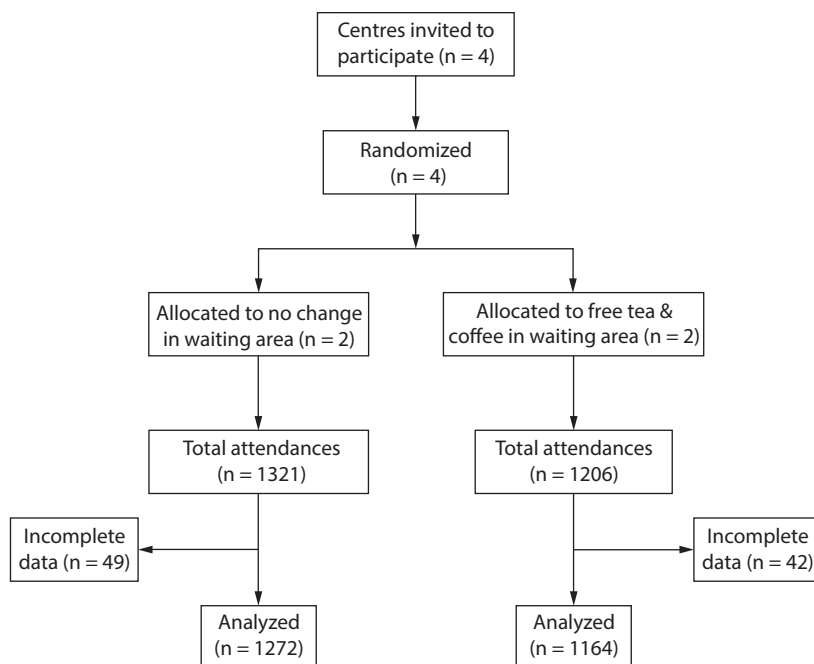


FIGURE A3.1 Patient flow diagram.

TABLE A3.2

Patients classified as left without being seen (LWBS)

Diagnostic Group	Total population		Intervention group		Control group		p-value
	Number	Percentage (95% CI)	Number	Percentage (95% CI)	Number	Percentage (95% CI)	
Medical	34	4.83	10	5.43	24	4.62	0.69
Surgical	17	6.75	8	7.41	9	6.25	0.80
Minors	77	5.20	29	3.33	48	7.89	<0.001
Overall	128	5.25	47	4.04	81	6.37	0.013

The time that patients waited to be seen by a doctor or nurse practitioner was slightly longer in the control group (79 minutes, 95% CI 78.1–80.5) than in the intervention group (75 minutes, 95% CI 72.7–76.3) and although small, and thus unlikely to be of importance, this difference was statistically significant ($p < 0.001$) (Table A3.3).

Patient satisfaction scores were obtained from a randomly selected subgroup of 200 patients (100 patients in each group, 20 medical, 20 surgical and 60 minors). There was a statistically significant improvement in patient satisfaction for those allocated to the intervention arm of the study ($p = 0.007$) and again this was most marked in the patients classified as ‘minors’ (Table A3.4).

TABLE A3.3
Waiting time to be seen by decision maker (minutes)

	Total population (mean)	Intervention group (mean)	Control group (mean)	p-value
Medical	83	93	79	<0.001
Surgical	89	97	83	<0.001
Minors	72	68	78	<0.001
Overall	77	75	79	<0.001

TABLE A3.4
Patient satisfaction – Hartland score

	Total population (mean)	Intervention group (mean)	Control group (mean)	p-value
Medical	8.3	8.2	8.5	0.32
Surgical	8.1	7.9	8.2	0.39
Minors	8.1	8.5	7.6	<0.001
Overall	8.0	8.3	7.9	0.007

Discussion

Our study showed that implementation of access to free tea and coffee in the waiting areas of UK DGH EDs significantly reduced the rates of patients leaving before being seen from 6.4% to 4.0% in the two groups, a reduction that would see hospitals achieve the NHS target of a LWBS rate of less than 5%. Although the intervention had its largest effect in those classified as having ‘minor’ injuries, it is important to note that the statistical significance of the result was maintained for all patients.

Although LWBS rates are perhaps intuitively clinically unimportant or certainly less important than clinical outcome measures such as pain scores and mortality, as has been highlighted by the UK Government, they can be seen to be reliable surrogates for the standard of care that EDs provide. Nevertheless, it is reassuring to note that not only were LWBS rates reduced by the intervention we studied but also that patient satisfaction scores in the intervention group were statistically higher than those in the control group. This suggests that this intervention has the rare ability to please clinicians, administrators and patients!

We believe this is the first prospective randomized controlled study to look at the impact of changes to improve comfort in the ED waiting area. Parker et al.⁵ previously performed a single centre pilot study with a retrospective control group looking at more general improvements in factors associated with patient comfort in the waiting area and demonstrated a 3% absolute reduction in LWBS rates between their before and after groups. The problems of retrospective and single centre studies are widely acknowledged but it is interesting to note that we found a similar absolute reduction in LWBS rates after implementation of a measure to improve the patient and relative experience.

The trial has a few limitations that are worthy of discussion. The lack of sufficient funding did not enable us to capture patient satisfaction scores at the time of their visit and thus these had to be done retrospectively by telephone. Equally, the size of the trial meant that it was not possible to follow up every patient and thus we had to be satisfied only with a small sample of the population studied rather than the whole population as would have been done in an ideal world. Finally, there is an inherent problem in offering food or drink to patients who may have needed surgical intervention early during their hospital visit. As such, it was left to the discretion of the individual ED to manage how best to advise this subset of patients attending the ED if free beverages were available. It was not possible either to identify these patients or to exclude them from the analysis but we felt this was justified as part of the pragmatic study design.

Conclusion

This multi-centre randomized controlled trial demonstrated a significant reduction in LWBS rates and improvements in patient satisfaction after implementation of free tea and coffee in the waiting area of UK EDs, particularly those classified as having 'minor injuries'. Such interventions are likely to not only improve the patient experience but to improve the ability of UK EDs to demonstrate their high standard of care through achievement of the LWBS QI and thus should be recommended to Trusts.

References

1. A&E Clinical Quality Indicators Implementation Guidance. Department of Health. Downloaded from http://www.dh.gov.uk/prod_consum_dh/groups/dh_digitalassets/@dh/@en/@ps/documents/digitalasset/dh_123055.pdf on May 3, 2012.
2. Galloway R, Coughlan E, Bootland D. An audit of Left Without Being Seen (LWBS) rates in six hospitals in the South of England. *J Crit Apprais* 2011;1:52–54.
3. Body D, Coughlan E. Telephone survey of those leaving the emergency department without being seen. *UK Med J* 2011;88:116–117.
4. Houghtan T, Smith J. The Hartland patient satisfaction scale – a validation study in a UK emergency department. *J Inv Med* 2006;26:14–18.
5. Parker S, Sadek S, O'Riordan S. Do improvements in comfort in the waiting area result in fewer patients leaving without being seen? A pilot study. *J Crit Apprais* 2011;2:102–105.

PRACTICE PAPER 1: QUESTIONS

Coughlan et al. Provision of tea & coffee in the ED to reduce 'Left Without Being Seen' rates

Question 1

Provide a no more than 200 word summary of this paper in the box provided. Only the first 200 words will be considered – short bullet points are acceptable. Maximum of seven marks.

Question 2

- i) What is cluster randomization? (1 mark)
- ii) Give one advantage and one disadvantage of cluster randomization. (2 marks)

Question 3

The following excerpt is taken from the text:

Previous work by Parker et al.⁵ demonstrated a 3% absolute difference in LWBS after an intervention to improve comfort in the ED waiting area. We calculated that a total of 2400 patients would be required to demonstrate a reduction from 7% in the control arm to 4% in the intervention arm, with a power of 0.80 and a type 1 error rate of 0.05.

What is meant by the power of a study? (1 mark)

Do you think the power of this study was appropriate? (1 mark)

What is meant by type 1 and type 2 error? Give an example of each. (2 marks)

Question 4

- i) What is the difference between a primary and secondary outcome measure? (1 mark)
- ii) What is the difference between qualitative and quantitative data? (1 mark)
- iii) In this study, are the secondary outcome measures quantitative or qualitative data? Explain your answer. (3 marks)

Question 5

- i) What does a ‘p-value’ represent? (2 marks)
- ii) What is meant by CONSORT statement? (2 marks)

Question 6

The following table was produced from a follow-up study in different hospitals but run with identical methodology. Calculate the absolute risk reduction, the relative risk, the relative risk reduction and the number needed to treat. (5 marks)

	LWBS	Waited to be seen
Free tea and coffee	330	610
No free tea and coffee	410	530

Question 7

Your chief executive has read this paper and tells you she wants you to use some of the ED budget to provide free tea and coffee in the waiting room.

What will you say to her? (6 marks)

PRACTICE PAPER 1: ANSWERS

Coughlan et al. Provision of tea & coffee in the ED to reduce 'Left Without Being Seen' rates

Question 1

Provide a no more than 200 word summary of this paper in the box provided. Only the first 200 words will be considered – short bullet points are acceptable. Maximum of seven marks.

Objectives: To see if provision of free tea and coffee in the waiting room reduces LWBS rate and improves the service experience of patients attending the ED

Design: A therapeutic, randomized controlled trial

Setting: Four district general hospital EDs in the South of England during one week in August 2012

Population: All adult (≥ 18 years) patients presenting to the ED; patients admitted to the resuscitation area were excluded

Methods: Patients were randomized by cluster with receiving hospital being the unit of randomization. Neither patients nor treating clinicians were blinded to the interventions.

Intervention/Treatment: Automated free tea and coffee machine available in the waiting room

Control: Standard waiting room

Outcome measures: Primary outcome measure was the rate of patients who left without having been seen by a clinical decision maker

Results: The overall LWBS rate was statistically different between the two groups (Relative Risk 0.63, 95% CI 0.45–0.90, $p = 0.04$) with a LWBS rate in the control group of 6.37% and in the intervention group 4.04%.

Conclusions: This trial demonstrated a significant reduction in LWBS rates and improvements in patient satisfaction after implementation of free tea and coffee in the waiting area.

Question 2

- i) What is cluster randomization? (1 mark)
- ii) Give one advantage and one disadvantage of cluster randomization. (2 marks)

- i) Cluster randomization is when the unit of randomization is not an individual study participant but rather a group of individuals treated in the same place or by the same person

ii)

Advantages:

Cluster randomization allows us to study interventions, which are applied to groups rather than individuals such as process changes such as the introduction of a Clinical Decision Unit (CDU). If we were to try having two different processes operating in the same clinical area (e.g. some patients having access to a CDU and others not), there is the potential for confusion, protocol violation and bias; cluster randomization works to avoid this.

Disadvantages:

Cluster randomization tends to require a greater complexity of design and analysis

Question 3

The following excerpt is taken from the text:

Previous work by Parker et al.⁵ demonstrated a 3% absolute difference in LWBS after an intervention to improve comfort in the ED waiting area. We calculated that a total of 2400 patients would be required to demonstrate a reduction from 7% in the control arm to 4% in the intervention arm, with a power of 0.80 and a type 1 error rate of 0.05.

What is meant by the power of a study? (1 mark)

Do you think the power of this study was appropriate? (1 mark)

What is meant by type 1 and type 2 error? Give an example of each. (2 marks)

The power of a study defines its ability to demonstrate an association or causal relationship between two variables, given that an association exists i.e. when there is a relationship, the ability of the study to actually show it. A more formal definition would be: the power of a statistical test is the probability that the test will reject a false null hypothesis (i.e. it will not make a type 2 error).

Yes. 80% power is the lowest conventionally accepted level, giving an 80% chance of ending up with a p-value ≤ 0.05 (i.e. a statistically significant result) if there really is a difference.

In this paper, the null hypothesis is that providing free tea and coffee will make no difference to LWBS rate.

Type 1 Error (α) = False positive. Falsely rejecting the null hypothesis when it is true; i.e. getting a result that makes us think that free tea and coffee does make a difference when in reality it doesn't.

Type 2 Error (β) = False negative. Falsely accepting (technically *not rejecting*) the null hypothesis when it is wrong; i.e. getting a result that makes us think tea and coffee doesn't make a difference when in reality it does.

Question 4

- i) What is the difference between a primary and secondary outcome measure? (1 mark)
- ii) What is the difference between qualitative and quantitative data? (1 mark)
- iii) In this study, are the secondary outcome measures quantitative or qualitative data? Explain your answer. (3 marks)

- i) The end point is the result that you set out to look for. The main end point that the researchers are looking for (and the end point that the trial is powered to look at) is termed the primary end point. In this study this would be LWBS rate. Other end points that are collected but that are not the primary end point are termed secondary end points, in this study this would be time to treatment and satisfaction scores.
- ii) Quantitative data is where a numerical value can be given to the data and this can be either discrete (i.e. with the number of people in a club it can be only 3 or 4 not 3.5) or continuous where it can be anywhere on a scale. Qualitative data (also called categorical data) is where there are a fixed number of defined categories into which data is organized and this can either be nominal (i.e. given a name) such as hair colour or ordinal (i.e. where there is a rank to the categories) such as first, second and third places in a race.
- iii) In the paper there are two secondary outcome measures. The waiting times are a quantitative result as they are a numerical value (time), the patient satisfaction results are also quantitative as patients have been asked to rate their satisfaction on a scale (giving a numerical result for satisfaction).

Question 5

- i) What does a 'p-value' represent? (2 marks)
- ii) What is meant by CONSORT statement? (2 marks)

- i) A p-value is the probability of obtaining a result at least as extreme as the one observed, assuming the null hypothesis is true (i.e. assuming there is no effect in the population). The lower the p-value, the less likely the result is under the assumption of the null hypothesis (i.e. the less likely the result is to have occurred by chance).

$P < 0.05$ means that the probability of obtaining the result by chance was less than 1 in 20. By convention, a p-value of less than 0.05 is the accepted threshold for statistical significance and is the level at which the null hypothesis can be rejected.

In this paper the primary outcome had a p-value of 0.04, meaning it was a statistically significant result.

- ii) A CONSORT statement is an evidence based minimum set of recommendations for reporting randomized controlled trials. It offers a standard way for authors to prepare reports of trial findings, facilitating their complete and transparent reporting, and aiding their critical appraisal and interpretation.

Question 6

The following table was produced from a follow-up study in different hospitals but run with identical methodology. Calculate the absolute risk reduction, the relative risk, the relative risk reduction and the number needed to treat. (5 marks)

	LWBS	Waited to be seen
Free tea and coffee	330	610
No free tea and coffee	410	530

Control event ratio (CER) = $410/(410 + 530) = 410/940 = 0.44$

Experimental event ratio (EER) = $330/(330 + 610) = 33/94 = 0.35$

Absolute risk reduction (ARR) = $0.44 - 0.35 = 0.09$

Relative risk (RR) = $0.35/0.44 = 0.80$

Relative risk reduction (RRR) = $0.09/0.44 = 0.20$

Number needed to treat (NNT) = $1/ARR = 1/0.09 = 11$

Question 7

Your chief executive has read this paper and tells you she wants you to use some of the ED budget to provide free tea and coffee in the waiting room. What will you say to her? (6 marks)

I would recommend to the chief executive that we should not be making this change at this point.

(Not all of the below would be required for all the marks!)

I have concerns about some aspects of the trial methodology and hence the internal validity of the trial (i.e. how well it was designed and run).

- This was an unblinded study and although given the type of trial it is not unreasonable that the patients and clinicians were not blinded we might have expected those collecting the data to be blinded. This leads to an increased risk of bias.

- Although the paper states that free tea and coffee were not available in the control group, it would have been valuable to know what refreshments were available to people waiting.
- When looking at the results, there were statistically significantly more medical patients and fewer minors patients in the control group than the intervention group. This raises the possibility that there were important differences in the populations between the control and intervention group. This may be an unrecognized confounder and the difference in LWBS may be due to the differences in minors patients in the different groups, rather than tea and coffee availability.
- The satisfaction scores were recorded retrospectively and it is unclear how the subset who were contacted were selected suggesting the possibility of recall and reporter biases respectively.
- The authors describe a 'significant' Hartland score change as 1 point but their results show a difference of only 0.4 – highlighting the difference between statistical and clinical significance.

The external validity of the trial (i.e. how much I can generalize these results to the patients of other EDs) concerns me. This is a single trial run across only four hospitals in one part of the UK and it was run over two weeks in a summer period (which incidentally coincided with the London Olympics) and therefore this population may not reflect the wider UK ED population.

Applicability: Finally, there is no mention of the cost of the intervention, making an assessment of its value difficult to judge. It is also unclear how centres restricted this from surgical patients who may not have been allowed to eat or drink prior to assessment. The treatment itself may have risks attached to it, which the authors have not discussed such as delaying surgery.

PRACTICE PAPER 2: JOURNAL ARTICLE

ρ -BRB as a diagnostic marker of bacterial infection in the emergency department

Coughlan EP, Galloway RI, Bootland DJ

Department of Emergency Medicine, Brighton

Introduction

Sepsis is a significant cause of death, causing 37,000 deaths in the UK annually¹. The Surviving Sepsis Guidelines 2012 recommend early recognition, resuscitation and administration of antibiotics². Early recognition and intervention saves lives. However, some patients with an infection have minimal or even no symptoms or signs, making early recognition difficult³.

Diagnostic tests such as white cell count (WCC), C-reactive protein (CRP) and procalcitonin have all been used to try to identify septic patients early and provide information on prognosis^{4,5,6}.

Serum ρ -BRB (rho-B cell reticular bioprotein) is a two hundred and twenty five amino acid peptide, produced by B-cells on release of antibodies. The rho isomer of the B-cell reticular protein is produced when antibodies are produced against bacterial infection. The level of this in the bloodstream should correspond to the degree of the body's response and hence to the severity of the infection. Presence of high levels of this protein in the bloodstream should indicate a need for antibiotics.

We assessed the diagnostic ability of serum ρ -BRB to diagnose an infection (as determined by positive sputum or blood cultures). We hypothesized that for patients who present to the emergency department, ρ -BRB would be more sensitive than CRP and WCC as a marker of the presence of bacterial infection.

Methods

The study was performed between March and September 2013 in the emergency department (ED) of a tertiary referral hospital with 150,000 attendances per year. This was a prospective study of patients presenting to a tertiary university hospital. Consecutive adult patients (aged over eighteen years) who presented to our ED with signs and symptoms suggestive of a respiratory tract infection and who required phlebotomy were enrolled into the study after consent was obtained by one of the study personnel. Signs and symptoms of infection included cough, wheeze, productive sputum, shortness of breath, pyrexia and chest pain. Patients were excluded if they had a recently diagnosed infection, if they were already on antibiotics or if they had a recognized cause of reduced immunity (e.g. autoimmune disease, diabetes, taking chemotherapeutic agents, taking steroids etc.).

Baseline characteristics were obtained for all patients enrolled. Blood was also taken so that it could be tested for full blood count, C-reactive protein and serum ρ -BRB. Samples of blood and sputum were taken for culture. Staff responsible for processing of cultures and blood tests did so independently without knowledge of the

results of other samples. Treating clinicians were kept unaware of the results of the serum ρ -BRB throughout the study. The study was approved by the hospital ethics committee and verbal consent was obtained from patients for the additional sampling.

Data analysis

The best cutoff value was chosen using Youden’s index. During statistical analysis the Mann-Whitney U test was used to compare independent samples and the chi-squared test was used to compare proportions. All tests were two-sided and a p-value of 0.05 was considered significant. Statistical analysis was performed using StatsforU (Rahaman Software, Brighton, UK).

Results

Over the six-month period, 700 patients were enrolled into the study. Patients ranged in age from 19 to 103 years of age with a mean of 64 years and a median age of 65 years. Sixty percent of patients were male. Twenty seven percent of patients in total had a positive culture (either sputum or blood). See Table A3.5.

To assess the ρ -BRB test, we drew a ROC curve (see Figure A3.2). Cutoff values were used every 20 ng/l from 20 ng/l to 160 ng/l (see Table A3.6) and the sensitivity and specificity of the test was assessed at each of these levels.

Figure A3.2 shows a ROC curve illustrating the sensitivity and (1-specificity) of ρ -BRB. The area under the ROC curve was 0.93. After assessing each of the cutoff values, the optimum level for the ρ -BRB test was found to be 40 ng/l.

The threshold values of a white cell count greater than 12,000/mm³ and less than 4,000/mm³ were chosen to mirror the diagnostic criteria of SIRS in an attempt to improve its sensitivity. At that value, in this study, white cell count had a sensitivity of 36% and a specificity of 30% for diagnosing the presence of bacterial infection.

The normal CRP value of <10 mg/l was chosen as this is the current accepted value from the literature. Using this cutoff, CRP had a sensitivity of 66% and a specificity of 84%.

From the ROC curve we calculated the optimum value for ρ -BRB to be 40 ng/l. At this cutoff level, of the 189 positive culture results, 178 people had a positive ρ -BRB (>40 ng/l). Of 511 patients in the negative culture group, 433 had

TABLE A3.5
Baseline characteristics of patient in respect to culture results

	Positive cultures (95% CI) n = 189	Negative cultures (95% CI) n = 511
Age (years)	76.6 (74.0–79.2)	59.8 (57.4–62.2)
Sex (M:F)	115:74	304:207
WCC (/mm ³)	11.6 (9.7–13.6)	10.2 (8.3–12.2)
CRP	56.3 (47.2–65.4)	14.8 (11.7–17.8)
ρ -BRB	124 (114.5–133.5)	30.4 (26.1–34.6)

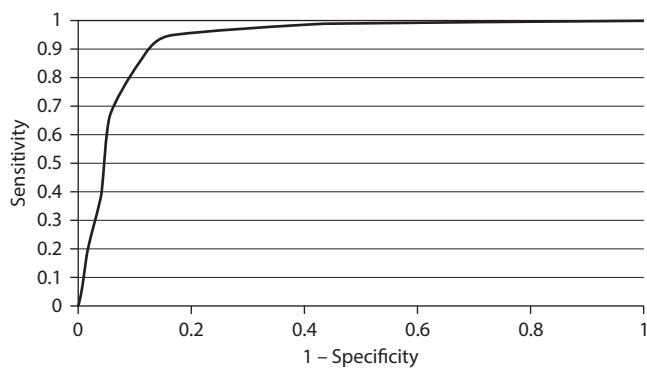


FIGURE A3.2 The ROC curve for ρ -BRB.

TABLE A3.6
Sensitivity and specificity for different cutoff values for ρ -BRB

Cutoff levels (ng/l)	Sensitivity	Specificity
20	0.99	0.60
40	0.94	0.85
60	0.84	0.90
80	0.70	0.94
100	0.61	0.95
120	0.44	0.96
140	0.32	0.97
160	0.27	0.97

a negative ρ -BRB (40 ng/l). This gives us a sensitivity of 94% and a specificity of 85% (Table A3.7).

TABLE A3.7
Sensitivity and specificity for white cell counts, C-reactive protein and ρ -BRB

	Sensitivity%	Specificity%
WCC (>12,000 or <4,000/mm ³)	36 (68/189) (95% CI 30–42)	70 (358/511) (95% CI 64–76)
CRP (>10 mg/l)	66 (126/189) (95% CI 60–72)	84 (429/511) (95% CI 80–88)
ρ -BRB (> 40 ng/l)	94 (178/189) (95% CI 90–98)	85 (433/511) (95% CI 81–89)

Discussion

The diagnosis of bacterial infection even in the acutely unwell patient is extremely difficult. The balance between prescribing antibiotics for those with a bacterial infection, who if left untreated may develop sepsis, and not prescribing antibiotics for those who do not have a bacterial infection is a difficult one. There are significant

potential risks for the individual patient and indeed to society if the wrong decision is made. The use of clinical history, examination and clinician gestalt has long been recognized as being insufficient if used in isolation³.

Many tests have been used to aid the identification of patients who required treatment with antibiotics. White cell counts are used daily in every hospital to aid in the diagnosis. In the meta-analysis by Ray et al., white cell count was found to have a low sensitivity (36%) but to improve the specificity (81%) of clinical examination and clinician gestalt⁴. These findings have been mirrored in this study.

C-reactive protein was initially discovered in 1930. It is a non-specific marker of inflammation, which rises in the presence of bacterial infection. It is often used in combination with clinical examination and clinician gestalt to diagnose bacterial infection. Although quite specific, it suffers from a lower sensitivity. The values we found (sensitivity of 66% and a specificity of 84%) for CRP are similar to those of other studies.

p-BRB is a new test for the diagnosis of bacterial infection. In this study we found it to have a sensitivity of 94% (90–98 95% CI) and a specificity of 85% (81–89 95% CI) when the cutoff value of 40 ng/l is used. The high sensitivity of this test compared to other markers of bacterial infection makes it extremely helpful in ruling out the need for antibiotics with a serum p-BRB of less than 40.

Further studies should focus on the use of this test in conjunction with other tests such as procalcitonin to hopefully improve the sensitivity so that the combination can be used as a rule out test.

Conclusions

The high sensitivity of this test may potentially allow us to rule out bacterial infection in patients and reduce the inappropriate prescription of antibiotics.

Competing interests

Dr. Coughlan has received a bursary from Apex pharmaceuticals to present this work at the ACEP conference in California.

References

1. O’Riordan S, Finn M. The epidemiology of sepsis. *BJM* 2009, 21(3): 431–435.
2. Surviving sepsis campaign: International guidelines for management of severe sepsis and septic shock. *Critical Care Medicine*, 2013, 41(2): 580–637.
3. Ransom P, McKenzie M, Duff M. Sepsis presentations to a tertiary referral center: Clinician gestalt is not enough. *JoS* 2010, 15(4): 2130–2138.
4. Wallman P, Beadsworth A, Bryant G. The use of White Cell Count as a prognostic indicator of infection: A meta-analysis. *JoS* 2008, 13 (6): 2860–2873.
5. Nelson M, Cottingham R. C-reactive protein in bacterial infection, a review of the current literature. *Jos* 2011, 16(1): 540–561.
6. Ghani R, Rahaman A. Procalcitonin as a marker of significant bacterial infection presenting to the Emergency department. *JoS* 2012, 17(6): 3142–3161.

PRACTICE PAPER 2: QUESTIONS

Bootland et al. p-BRB as a diagnostic marker of bacterial infection in the emergency department

Question 1

Provide a no more than 200 word summary of this paper in the box provided. Only the first 200 words will be considered – short bullet points are acceptable. Maximum of seven marks.

Question 2

- i) Describe your understanding of a confidence interval. (2 marks)
- ii) Is this more or less useful than a p-value? (2 marks)
- iii) The mean and 95% CI for age in the positive culture group was 76.6 (74.0–79.2) and in the negative culture group was 59.8 (57.4–62.2). If the CIs do not meet, what does this mean? (2 marks)

Question 3

Give three aspects of the methodology of this study that you consider were good. (3 marks)

Question 4

What is meant by the reference (gold) standard test? What do you think of the gold standard test in this study? (2 marks)

Question 5

What do you understand as sensitivity and specificity? What is meant by the area under the curve in respect to ROC curves? (2 marks)

Question 6

From the following statement 'Of the 189 positive culture results, 178 people had a positive ρ -BRB (> 40 ng/l). Of 511 patients in the negative culture group, 433 had a negative ρ -BRB (40 ng/l)' draw a 2×2 table and showing all your workings, calculate the following (8 marks):

- i) Sensitivity
- ii) Specificity
- iii) Positive predictive value
- iv) Negative predictive value
- v) Positive likelihood ratio
- vi) Negative likelihood ratio

Please show your workings.

Question 7

Your chief executive read this paper. He feels that we should get this test into the laboratory and that it will save the trust money by reducing the prescription of antibiotics. What will you say to him? (4 marks)

PRACTICE PAPER 2: ANSWERS

Bootland et al. ρ -BRB as a diagnostic marker of bacterial infection in the Emergency department

Question 1

Provide a no more than 200 word summary of this paper in the box provided. Only the first 200 words will be considered – short bullet points are acceptable. Maximum of seven marks.

Objectives: To assess the diagnostic capability of serum ρ -BRB in making a diagnosis of bacterial infection and to compare it to white cell count (WCC) and C-reactive protein (CRP)

Design: Diagnostic observation study

Setting: Emergency department of a tertiary referral hospital with 150,000 attendances per year

Population: Consecutive adults over eighteen with signs and symptoms suggestive of a respiratory tract infection and who required phlebotomy. Patients were excluded if they had a recently diagnosed infection, were already on treatment or had a recognized cause of reduced immunity.

Test under investigation: ρ -BRB

Gold standard: Positive culture of blood or sputum

Results: A ROC curve was performed for ρ -BRB with the area under the curve of 0.93 and optimal cutoff value of 40 ng/l. At this cutoff level ρ -BRB had a sensitivity of 94% (90–98 95% CI) and a specificity of 85% (81–89 95% CI). The sensitivity of ρ -BRB was better than both WCC (36%) and CRP (66%).

Conclusions: The high sensitivity of this test may potentially allow us to rule out bacterial infection in patients and reduce the inappropriate prescription of antibiotics.

180 words

Question 2

- i) Describe your understanding of a confidence interval. (2 marks)
- ii) Is this more or less useful than a p-value? (2 marks)
- iii) The mean and 95% CI for age in the positive culture group was 76.6 (74.0–79.2) and in the negative culture group was 59.8 (57.4–62.2). If the CIs do not meet, what does this mean? (2 marks)

- i) Confidence intervals: This is a range in which the population value is likely to lie.
- ii) P-values give you a yes or no answer to whether a result is statistically significant or not whereas a CI gives you a range in which the population value is likely to lie. Confidence intervals are generally considered to be more

clinically valuable as they allow you to interpret the likely size of the effect.

- iii) It means that there is a statistically significant difference between the two groups in respect to age because the CIs do not overlap.

Question 3

Give three aspects of the methodology of this study that you consider were good. (3 marks)

(In the exam, you would only need to give three points – but we have listed all possibilities.)

- It is a prospective study, which is generally regarded as higher quality than retrospective studies.
- Consecutive patients were enrolled reducing the chance for selection bias, i.e. reducing the chance that researchers influence which patients get enrolled. This form of sampling is more likely than convenience sampling to be representative of your population.
- The inclusion and exclusion criteria are both pragmatic and yet excluded important subgroups that might have influenced the test.
- All patients had the gold standard test and the test under investigation and they were assessed independently of each other.
- A ROC curve was drawn to assess the most appropriate cutoff value for the test.

Question 4

What is meant by the reference (gold) standard test? What do you think of the gold standard test in this study? (2 marks)

The gold standard test is the best test available to diagnose a condition, it should be used as the comparison test in diagnostic studies looking at the ability of a novel test to detect the same disease.

The gold standard test was a culture sample and was inappropriate in this study. The absence of a positive culture does not exclude bacterial infection. As such many of the patients who have a negative culture may have had a bacterial cause for their symptoms.

A better gold standard would be discharge diagnosis (which would be based on clinical examination and history, imaging, relevant blood and pathology tests and a period of observation).

Question 5

What do you understand as sensitivity and specificity? What is meant by the area under the curve in respect to ROC curves? (2 marks)

Sensitivity: the proportion of people with the disease who test positive. It is useful for those who are interested in ‘how good is this test’, e.g. clinical biochemists. A high sensitivity tells us a test is good at ruling out disease in a patient who tests negative (SnOut).

Specificity: the proportion of people without the disease who test negative. It is useful for those who are interested in ‘how good is this test’, e.g. clinical biochemists. A high specificity tells us a test is good at ruling in a disease in a patient who tests positive (SpIn).

The area under the curve of a ROC curve is a measure of how accurate a test is. It is measured between 0 and 1, the closer to 1 the better the test.

Question 6

From the following statement ‘Of the 189 positive culture results, 178 people had a positive ρ-BRB (> 40ng/l). Of 511 patients in the negative culture group, 433 had a negative ρ-BRB (40ng/l)’ draw a 2 × 2 table and showing all your workings, calculate the following (8 marks):

- i) Sensitivity
- ii) Specificity
- iii) Positive predictive value
- iv) Negative predictive value
- v) Positive likelihood ratio
- vi) Negative likelihood ratio

Please show your workings.

		Gold Standard	
		Culture positive	Culture negative
Test under investigation	>40 ng/l	178 (a)	78 (b)
	<40 ng/l	11 (c)	433 (d)

- i) Sensitivity: the proportion of people with the disease who test positive
= $a/a + c$
= $178/(178 + 11) = 178/189 = 0.94$
- ii) Specificity: the proportion of people without the disease who test negative
= $d/(b + d)$
= $433/(78 + 433) = 433/511 = 0.85$
- iii) PPV: the proportion of people with a positive test who have the disease
= $a/(a + b)$
= $178/(178 + 78) = 178/256 = 0.70$

- iv) NPV: the proportion of people with a negative test who don't have the disease
 $= d/(c + d)$
 $= 433/(11 + 433) = 433/444 = 0.98$
- v) Positive likelihood ratio
 $=$ How much more likely is a positive test to be found in a person with the disease (sensitivity) compared to without the disease (False positive = $1 - \text{Specificity}$)
 $= \text{Sensitivity}/1 - \text{Specificity}$
 $= 0.94/(1 - 0.85) = 0.94/0.15 = 6.3$
(or from first principles $= [a/(a + c)]/[b/(b + d)]$)
- vi) Negative likelihood ratio
 $=$ How much more likely is a -ve test in a person with the disease (False negative = $1 - \text{Sensitivity}$) compared to a -ve test without the disease (specificity)
 $= (1 - \text{Sensitivity})/\text{Specificity}$
 $= (1 - 0.94)/0.85 = 0.06/0.85 = 0.07$
(or from first principles $= [c/(c + a)]/[d/(d + b)]$)

In the actual exam you will get mathematics, which are numerically easier to do - you would never fail for a lack of multiple division knowledge.

Question 7

Your chief executive read this paper. He feels that we should get this test into the laboratory and that it will save the trust money by reducing the prescription of antibiotics. What will you say to him? (4 marks)

Politely decline this suggestion on the basis of:

- Single centre study and as such external validity is not assured.
- Although the gold standard used should diagnose a positive infection in patients, there will be a large percentage of the patients with a negative culture who also had a bacterial infection and as such should have antibiotics to prevent worsening of their illness, i.e. the internal validity is not good as the gold standard used in this paper is not appropriate.
- The sensitivity is too low to use it as a rule out test, hence you could not safely choose not to give antibiotics to these patients (poor applicability).
- No cost analysis was performed and as such we cannot know the cost of introducing the test and the daily cost of doing the test would reduce antibiotic costs even if the test did have the sensitivity to rule out the need for antibiotics (poor applicability).

Critical Appraisal for FCEM

The ability to read a paper, judge its quality, the importance of its results, and make a decision about whether to change practice based on the information given, is a core skill for all doctors. To be able to do this quickly and efficiently is, without a doubt, a skill needed by all time-pressured emergency doctors and one which is tested in the Fellowship of the College of Emergency Medicine (FCEM) examination.

Critical Appraisal for FCEM is the essential revision source for all those who want to pass the critical appraisal section of this exam. It is also required reading for those who want to incorporate evidence-based medicine into their everyday clinical practice.

Features

- Helps you become truly competent in critical appraisal
- Provides information in "Spod's Corner," which helps you reach the next level and excel
- Prepares you for the Critical Topic Review
- Contains two fictional practice papers to test and practise your knowledge

With its relaxed conversational style – yet crammed with essential information, key tips, and advice – this book is indispensable for all those wanting to achieve success in their FCEM and MCEM examinations.

Authors

Duncan Bootland MBBS, BSc, FCEM, Consultant in Emergency Medicine, Brighton and Sussex University Hospitals NHS Trust

Evan Coughlan MB BCh NUI, MRCSEd, DSEM, FCEM, DipTox, Emergency Medicine Consultant, Brighton and Sussex University Hospitals NHS Trust

Robert Galloway MBBS, BSc, MRCP, FCEM, PGCMedED, Emergency Medicine Consultant, Brighton and Sussex University Hospitals NHS Trust; Year Five Sub-Dean and Lead for Undergraduate Emergency Medicine, Brighton and Sussex Medical School

Stephanie Goubet MSc, Medical Statistician for Brighton and Sussex University Hospitals NHS Trust and Brighton and Sussex Medical School

K17686



CRC Press
Taylor & Francis Group
an informa business

www.crcpress.com

6000 Broken Sound Parkway, NW
Suite 300, Boca Raton, FL 33487
711 Third Avenue
New York, NY 10017
2 Park Square, Milton Park
Abingdon, Oxon OX14 4RN, UK

ISBN: 978-1-4441-8648-2



9 781444 186482