

01
02
03
04
05
06
07
08
09
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

Introduction to Meta-Analysis

Michael Borenstein

Biostat, Inc, New Jersey, USA.

Larry V. Hedges

Northwestern University, Evanston, USA.

Julian P. T. Higgins

MRC, Cambridge, UK.

Hannah R. Rothstein

Baruch College, New York, USA.



A John Wiley and Sons, Ltd., Publication

01 This edition first published 2009
02 © 2009 John Wiley & Sons, Ltd

03 *Registered office*

04 John Wiley & Sons Ltd, The Atrium, Southern Gate, Chichester, West Sussex, PO19 8SQ, United Kingdom

05 For details of our global editorial offices, for customer services and for information about how to apply for
06 permission to reuse the copyright material in this book please see our website at www.wiley.com.

07 The right of the author to be identified as the author of this work has been asserted in accordance with the
08 Copyright, Designs and Patents Act 1988.

09 All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted,
10 in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, except as
11 permitted by the UK Copyright, Designs and Patents Act 1988, without the prior permission of the publisher.

12 Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be
13 available in electronic books.

14 Designations used by companies to distinguish their products are often claimed as trademarks. All brand
15 names and product names used in this book are trade names, service marks, trademarks or registered
16 trademarks of their respective owners. The publisher is not associated with any product or vendor mentioned
17 in this book. This publication is designed to provide accurate and authoritative information in regard to the
18 subject matter covered. It is sold on the understanding that the publisher is not engaged in rendering
19 professional services. If professional advice or other expert assistance is required, the services of a competent
20 professional should be sought.

21 *Library of Congress Cataloguing-in-Publication Data*

22 Introduction to meta-analysis / Michael Borenstein . . . [et al.].

23 p. ; cm.

24 Includes bibliographical references and index.

25 ISBN 978-0-470-05724-7 (cloth)

26 1. Meta-analysis. I. Borenstein, Michael.

27 [DNLM: 1. Meta-Analysis as Topic. WA 950 I614 2009].

28 R853.M48I58 2009

29 610.72—dc22

2008043732

30 A catalogue record for this book is available from the British Library.

31 ISBN: 978-0-470-05724-7

32 Set in 10.5/13pt Times by Integra Software Services Pvt. Ltd, Pondicherry, India

33 Printed in the UK by TJ International, Padstow, Cornwall

34
35
36
37
38
39
40
41
42
43

Contents

01
02
03
04
05
06
07
08
09
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

List of Tables	xiii
List of Figures	xv
Acknowledgements	xix
Preface	xxi
Web site	xxix

PART 1: INTRODUCTION

1 HOW A META-ANALYSIS WORKS	3
Introduction	3
Individual studies	3
The summary effect	5
Heterogeneity of effect sizes	6
Summary points	7
2 WHY PERFORM A META-ANALYSIS	9
Introduction	9
The streptokinase meta-analysis	10
Statistical significance	11
Clinical importance of the effect	12
Consistency of effects	12
Summary points	14

PART 2: EFFECT SIZE AND PRECISION

3 OVERVIEW	17
Treatment effects and effect sizes	17
Parameters and estimates	18
Outline of effect size computations	19
4 EFFECT SIZES BASED ON MEANS	21
Introduction	21
Raw (unstandardized) mean difference D	21
Standardized mean difference, d and g	25
Response ratios	30
Summary points	32

01	5 EFFECT SIZES BASED ON BINARY DATA (2×2 TABLES)	33
02	Introduction	33
03	Risk ratio	34
04	Odds ratio	36
05	Risk difference	37
06	Choosing an effect size index	38
07	Summary points	39
08		
09	6 EFFECT SIZES BASED ON CORRELATIONS	41
10	Introduction	41
11	Computing r	41
12	Other approaches	43
13	Summary points	43
14		
15	7 CONVERTING AMONG EFFECT SIZES	45
16	Introduction	45
17	Converting from the log odds ratio to d	47
18	Converting from d to the log odds ratio	47
19	Converting from r to d	48
20	Converting from d to r	48
21	Summary points	49
22		
23		
24	8 FACTORS THAT AFFECT PRECISION	51
25	Introduction	51
26	Factors that affect precision	52
27	Sample size	52
28	Study design	53
29	Summary points	55
30		
31	9 CONCLUDING REMARKS	57
32		
33	PART 3: FIXED-EFFECT VERSUS RANDOM-EFFECTS MODELS	
34		
35	10 OVERVIEW	61
36	Introduction	61
37	Nomenclature	62
38		
39	11 FIXED-EFFECT MODEL	63
40	Introduction	63
41	The true effect size	63
42	Impact of sampling error	63
43		

01	Performing a fixed-effect meta-analysis	65
02	Summary points	67
03		
04	12 RANDOM-EFFECTS MODEL	69
05	Introduction	69
06	The true effect sizes	69
07	Impact of sampling error	70
08	Performing a random-effects meta-analysis	72
09	Summary points	74
10		
11	13 FIXED-EFFECT VERSUS RANDOM-EFFECTS MODELS	77
12	Introduction	77
13	Definition of a summary effect	77
14	Estimating the summary effect	78
15	Extreme effect size in a large study or a small study	79
16	Confidence interval	80
17	The null hypothesis	83
18	Which model should we use?	83
19	Model should not be based on the test for heterogeneity	84
20	Concluding remarks	85
21	Summary points	85
22		
23	14 WORKED EXAMPLES (PART 1)	87
24	Introduction	87
25	Worked example for continuous data (Part 1)	87
26	Worked example for binary data (Part 1)	92
27	Worked example for correlational data (Part 1)	97
28	Summary points	102
29		
30		
31	PART 4: HETEROGENEITY	
32		
33	15 OVERVIEW	105
34	Introduction	105
35	Nomenclature	106
36	Worked examples	106
37		
38	16 IDENTIFYING AND QUANTIFYING HETEROGENEITY	107
39	Introduction	107
40	Isolating the variation in true effects	107
41	Computing Q	109
42	Estimating τ^2	114
43	The I^2 statistic	117

01	Comparing the measures of heterogeneity	119
02	Confidence intervals for τ^2	122
03	Confidence intervals (or uncertainty intervals) for I^2	124
04	Summary points	125
05		
06	17 PREDICTION INTERVALS	127
07	Introduction	127
08	Prediction intervals in primary studies	127
09	Prediction intervals in meta-analysis	129
10	Confidence intervals and prediction intervals	131
11	Comparing the confidence interval with the prediction interval	132
12	Summary points	133
13		
14	18 WORKED EXAMPLES (PART 2)	135
15	Introduction	135
16	Worked example for continuous data (Part 2)	135
17	Worked example for binary data (Part 2)	139
18	Worked example for correlational data (Part 2)	143
19	Summary points	147
20		
21	19 SUBGROUP ANALYSES	149
22	Introduction	149
23	Fixed-effect model within subgroups	151
24	Computational models	161
25	Random effects with separate estimates of τ^2	164
26	Random effects with pooled estimate of τ^2	171
27	The proportion of variance explained	179
28	Mixed-effects model	183
29	Obtaining an overall effect in the presence of subgroups	184
30	Summary points	186
31		
32	20 META-REGRESSION	187
33	Introduction	187
34	Fixed-effect model	188
35	Fixed or random effects for unexplained heterogeneity	193
36	Random-effects model	196
37	Summary points	203
38		
39	21 NOTES ON SUBGROUP ANALYSES AND META-REGRESSION	205
40	Introduction	205
41	Computational model	205
42	Multiple comparisons	208
43	Software	209
	Analyses of subgroups and regression analyses are observational	209

01	Statistical power for subgroup analyses and meta-regression	210
02	Summary points	211
03		
04	PART 5: COMPLEX DATA STRUCTURES	
05		
06	22 OVERVIEW	215
07		
08	23 INDEPENDENT SUBGROUPS WITHIN A STUDY	217
09	Introduction	217
10	Combining across subgroups	218
11	Comparing subgroups	222
12	Summary points	223
13		
14	24 MULTIPLE OUTCOMES OR TIME-POINTS WITHIN A STUDY	225
15	Introduction	225
16	Combining across outcomes or time-points	226
17	Comparing outcomes or time-points within a study	233
18	Summary points	238
19		
20	25 MULTIPLE COMPARISONS WITHIN A STUDY	239
21	Introduction	239
22	Combining across multiple comparisons within a study	239
23	Differences between treatments	240
24	Summary points	241
25		
26	26 NOTES ON COMPLEX DATA STRUCTURES	243
27	Introduction	243
28	Summary effect	243
29	Differences in effect	244
30		
31	PART 6: OTHER ISSUES	
32		
33	27 OVERVIEW	249
34		
35	28 VOTE COUNTING – A NEW NAME FOR AN OLD PROBLEM	251
36	Introduction	251
37	Why vote counting is wrong	252
38	Vote counting is a pervasive problem	253
39	Summary points	255
40		
41	29 POWER ANALYSIS FOR META-ANALYSIS	257
42	Introduction	257
43	A conceptual approach	257
	In context	261
	When to use power analysis	262

01	Planning for precision rather than for power	263
02	Power analysis in primary studies	263
03	Power analysis for meta-analysis	267
04	Power analysis for a test of homogeneity	272
05	Summary points	275
06		
07	30 PUBLICATION BIAS	277
08	Introduction	277
09	The problem of missing studies	278
10	Methods for addressing bias	280
11	Illustrative example	281
12	The model	281
13	Getting a sense of the data	281
14	Is there evidence of any bias?	283
15	Is the entire effect an artifact of bias?	284
16	How much of an impact might the bias have?	286
17	Summary of the findings for the illustrative example	289
18	Some important caveats	290
19	Small-study effects	291
20	Concluding remarks	291
21	Summary points	291
22		
23	PART 7: ISSUES RELATED TO EFFECT SIZE	
24		
25	31 OVERVIEW	295
26		
27	32 EFFECT SIZES RATHER THAN p-VALUES	297
28	Introduction	297
29	Relationship between p -values and effect sizes	297
30	The distinction is important	299
31	The p -value is often misinterpreted	300
32	Narrative reviews vs. meta-analyses	301
33	Summary points	302
34		
35	33 SIMPSON'S PARADOX	303
36	Introduction	303
37	Circumcision and risk of HIV infection	303
38	An example of the paradox	305
39	Summary points	308
40		
41	34 GENERALITY OF THE BASIC INVERSE-VARIANCE METHOD	311
42	Introduction	311
43	Other effect sizes	312
44	Other methods for estimating effect sizes	315
	Individual participant data meta-analyses	316

01	Bayesian approaches	318
02	Summary points	319
03		
04	PART 8: FURTHER METHODS	
05	35 OVERVIEW	323
06		
07	36 META-ANALYSIS METHODS BASED ON DIRECTION AND p-VALUES	325
08	Introduction	325
09	Vote counting	325
10	The sign test	325
11	Combining p -values	326
12	Summary points	330
13		
14	37 FURTHER METHODS FOR DICHOTOMOUS DATA	331
15	Introduction	331
16	Mantel-Haenszel method	331
17	One-step (Peto) formula for odds ratio	336
18	Summary points	339
19		
20	38 PSYCHOMETRIC META-ANALYSIS	341
21	Introduction	341
22	The attenuating effects of artifacts	342
23	Meta-analysis methods	344
24	Example of psychometric meta-analysis	346
25	Comparison of artifact correction with meta-regression	348
26	Sources of information about artifact values	349
27	How heterogeneity is assessed	349
28	Reporting in psychometric meta-analysis	350
29	Concluding remarks	351
30	Summary points	351
31		
32	PART 9: META-ANALYSIS IN CONTEXT	
33	39 OVERVIEW	355
34		
35	40 WHEN DOES IT MAKE SENSE TO PERFORM A META-ANALYSIS?	357
36	Introduction	357
37	Are the studies similar enough to combine?	358
38	Can I combine studies with different designs?	359
39	How many studies are enough to carry out a meta-analysis?	363
40	Summary points	364
41		
42	41 REPORTING THE RESULTS OF A META-ANALYSIS	365
43	Introduction	365
	The computational model	366

01	Forest plots	366
02	Sensitivity analysis	368
03	Summary points	369
04		
05	42 CUMULATIVE META-ANALYSIS	371
06	Introduction	371
07	Why perform a cumulative meta-analysis?	373
08	Summary points	376
09		
10	43 CRITICISMS OF META-ANALYSIS	377
11	Introduction	377
12	One number cannot summarize a research field	378
13	The file drawer problem invalidates meta-analysis	378
14	Mixing apples and oranges	379
15	Garbage in, garbage out	380
16	Important studies are ignored	381
17	Meta-analysis can disagree with randomized trials	381
18	Meta-analyses are performed poorly	384
19	Is a narrative review better?	385
20	Concluding remarks	386
21	Summary points	386
22		
23	PART 10: RESOURCES AND SOFTWARE	
24		
25	44 SOFTWARE	391
26	Introduction	391
27	The software	392
28	Three examples of meta-analysis software	393
29	Comprehensive Meta-Analysis (CMA) 2.0	395
30	RevMan 5.0	398
31	Stata macros with Stata 10.0	400
32	Summary points	403
33		
34	45 BOOKS, WEB SITES AND PROFESSIONAL ORGANIZATIONS	405
35	Books on systematic review methods	405
36	Books on meta-analysis	405
37	Web sites	406
38		
39	REFERENCES	409
40	INDEX	415
41		
42		
43		

Preface

01
02
03
04
05
06
07
08
09 In his best-selling book *Baby and Child Care*, Dr. Benjamin Spock wrote ‘I think it
10 is preferable to accustom a baby to sleeping on his stomach from the beginning if he
11 is willing’. This statement was included in most editions of the book, and in most of
12 the 50 million copies sold from the 1950s into the 1990s. The advice was not
13 unusual, in that many pediatricians made similar recommendations at the time.

14 During this same period, from the 1950s into the 1990s, more than 100,000 babies
15 died of sudden infant death syndrome (SIDS), also called *crib death* in the United
16 States and *cot death* in the United Kingdom, where a seemingly healthy baby goes
17 to sleep and never wakes up.

18 In the early 1990s, researchers became aware that the risk of SIDS decreased by at
19 least 50% when babies were put to sleep on their backs rather than face down.
20 Governments in various countries launched educational initiatives such as the *Back
21 to sleep* campaigns in the UK and the US, which led to an immediate and dramatic
22 drop in the number of SIDS deaths.

23 While the loss of more than 100,000 children would be unspeakably sad in any
24 event, the real tragedy lies in the fact that many of these deaths could have been
25 prevented. Gilbert *et al.* (2005) write

26 ‘Advice to put infants to sleep on the front for nearly half a century was contrary to
27 evidence available from 1970 that this was likely to be harmful. Systematic review of
28 preventable risk factors for SIDS from 1970 would have led to earlier recognition of
29 the risks of sleeping on the front and might have prevented over 10,000 infant deaths
30 in the UK and at least 50,000 in the Europe, the USA and Australasia.’

31 32 AN ETHICAL IMPERATIVE

33
34 This example is one of several cited by Sir Iain Chalmers in a talk entitled *The
35 scandalous failure of scientists to cumulate scientifically* (Chalmers, 2006). The
36 theme of this talk was that we live in a world where the utility of almost any
37 intervention will be tested repeatedly, and that rather than looking at any study in
38 isolation, we need to look at the body of evidence. While not all systematic reviews
39 carry the urgency of SIDS, the logic of looking at the body of evidence, rather than
40 trying to understand studies in isolation, is always compelling.

41 Meta-analysis refers to the statistical synthesis of results from a series of studies.
42 While the statistical procedures used in a meta-analysis can be applied to any set of
43 data, the synthesis will be meaningful only if the studies have been collected

01 systematically. This could be in the context of a systematic review, the process of
02 systematically locating, appraising, and then synthesizing data from a large number
03 of sources. Or, it could be in the context of synthesizing data from a select group of
04 studies, such as those conducted by a pharmaceutical company to assess the efficacy
05 of a new drug.

06 If a treatment effect (or effect size) is consistent across the series of studies, these
07 procedures enable us to report that the effect is robust across the kinds of popula-
08 tions sampled, and also to estimate the magnitude of the effect more precisely than
09 we could with any of the studies alone. If the treatment effect varies across the series
10 of studies, these procedures enable us to report on the range of effects, and may
11 enable us to identify factors associated with the magnitude of the effect size.

13 FROM NARRATIVE REVIEWS TO SYSTEMATIC REVIEWS

14
15 Prior to the 1990s, the task of combining data from multiple studies had been
16 primarily the purview of the narrative review. An expert in a given field would
17 read the studies that addressed a question, summarize the findings, and then arrive at
18 a conclusion – for example, that the treatment in question was, or was not, effective.
19 However, this approach suffers from some important limitations.

20 One limitation is the subjectivity inherent in this approach, coupled with the lack
21 of transparency. For example, different reviewers might use different criteria for
22 deciding which studies to include in the review. Once a set of studies has been
23 selected, one reviewer might give more credence to larger studies, while another
24 gives more credence to ‘quality’ studies and yet another assigns a comparable
25 weight to all studies. One reviewer may require a substantial body of evidence
26 before concluding that a treatment is effective, while another uses a lower threshold.
27 In fact, there are examples in the literature where two narrative reviews come to
28 opposite conclusions, with one reporting that a treatment is effective while the other
29 reports that it is not. As a rule, the narrative reviewer will not articulate (and may not
30 even be fully aware of) the decision-making process used to synthesize the data and
31 arrive at a conclusion.

32 A second limitation of narrative reviews is that they become *less useful as more*
33 *information becomes available*. The thought process required for a synthesis requires
34 the reviewer to capture the finding reported in each study, to assign an appropriate
35 *weight* to that finding, and then to synthesize these findings across all studies in the
36 synthesis. While a reviewer may be able to synthesize data from a few studies in their
37 head, the process becomes difficult and eventually untenable as the number of studies
38 increases. This is true even when the treatment effect (or effect size) is consistent from
39 study to study. Often, however, the treatment effect will vary as a function of study-
40 level covariates, such as the patient population, the dose of medication, the outcome
41 variable, and other factors. In these cases, a proper synthesis requires that the
42 researcher be able to understand how the treatment effect varies as a function of
43 these variables, and the narrative review is poorly equipped to address these kinds of
issues.

THE SYSTEMATIC REVIEW AND META-ANALYSIS

For these reasons, beginning in the mid 1980s and taking root in the 1990s, researchers in many fields have been moving away from the narrative review, and adopting systematic reviews and meta-analysis.

For systematic reviews, a clear set of rules is used to search for studies, and then to determine which studies will be included in or excluded from the analysis. Since there is an element of subjectivity in setting these criteria, as well as in the conclusions drawn from the meta-analysis, we cannot say that the systematic review is entirely objective. However, because all of the decisions are specified clearly, the mechanisms are transparent.

A key element in most systematic reviews is the statistical synthesis of the data, or the meta-analysis. Unlike the narrative review, where reviewers implicitly assign some level of importance to each study, in meta-analysis the weights assigned to each study are based on mathematical criteria that are specified in advance. While the reviewers and readers may still differ on the substantive meaning of the results (as they might for a primary study), the statistical analysis provides a transparent, objective, and replicable framework for this discussion.

The formulas used in meta-analysis are extensions of formulas used in primary studies, and are used to address similar kinds of questions to those addressed in primary studies. In primary studies we would typically report a mean and standard deviation for the subjects. If appropriate, we might also use analysis of variance or multiple regression to determine if (and how) subject scores were related to various factors. Similarly, in a meta-analysis, we might report a mean and standard deviation for the treatment effect. And, if appropriate, we would also use procedures analogous to analysis of variance or multiple regression to assess the relationship between the effect and study-level covariates.

Meta-analyses are conducted for a variety of reasons, not only to synthesize evidence on the effects of interventions or to support evidence-based policy or practice. The purpose of the meta-analysis, or more generally, the purpose of any research synthesis has implications for *when* it should be performed, what model should be used to analyze the data, what sensitivity analyses should be undertaken, and how the results should be interpreted. Losing sight of the fact that meta-analysis is a tool with multiple applications causes confusion and leads to pointless discussions about *what is the right way to perform a research synthesis*, when there is no single right way. It all depends on the purpose of the synthesis, and the data that are available. Much of this book will expand on this idea.

META-ANALYSIS IS USED IN MANY FIELDS OF RESEARCH

In medicine, systematic reviews and meta-analysis form the core of a movement to ensure that medical treatments are based on the best available empirical data. For example, The Cochrane Collaboration has published the results of over 3700 meta-analyses (as of January 2009) which synthesize data on treatments in all areas of

01 health care including headaches, cancer, allergies, cardiovascular disease, pain pre-
02 ventions, and depression. The reviews look at interventions relevant to neo-natal care,
03 childbirth, infant and childhood diseases, as well as diseases common in adolescents,
04 adults, and the elderly. The kinds of interventions assessed include surgery, drugs,
05 acupuncture, and social interventions. BMJ publishes a series of journals on Evidence
06 Based Medicine, built on the results from systematic reviews. Systematic reviews and
07 meta-analyses are also used to examine the performance of diagnostic tests, and of
08 epidemiological associations between exposure and disease prevalence, among other
09 topics.

10 Pharmaceutical companies usually conduct a series of studies to assess the
11 efficacy of a drug. They use meta-analysis to synthesize the data from these studies,
12 yielding a more powerful test (and more precise estimate) of the drug's effect.
13 Additionally, the meta-analysis provides a framework for evaluating the series of
14 studies as a whole, rather than looking at each in isolation. These analyses play a
15 role in internal research, in submissions to governmental agencies, and in market-
16 ing. Meta-analyses are also used to synthesize data on adverse events, since these
17 events are typically rare and we need to accumulate information over a series of
18 studies to properly assess the risk of these events.

19 In the field of education, meta-analysis has been applied to topics as diverse as
20 the comparison of distance education with traditional classroom learning, assess-
21 ment of the impact of schooling on developing economies, and the relationship
22 between teacher credentials and student achievement. Results of these and similar
23 meta-analyses have influenced practice and policy in various locations around the
24 world.

25 In psychology, meta-analysis has been applied to basic science as well as in
26 support of evidence-based practice. It has been used to assess personality change
27 over the life span, to assess the influence of media violence on aggressive
28 behavior, and to examine gender differences in mathematics ability, leadership,
29 and nonverbal communication. Meta-analyses of psychological interventions have
30 been use to compare and select treatments for psychological problems, including
31 obsessive-compulsive disorder, impulsivity disorder, bulimia nervosa, depression,
32 phobias, and panic disorder.

33 In the field of criminology, government agencies have funded meta-analyses to
34 examine the relative effectiveness of various programs in reducing criminal beha-
35 vior. These include initiatives to prevent delinquency, reduce recidivism, assess the
36 effectiveness of different strategies for police patrols, and for the use of special
37 courts to deal with drug-related crimes.

38 In business, meta-analyses of the predictive validity of tests that are used as part
39 of the hiring process, have led to changes in the types of tests that are used to select
40 employees in many organizations. Meta-analytic results have also been used to
41 guide practices for the reduction of absenteeism, turnover, and counterproductive
42 behavior, and to assess the effectiveness of programs used to train employees.

43 In the field of ecology, meta-analyses are being used to identify the environmental
impact of wind farms, biotic resistance to exotic plant invasion, the effects of changes

01 in the marine food chain, plant reactions to global climate change, the effectiveness of
02 conservation management interventions, and to guide conservation efforts.

04 **META-ANALYSIS AS PART OF THE RESEARCH PROCESS**

05
06 Systematic reviews and meta-analyses are used to synthesize the available evidence
07 for a given question to inform policy, as in the examples cited above from medicine,
08 social science, business, ecology, and other fields. While this is probably the most
09 common use of the methodology, meta-analysis can also play an important role in
10 other parts of the research process.

11 Systematic reviews and meta-analyses can play a role in designing new research.
12 As a first step, they can help determine whether the planned study is necessary.
13 It may be possible to find the required information by synthesizing data from prior
14 studies, and in this case, the research should not be performed. Iain Chalmers (2007)
15 made this point in an article entitled *The lethal consequences of failing to make use*
16 *of all relevant evidence about the effects of medical treatments: the need for*
17 *systematic reviews*.

18 In the event that the new study *is needed*, the meta-analysis may be useful in
19 helping to design that study. For example, the meta-analysis may show that in the
20 prior studies one outcome index had proven to be more sensitive than others, or that
21 a specific mode of administration had proven to be more effective than others, and
22 should be used in the planned study as well.

23 For these reasons, various government agencies, including institutes of health in
24 various countries, have been encouraging (or requiring) researchers to conduct a
25 meta-analysis of existing research prior to undertaking new funded studies.

26 The systematic review can also play a role in the publication of any new primary
27 study. In the introductory section of the publication, a systematic review can help to
28 place the new study in context by describing what we knew before, and what we
29 hoped to learn from the new study. In the discussion section of the publication, a
30 systematic review allows us to address not only the information provided by the new
31 study, but the body of evidence as enhanced by the new study. Iain Chalmers and
32 Michael Clarke (1998) see this approach as a way to avoid studies being reported
33 without context, which they refer to as 'Islands in Search of Continents'. Systematic
34 reviews would provide this context in a more rigorous and transparent manner than
35 the narrative reviews that are typically used for this purpose.

37 **THE INTENDED AUDIENCE FOR THIS BOOK**

38
39 Since meta-analysis is a relatively new field, many people, including those who
40 actually use meta-analysis in their work, have not had the opportunity to learn about
41 it systematically. We hope that this volume will provide a framework that allows
42 them to understand the logic of meta-analysis, as well as how to apply and interpret
43 meta-analytic procedures properly.

01 This book is aimed at researchers, clinicians, and statisticians. Our approach is
02 primarily conceptual. The reader will be able to skip the formulas and still under-
03 stand, for example, the differences between fixed-effect and random-effects analy-
04 sis, and the mechanisms used to assess the dispersion in effects from study to study.
05 However, for those with a statistical orientation, we include all the relevant for-
06 mulas, along with worked examples. Additionally, the spreadsheets and data files
07 can be downloaded from the web at www.Meta-Analysis.com.

08 This book can be used as the basis for a course in meta-analysis. Supplementary
09 materials and exercises are posted on the book's web site.

10 This volume is intended for readers from various substantive fields, including
11 medicine, epidemiology, social science, business, ecology, and others. While we
12 have included examples from many of these disciplines, the more important mes-
13 sage is that meta-analytic methods that may have developed in any one of these
14 fields have application to all of them.

15 Since our goal in using these examples is to explain the meta-analysis itself rather
16 than to address the substantive issues, we provide only the information needed for
17 this purpose. For example, we may present an analysis showing that a treatment
18 reduces pain, while ignoring other analyses that show the same treatment increases
19 the risk of adverse events. Therefore, any reader interested in the substantive issues
20 addressed in an example should not rely on this book for that purpose.

21 22 **AN OUTLINE OF THIS BOOK'S CONTENTS**

23
24 Part 1 is an introduction to meta-analysis. We present a completed meta-analysis to
25 serve as an example, and highlight the elements of this analysis – the effect size for
26 each study, the summary effect, the dispersion of effects across studies, and so on.
27 Our intent is to show where each element fits into the analysis, and thus provide the
28 reader with a context as they move on to the subsequent parts of the book where
29 each of the elements is explored in detail.

30 Part 2 introduces the effect sizes, such as the standardized mean difference or the
31 risk ratio, that are computed for each study, and that serve as the unit of currency in
32 the meta-analysis. We also discuss factors that determine the variance of an effect
33 size and show how to compute the variance for each study, since this affects the
34 weight assigned to that study in the meta-analysis.

35 Part 3 discusses the two computational models used in the vast majority of meta-
36 analyses, the fixed-effect model and the random-effects model. We discuss the
37 conceptual and practical differences between the two, and show how to compute a
38 summary effect using either one.

39 Part 4 focuses on the issue of dispersion in effect sizes, the fact that the effect size
40 varies from one study to the next. We discuss methods to quantify the heterogeneity,
41 to test it, to incorporate it in the weighting scheme, and to understand it in a
42 substantive as well as a statistical context. Then, we discuss methods to explain
43 the heterogeneity. These include subgroup analyses to compare the effect in

01 different subgroups of studies (analogous to analysis of variance in primary stu-
02 dies), and meta-regression (analogous to multiple regression).

03 Part 5 shows how to work with complex data structures. These include studies
04 that report an effect size for two or more independent subgroups, for two or more
05 outcomes or time-points, and for two or more comparison groups (such as two
06 treatments being compared with the same control).

07 Part 6 is used to address three separate issues. One chapter discusses the proce-
08 dure called vote counting, common in narrative reviews, and explains the problems
09 with this approach. One chapter discusses statistical power for a meta-analysis. We
10 show how meta-analysis often (but not always) yields a more powerful test of the
11 null than do any of the included studies. Another chapter addresses the question of
12 publication bias. We explain what this is, and discuss methods that have been
13 developed to assess its potential impact.

14 Part 7 focuses on the issue of why we work with effect sizes in a meta-analysis. In
15 one chapter we explain why we work with effect sizes rather than p -values. In
16 another we explain why we compute an effect size for each study, rather than
17 summing data over all studies and then computing an effect size for the summed
18 data. The final chapter in this part shows how the use of inverse-variance weights
19 can be extended to other applications including Bayesian meta-analysis and ana-
20 lyses based on individual participant data.

21 Part 8 includes chapters on methods that are sometimes used in meta-analysis but
22 that fall outside the central narrative of this volume. These include meta-analyses
23 based on p -values, alternate approaches (such as the Mantel-Haenszel method) for
24 assigning study weights, and options sometimes used in psychometric meta-analyses.

25 Part 9 is dedicated to a series of general issues related to meta-analysis. We
26 address the question of when it makes sense to perform a meta-analysis. This Part is
27 also the location for a series of chapters on separate issues such as reporting the
28 results of a meta-analysis, and the proper use of cumulative meta-analysis. Finally,
29 we discuss some of the criticisms of meta-analysis and try to put them in context.

30 Part 10 is a discussion of resources for meta-analysis and systematic reviews.
31 This includes an overview of several computer programs for meta-analysis. It also
32 includes a discussion of organizations that promote the use of systematic reviews
33 and meta-analyses in specific fields, and a list of useful web sites.

34 35 36 **WHAT THIS BOOK DOES NOT COVER**

37 **Other elements of a systematic review**

38
39 This book deals only with meta-analysis, the statistical formulas and methods used
40 to synthesize data from a set of studies. A meta-analysis can be applied to any data,
41 but if the goal of the analysis is to provide a synthesis of a body of data from various
42 sources, then it is usually imperative that the data be compiled as part of a
43 systematic review.

01 A systematic review incorporates many components, such as specification of
02 the question to be addressed, determination of methods to be used for searching
03 the literature and for including or excluding studies, specification of mechanisms
04 to appraise the validity of the included studies, specification of methods to be
05 used for performing the statistical analysis, and a mechanism for disseminating
06 the results.

07 If the entire review is performed properly, so that the search strategy matches the
08 research question, and yields a reasonably complete and unbiased collection of
09 the relevant studies, then (providing that the included studies are themselves valid)
10 the meta-analysis will also be addressing the intended question. On the other hand,
11 if the search strategy is flawed in concept or execution, or if the studies are
12 providing biased results, then problems exist in the review that the meta-analysis
13 cannot correct.

14 In Part 10 we include an annotated listing of suggested readings for the other
15 components in the systematic review, but these components are not otherwise
16 addressed in this volume.

17 18 19 **Other meta-analytic methods**

20 In this volume we focus primarily on meta-analyses of effect sizes. That is, analyses
21 where each study yields an estimate of some statistic (a standardized mean differ-
22 ence, a risk ratio, a prevalence, and so on) and our goal is to assess the dispersion in
23 these effects and (if appropriate) compute a summary effect. The vast majority of
24 meta-analyses performed use this approach. We deal only briefly (see Part 8) with
25 other approaches, such as meta-analyses that combine *p*-values rather than effect
26 sizes. We do not address meta-analysis of diagnostic tests.

27 28 29 **Further Reading**

- 30
31 Chalmers, I. (2007). The lethal consequences of failing to make use of all relevant evidence about
32 the effects of medical treatments: the need for systematic reviews. In P. Rothwell(ed.),
33 *Treating Individuals*, ed. London: Lancet: 37–58.
- 34 Chalmers, I., Hedges, L.V. & Cooper, H. (2002). A brief history of research synthesis. *Evaluation*
35 *in the Health Professions*. 25(1): 12–37.
- 36 Clarke, M, Hopewell, S. & Chalmers, I. (2007). Reports of clinical trials should begin and end
37 with up-to-date systematic reviews of other relevant evidence: a status report. *Journal of the*
38 *Royal Society of Medicine*. 100: 187–190.
- 39 Hunt, M. (1999). *How Science Takes Stock: The Story of Meta-analysis*. New York: Russell Sage
40 Foundation.
- 41 Sutton, A.J. & Higgins, J.P.T. (2008). Recent developments in meta-analysis. *Statistics in*
42 *Medicine* 27: 625–650.
- 43

Fixed-Effect Versus Random-Effects Models

01
02
03
04
05
06
07
08
09
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

Overview

01
02
03
04
05
06
07
08 Introduction
09 Nomenclature
10
11

INTRODUCTION

12
13
14
15 Most meta-analyses are based on one of two statistical models, the fixed-effect
16 model or the random-effects model.

17
18 Under the fixed-effect model we assume that there is one *true effect size* (hence
19 the term *fixed effect*) which underlies all the studies in the analysis, and that all
20 differences in observed effects are due to sampling error. While we follow the
21 practice of calling this a fixed-effect model, a more descriptive term would be a
22 *common-effect* model. In either case, we use the singular (*effect*) since there is only
23 one true effect.

24
25 By contrast, under the random-effects model we allow that the true effect could
26 vary from study to study. For example, the effect size might be higher (or lower) in
27 studies where the participants are older, or more educated, or healthier than in others,
28 or when a more intensive variant of an intervention is used, and so on. Because studies
29 will differ in the mixes of participants and in the implementations of interventions,
30 among other reasons, there may be *different effect sizes* underlying different studies.
31 If it were possible to perform an infinite number of studies (based on the inclusion
32 criteria for our analysis), the true effect sizes for these studies would be distributed
33 about some mean. The effect sizes in the studies that actually *were performed* are
34 assumed to represent a random sample of these effect sizes (hence the term *random*
35 *effects*). Here, we use the plural (*effects*) since there is an array of true effects.

36
37 In the chapters that follow we discuss the two models and show how to compute a
38 summary effect using each one. Because the computations for a summary effect are
39 not always intuitive, it helps to keep in mind that the summary effect is nothing
40 more than the mean of the effect sizes, with more weight assigned to the more
41 precise studies. We need to consider what we mean by the *more precise* studies and

	True effect	Observed effect
Study	●	■
Combined	▼	◆

Figure 10.1 Symbols for true and observed effects.

how this translates into a study weight (this depends on the model), but not lose track of the fact that we are simply computing a weighted mean.

NOMENCLATURE

Throughout this Part we distinguish between a true effect size and an observed effect size. A study's *true effect size* is the effect size in the underlying population, and is the effect size that we would observe if the study had an infinitely large sample size (and therefore no sampling error). A study's *observed effect size* is the effect size that is actually observed.

In the schematics we use different symbols to distinguish between true effects and observed effects. For individual studies we use a circle for the former and a square for the latter (see Figure 10.1). For summary effects we use a triangle for the former and a diamond for the latter.

Worked examples

In meta-analysis the same formulas apply regardless of the effect size being used. To allow the reader to work with an effect size of their choosing, we have separated the formulas (which are presented in the following chapters) from the worked examples (which are presented in Chapter 14). There, we provide a worked example for the standardized mean difference, one for the odds ratio, and one for correlations.

The reader is encouraged to select one of the worked examples and follow the details of the computations while studying the formulas. The three datasets and all computations are available as Excel spreadsheets on the book's web site.

Fixed-Effect Model

01
02
03
04
05
06
07
08 Introduction
09 The true effect size
10 Impact of sampling error
11 Performing a fixed-effect meta-analysis
12
13
14
15

INTRODUCTION

17 In this chapter we introduce the fixed-effect model. We discuss the assumptions of
18 this model, and show how these are reflected in the formulas used to compute a
19 summary effect, and in the meaning of the summary effect.
20

THE TRUE EFFECT SIZE

23 Under the fixed-effect model we assume that all studies in the meta-analysis share a
24 common (true) effect size. Put another way, all factors that could influence the
25 effect size are the same in all the studies, and therefore the true effect size is the
26 same (hence the label *fixed*) in all the studies. We denote the true (unknown) effect
27 size by theta (θ)

28 In Figure 11.1 the true overall effect size is 0.60 and this effect (represented by a
29 triangle) is shown at the bottom. The true effect for each study is represented by a
30 circle. Under the definition of a fixed-effect model the true effect size for each study
31 must also be 0.60, and so these circles are aligned directly above the triangle.
32

IMPACT OF SAMPLING ERROR

34 Since all studies share the same true effect, it follows that the observed effect size
35 varies from one study to the next only because of the random error inherent in each
36 study. If each study had an infinite sample size the sampling error would be zero and
37 the observed effect for each study would be the same as the true effect. If we were to
38 plot the observed effects rather than the true effects, the observed effects would
39 exactly coincide with the true effects.
40

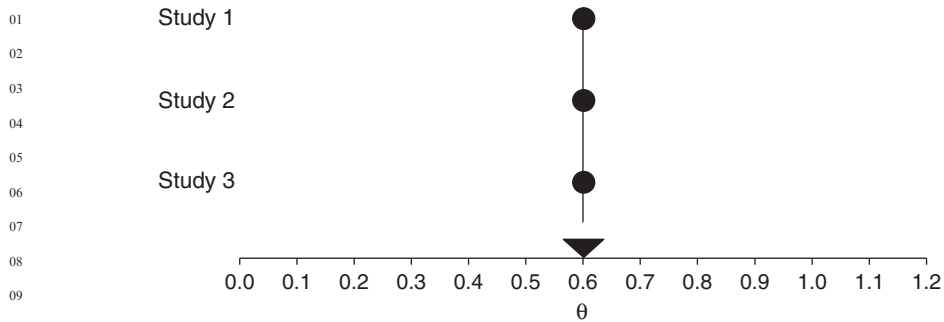


Figure 11.1 Fixed-effect model – true effects.

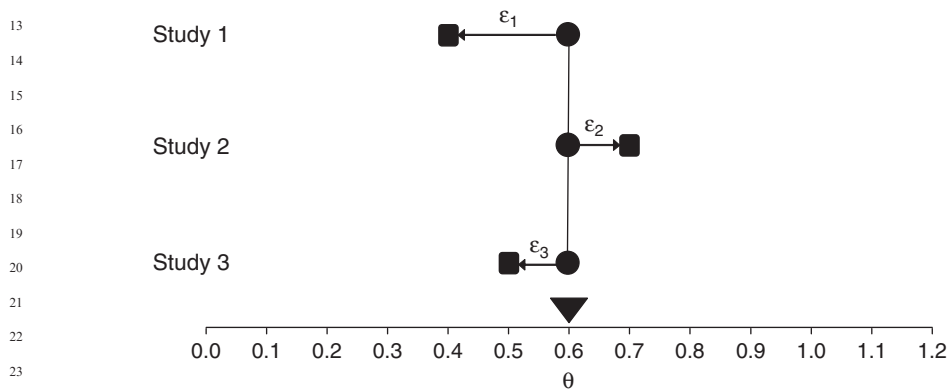


Figure 11.2 Fixed-effect model – true effects and sampling error.

In practice, of course, the sample size in each study is not infinite, and so there is sampling error and the effect observed in the study is not the same as the true effect. In Figure 11.2 the true effect for each study is still 0.60 (as depicted by the circles) but the observed effect (depicted by the squares) differs from one study to the next.

In Study 1 the sampling error (ϵ_1) is -0.20 , which yields an observed effect (Y_1) of

$$Y_1 = 0.60 - 0.20 = 0.40.$$

In Study 2 the sampling error (ϵ_2) is 0.10 , which yields an observed effect (Y_2) of

$$Y_2 = 0.60 + 0.10 = 0.70.$$

In Study 3 the sampling error (ϵ_3) is -0.10 , which yields an observed effect (Y_3) of

$$Y_3 = 0.60 - 0.10 = 0.50.$$

More generally, the observed effect Y_i for any study is given by the population mean plus the sampling error in that study. That is,

$$Y_i = \theta + \epsilon_i. \quad (11.1)$$

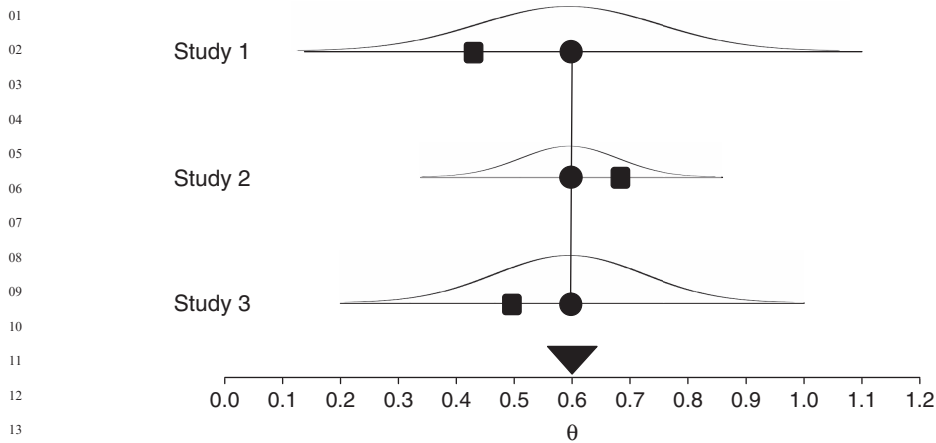


Figure 11.3 Fixed-effect model – distribution of sampling error.

While the error in any given study is random, we *can* estimate the sampling distribution of the errors. In Figure 11.3 we have placed a normal curve about the true effect size for each study, with the width of the curve being based on the variance in that study. In Study 1 the sample size was small, the variance large, and the observed effect is likely to fall anywhere in the relatively wide range of 0.20 to 1.00. By contrast, in Study 2 the sample size was relatively large, the variance is small, and the observed effect is likely to fall in the relatively narrow range of 0.40 to 0.80. (The width of the normal curve is based on the square root of the variance, or standard error).

PERFORMING A FIXED-EFFECT META-ANALYSIS

In an actual meta-analysis, of course, rather than starting with the population effect and making projections about the observed effects, we work backwards, starting with the observed effects and trying to estimate the population effect. In order to obtain the most precise estimate of the population effect (to minimize the variance) we compute a weighted mean, where the weight assigned to each study is the inverse of that study's variance. Concretely, the weight assigned to each study in a fixed-effect meta-analysis is

$$W_i = \frac{1}{V_{Y_i}}, \quad (11.2)$$

where V_{Y_i} is the within-study variance for study (i). The weighted mean (M) is then computed as

$$M = \frac{\sum_{i=1}^k W_i Y_i}{\sum_{i=1}^k W_i}, \quad (11.3)$$

that is, the sum of the products $W_i Y_i$ (effect size multiplied by weight) divided by the sum of the weights.

The variance of the summary effect is estimated as the reciprocal of the sum of the weights, or

$$V_M = \frac{1}{\sum_{i=1}^k W_i} \quad (11.4)$$

and the estimated standard error of the summary effect is then the square root of the variance,

$$SE_M = \sqrt{V_M}. \quad (11.5)$$

Then, 95% lower and upper limits for the summary effect are estimated as

$$LL_M = M - 1.96 \times SE_M \quad (11.6)$$

and

$$UL_M = M + 1.96 \times SE_M. \quad (11.7)$$

Finally, a Z-value to test the null hypothesis that the common true effect θ is zero can be computed using

$$Z = \frac{M}{SE_M}. \quad (11.8)$$

For a one-tailed test the p -value is given by

$$p = 1 - \Phi(\pm|Z|), \quad (11.9)$$

where we choose '+' if the difference is in the expected direction and '-' otherwise, and for a two-tailed test by

$$p = 2 \left[1 - \Phi(|Z|) \right], \quad (11.10)$$

where $\Phi(Z)$ is the standard normal cumulative distribution. This function is tabled in many introductory statistics books, and is implemented in Excel as the function =NORMSDIST(Z).

Illustrative example

We suggest that you turn to a worked example for the fixed-effect model before proceeding to the random-effects model. A worked example for the standardized

01 mean difference (Hedges' g) is on page 87, a worked example for the odds ratio is on
02 page 92, and a worked example for correlations is on page 97.
03
04

05 SUMMARY POINTS

- 06 • Under the fixed-effect model all studies in the analysis share a common true
07 effect.
- 08 • The summary effect is our estimate of this common effect size, and the null
09 hypothesis is that this common effect is zero (for a difference) or one (for a
10 ratio).
- 11 • All observed dispersion reflects sampling error, and study weights are
12 assigned with the goal of minimizing this within-study error.
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

01
02
03
04
05
06
07
08
09
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

Random-Effects Model

01	
02	
03	
04	
05	
06	
07	
08	Introduction
09	The true effect sizes
10	Impact of sampling error
11	Performing a random-effects meta-analysis
12	
13	

INTRODUCTION

In this chapter we introduce the random-effects model. We discuss the assumptions of this model, and show how these are reflected in the formulas used to compute a summary effect, and in the meaning of the summary effect.

THE TRUE EFFECT SIZES

The fixed-effect model, discussed above, starts with the assumption that the true effect size is the same in all studies. However, in many systematic reviews this assumption is implausible. When we decide to incorporate a group of studies in a meta-analysis, we assume that the studies have enough in common that it makes sense to synthesize the information, but there is generally no reason to assume that they are *identical* in the sense that the true effect size is *exactly the same* in all the studies.

For example, suppose that we are working with studies that compare the proportion of patients developing a disease in two groups (vaccinated versus placebo). If the treatment works we would expect the effect size (say, the risk ratio) to be *similar but not identical* across studies. The effect size might be higher (or lower) when the participants are older, or more educated, or healthier than others, or when a more intensive variant of an intervention is used, and so on. Because studies will differ in the mixes of participants and in the implementations of interventions, among other reasons, there may be *different effect sizes* underlying different studies.

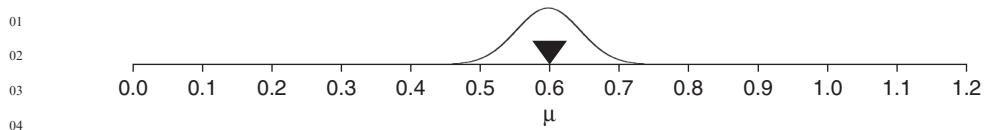


Figure 12.1 Random-effects model – distribution of true effects.

Or, suppose that we are working with studies that assess the impact of an educational intervention. The magnitude of the impact might vary depending on the other resources available to the children, the class size, the age, and other factors, which are likely to vary from study to study.

We might not have assessed these covariates in each study. Indeed, we might not even know what covariates actually are related to the size of the effect. Nevertheless, logic dictates that such factors do exist and will lead to variations in the magnitude of the effect.

One way to address this variation across studies is to perform a *random-effects* meta-analysis. In a random-effects meta-analysis we usually assume that the true effects are normally distributed. For example, in Figure 12.1 the mean of all true effect sizes is 0.60 but the individual effect sizes are distributed about this mean, as indicated by the normal curve. The width of the curve suggests that most of the true effects fall in the range of 0.50 to 0.70.

IMPACT OF SAMPLING ERROR

Suppose that our meta-analysis includes three studies drawn from the distribution of studies depicted by the normal curve, and that the true effects (denoted θ_1 , θ_2 , and θ_3) in these studies happen to be 0.50, 0.55 and 0.65 (see Figure 12.2).

If each study had an infinite sample size the sampling error would be zero and the observed effect for each study would be the same as the true effect for that study.

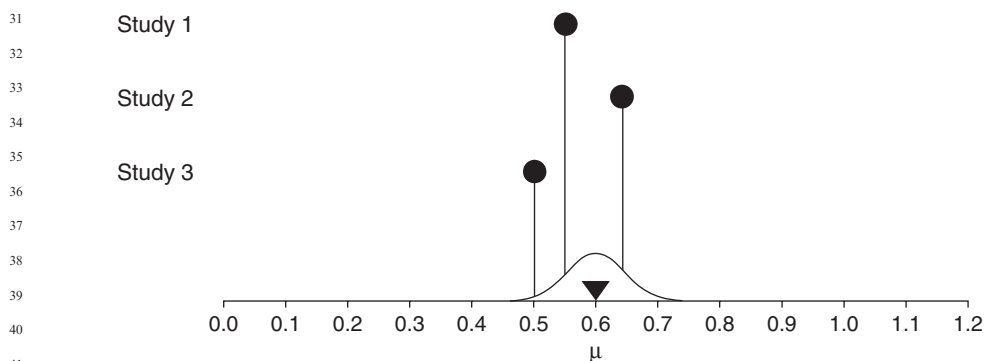


Figure 12.2 Random-effects model – true effects.

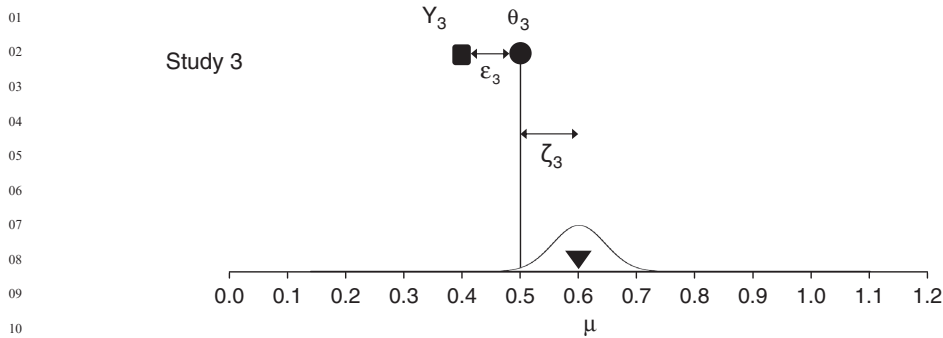


Figure 12.3 Random-effects model – true and observed effect in one study.

If we were to plot the observed effects rather than the true effects, the observed effects would exactly coincide with the true effects.

Of course, the sample size in any study is not infinite and therefore the sampling error is not zero. If the true effect size for a study is θ_i , then the observed effect for that study will be less than or greater than θ_i because of sampling error. For example, consider Study 3 in Figure 12.2. This study is the subject of Figure 12.3, where we consider the factors that control the observed effect. The true effect for Study 3 is 0.50 but the sampling error for this study is -0.10 , and the observed effect for this study is 0.40.

This figure also highlights the fact that the distance between the overall mean and the observed effect in any given study consists of two distinct parts: true variation in effect sizes (ζ_i) and sampling error (ε_i). In Study 3 the total distance from μ to Y_3 is -0.20 . The distance from μ to θ_3 (0.60 to 0.50) reflects the fact that the true effect size actually varies from one study to the next, while the distance from θ_3 to Y_3 (0.5 to 0.4) is sampling error.

More generally, the observed effect Y_i for any study is given by the grand mean, the deviation of the study's true effect from the grand mean, and the deviation of the study's observed effect from the study's true effect. That is,

$$Y_i = \mu + \zeta_i + \varepsilon_i. \quad (12.1)$$

Therefore, to predict how far the observed effect Y_i is likely to fall from μ in any given study we need to consider both the variance of ζ_i and the variance of ε_i .

The distance from μ (the triangle) to each θ_i (the circles) depends on the standard deviation of the distribution of the true effects across studies, called τ (tau) (or τ^2 for its variance). The same value of τ^2 applies to all studies in the meta-analysis, and in Figure 12.4 is represented by the normal curve at the bottom, which extends roughly from 0.50 to 0.70.

The distance from θ_i to Y_i depends on the sampling distribution of the sample effects about θ_i . This depends on the variance of the observed effect size from each study, V_{Y_i} , and so will vary from one study to the next. In Figure 12.4 the curve for Study 1 is relatively wide while the curve for Study 2 is relatively narrow.

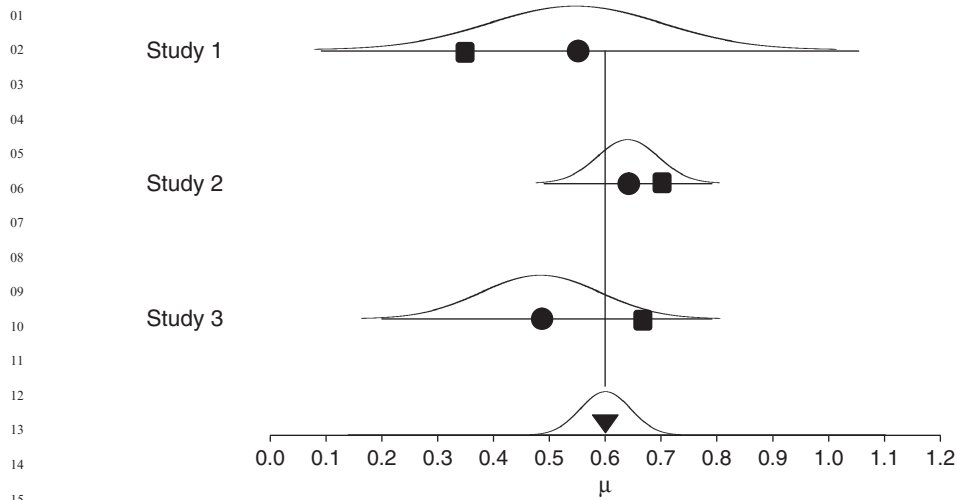


Figure 12.4 Random-effects model – between-study and within-study variance.

PERFORMING A RANDOM-EFFECTS META-ANALYSIS

In an actual meta-analysis, of course, rather than start with the population effect and make projections about the observed effects, we start with the observed effects and try to estimate the population effect. In other words our goal is to use the collection of Y_i to estimate the overall mean, μ . In order to obtain the most precise estimate of the overall mean (to minimize the variance) we compute a weighted mean, where the weight assigned to each study is the inverse of that study's variance.

To compute a study's variance under the random-effects model, we need to know both the within-study variance and τ^2 , since the study's total variance is the sum of these two values. Formulas for computing the within-study variance were presented in Part 3. A method for estimating the between-studies variance is given here so that we can proceed with the worked example, but a full discussion of this method is deferred to Part 4, where we shall pursue the issue of heterogeneity in some detail.

Estimating tau-squared

The parameter τ^2 (tau-squared) is the between-studies variance (the variance of the effect size parameters across the population of studies). In other words, if we somehow knew the *true* effect size for each study, and computed the variance of these effects sizes (across an infinite number of studies), this variance would be τ^2 . One method for estimating τ^2 is the method of moments (or the DerSimonian and Laird) method, as follows. We compute

$$T^2 = \frac{Q - df}{C}, \quad (12.2)$$

where

$$Q = \sum_{i=1}^k W_i Y_i^2 - \frac{\left(\sum_{i=1}^k W_i Y_i \right)^2}{\sum_{i=1}^k W_i}, \quad (12.3)$$

$$df = k - 1, \quad (12.4)$$

where k is the number of studies, and

$$C = \sum W_i - \frac{\sum W_i^2}{\sum W_i}. \quad (12.5)$$

Estimating the mean effect size

In the fixed-effect analysis each study was weighted by the inverse of its variance. In the random-effects analysis, too, each study will be weighted by the inverse of its variance. The difference is that the variance now includes the original (within-studies) variance plus the estimate of the between-studies variance, T^2 . In keeping with the book's convention, we use τ^2 to refer to the parameter and T^2 to refer to the sample estimate of that parameter.

To highlight the parallel between the formulas here (random effects) and those in the previous chapter (fixed effect) we use the same notations but add an asterisk (*) to represent the random-effects version. Under the random-effects model the weight assigned to each study is

$$W_i^* = \frac{1}{V_{Y_i}^*} \quad (12.6)$$

where $V_{Y_i}^*$ is the within-study variance for study i plus the between-studies variance, T^2 . That is,

$$V_{Y_i}^* = V_{Y_i} + T^2.$$

The weighted mean, M^* , is then computed as

$$M^* = \frac{\sum_{i=1}^k W_i^* Y_i}{\sum_{i=1}^k W_i^*} \quad (12.7)$$

that is, the sum of the products (effect size multiplied by weight) divided by the sum of the weights.

The variance of the summary effect is estimated as the reciprocal of the sum of the weights, or

$$V_{M^*} = \frac{1}{\sum_{i=1}^k W_i^*} \quad (12.8)$$

and the estimated standard error of the summary effect is then the square root of the variance,

$$SE_{M^*} = \sqrt{V_{M^*}}. \quad (12.9)$$

The 95% lower and upper limits for the summary effect would be computed as

$$LL_{M^*} = M^* - 1.96 \times SE_{M^*}, \quad (12.10)$$

and

$$UL_{M^*} = M^* + 1.96 \times SE_{M^*}. \quad (12.11)$$

Finally, a Z -value to test the null hypothesis that the mean effect μ is zero could be computed using

$$Z^* = \frac{M^*}{SE_{M^*}}. \quad (12.12)$$

For a one-tailed test the p -value is given by

$$p^* = 1 - \Phi(\pm|Z^*|), \quad (12.13)$$

where we choose '+' if the difference is in the expected direction or '-' otherwise, and for a two-tailed test by

$$p^* = 2[1 - (\Phi(|Z^*|))], \quad (12.14)$$

where $\Phi(Z^*)$ is the standard normal cumulative distribution. This function is tabled in many introductory statistics books, and is implemented in Excel as the function =NORMSDIST(Z^*).

Illustrative example

As before, we suggest that you turn to one of the worked examples in the next chapter before proceeding with this discussion.

SUMMARY POINTS

- Under the random-effects model, the true effects in the studies are assumed to have been sampled from a distribution of true effects.
- The summary effect is our estimate of the mean of all relevant true effects, and the null hypothesis is that the mean of these effects is 0.0 (equivalent to a ratio of 1.0 for ratio measures).

- Since our goal is to estimate the mean of the distribution, we need to take account of two sources of variance. First, there is within-study error in estimating the effect in each study. Second (even if we knew the true mean for each of our studies), there is variation in the true effects across studies. Study weights are assigned with the goal of minimizing both sources of variance.

01
02
03
04
05
06
07
08
09
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

01
02
03
04
05
06
07
08
09
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

Fixed-Effect Versus Random-Effects Models

01	
02	
03	
04	
05	
06	
07	
08	
09	
10	
11	Introduction
12	Definition of a summary effect
13	Estimating the summary effect
14	Extreme effect size in a large study or a small study
15	Confidence interval
16	The null hypothesis
17	Which model should we use?
18	Model should <i>not</i> be based on the test for heterogeneity
19	Concluding remarks
20	
21	
22	
23	
24	
25	
26	
27	
28	
29	
30	
31	
32	
33	
34	
35	
36	
37	
38	
39	
40	
41	
42	
43	

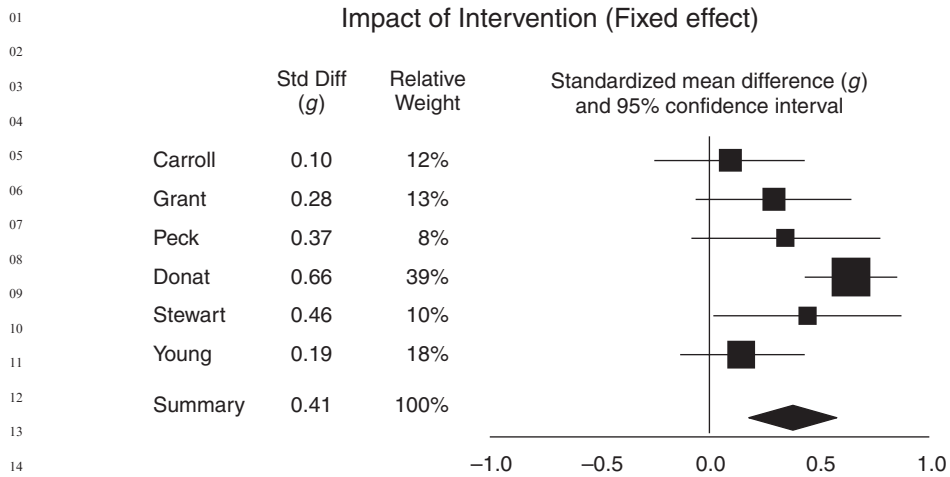
INTRODUCTION

In Chapter 11 and Chapter 12 we introduced the fixed-effect and random-effects models. Here, we highlight the conceptual and practical differences between them.

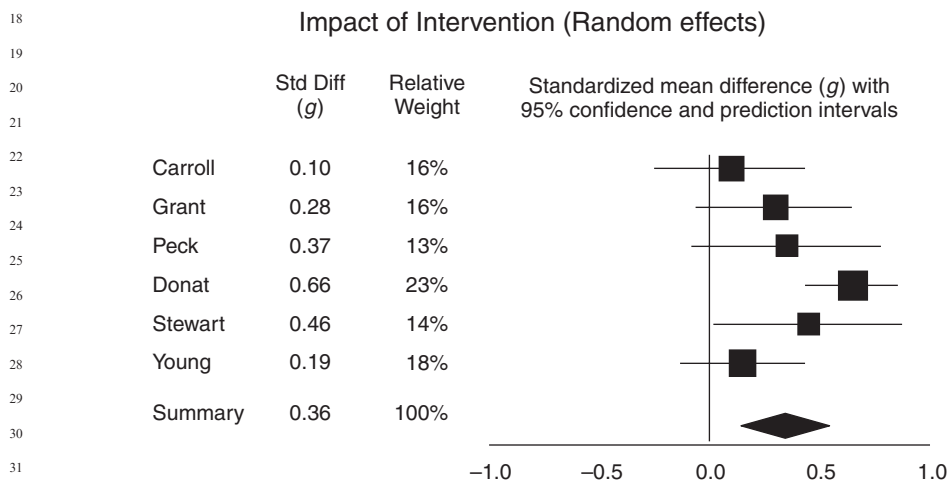
Consider the forest plots in Figures 13.1 and 13.2. They include the same six studies, but the first uses a fixed-effect analysis and the second a random-effects analysis. These plots provide a context for the discussion that follows.

DEFINITION OF A SUMMARY EFFECT

Both plots show a summary effect on the bottom line, but the meaning of this summary effect is different in the two models. In the fixed-effect analysis we assume that the true effect size is the same in all studies, and the summary effect is our estimate of this common effect size. In the random-effects analysis we assume that the true effect size varies from one study to the next, and that the studies in our analysis represent a random sample of effect sizes that could



16 **Figure 13.1** Fixed-effect model – forest plot showing relative weights.



33 **Figure 13.2** Random-effects model – forest plot showing relative weights.

34

35 have been observed. The summary effect is our estimate of the mean of these

36 effects.

37 ESTIMATING THE SUMMARY EFFECT

38

39 Under the fixed-effect model we assume that the true effect size for all studies

40 is identical, and the only reason the effect size varies between studies is

41 sampling error (error in estimating the effect size). Therefore, when assigning

42

43

01 weights to the different studies we can largely ignore the information in the
02 smaller studies since we have better information about the same effect size in
03 the larger studies.

04 By contrast, under the random-effects model the goal is not to estimate one true
05 effect, but to estimate the mean of a distribution of effects. Since each study
06 provides information about a different effect size, we want to be sure that all these
07 effect sizes are represented in the summary estimate. This means that we cannot
08 discount a small study by giving it a very small weight (the way we would in
09 a fixed-effect analysis). The estimate provided by that study may be imprecise, but
10 it is information about an effect that no other study has estimated. By the same
11 logic we cannot give too much weight to a very large study (the way we might in
12 a fixed-effect analysis). Our goal is to estimate the mean effect in a range of
13 studies, and we do not want that overall estimate to be overly influenced by any
14 one of them.

15 In these graphs, the weight assigned to each study is reflected in the size of the
16 box (specifically, the area) for that study. Under the fixed-effect model there is a
17 wide range of weights (as reflected in the size of the boxes) whereas under the
18 random-effects model the weights fall in a relatively narrow range. For example,
19 compare the weight assigned to the largest study (Donat) with that assigned to the
20 smallest study (Peck) under the two models. Under the fixed-effect model Donat is
21 given about five times as much weight as Peck. Under the random-effects model
22 Donat is given only 1.8 times as much weight as Peck.

25 EXTREME EFFECT SIZE IN A LARGE STUDY OR A SMALL STUDY

26 How will the selection of a model influence the overall effect size? In this example
27 Donat is the largest study, and also happens to have the highest effect size. Under
28 the fixed-effect model Donat was assigned a large share (39%) of the total weight
29 and pulled the mean effect up to 0.41. By contrast, under the random-effects model
30 Donat was assigned a relatively modest share of the weight (23%). It therefore had
31 less pull on the mean, which was computed as 0.36.

32 Similarly, Carroll is one of the smaller studies and happens to have the smallest
33 effect size. Under the fixed-effect model Carroll was assigned a relatively small
34 proportion of the total weight (12%), and had little influence on the summary effect.
35 By contrast, under the random-effects model Carroll carried a somewhat higher
36 proportion of the total weight (16%) and was able to pull the weighted mean toward
37 the left.

38 The operating premise, as illustrated in these examples, is that whenever τ^2 is
39 nonzero, the relative weights assigned under random effects will be *more balanced*
40 than those assigned under fixed effects. As we move from fixed effect to random
41 effects, extreme studies will lose influence if they are large, and will gain influence
42 if they are small.

CONFIDENCE INTERVAL

Under the fixed-effect model the only source of uncertainty is the within-study (sampling or estimation) error. Under the random-effects model there is this same source of uncertainty plus an additional source (between-studies variance). It follows that the variance, standard error, and confidence interval for the summary effect will always be larger (or wider) under the random-effects model than under the fixed-effect model (unless T^2 is zero, in which case the two models are the same). In this example, the standard error is 0.064 for the fixed-effect model, and 0.105 for the random-effects model.

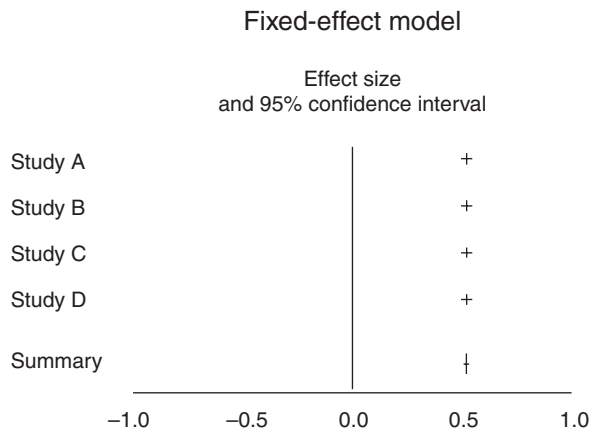


Figure 13.3 Very large studies under fixed-effect model.

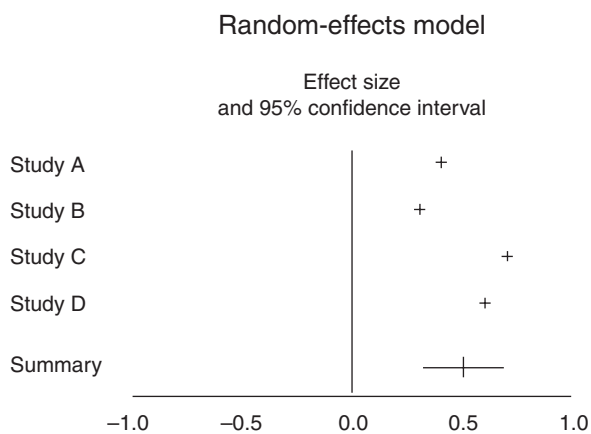


Figure 13.4 Very large studies under random-effects model.

01 Consider what would happen if we had five studies, and each study had an
 02 infinitely large sample size. Under either model the confidence interval for the
 03 effect size in each study would have a width approaching zero, since we know
 04 the effect size in that study with perfect precision. Under the fixed-effect
 05 model the summary effect would also have a confidence interval with a width
 06 of zero, since we know the common effect precisely (Figure 13.3). By con-
 07 trast, under the random-effects model the width of the confidence interval
 08 would not approach zero (Figure 13.4). While we know the effect in each
 09 study precisely, these effects have been sampled from a universe of possible
 10 effect sizes, and provide only an estimate of the mean effect. Just as the error
 11 within a study will approach zero only as the sample size approaches infinity,
 12 so too the error of these studies as an estimate of the mean effect will
 13 approach zero only as the number of studies approaches infinity.

14 More generally, it is instructive to consider what factors influence the standard
 15 error of the summary effect under the two models. The following formulas are
 16 based on a meta-analysis of means from k one-group studies, but the conceptual
 17 argument applies to all meta-analyses. The within-study variance of each mean
 18 depends on the standard deviation (denoted σ) of participants' scores and the
 19 sample size of each study (n). For simplicity we assume that all of the studies
 20 have the same sample size and the same standard deviation (see Box 13.1 for
 21 details).

22 Under the fixed-effect model the standard error of the summary effect is given by

$$23 \quad SE_M = \sqrt{\frac{\sigma^2}{k \times n}}. \quad (13.1)$$

26 It follows that with a large enough sample size the standard error will approach zero,
 27 and this is true whether the sample size is concentrated on one or two studies, or
 28 dispersed across any number of studies.

29 Under the random-effects model the standard error of the summary effect is
 30 given by

$$31 \quad SE_M = \sqrt{\frac{\sigma^2}{k \times n} + \frac{\tau^2}{k}}. \quad (13.2)$$

34 The first term is identical to that for the fixed-effect model and, again, with a
 35 large enough sample size, this term will approach zero. By contrast, the second
 36 term (which reflects the between-studies variance) will only approach zero as the
 37 number of studies approaches infinity. These formulas do not apply exactly in
 38 practice, but the conceptual argument does. Namely, increasing the sample size
 39 within studies is not sufficient to reduce the standard error beyond a certain point
 40 (where that point is determined by τ^2 and k). If there is only a small number of
 41 studies, then the standard error could still be substantial even if the total n is in the
 42 tens of thousands or higher.
 43

BOX 13.1 FACTORS THAT INFLUENCE THE STANDARD ERROR OF THE SUMMARY EFFECT.

To illustrate the concepts with some simple formulas, let us consider a meta-analysis of studies with the very simplest design, such that each study comprises a single sample of n observations with standard deviation σ . We combine estimates of the mean in a meta-analysis. The variance of each estimate is

$$V_{Y_i} = \frac{\sigma^2}{n}$$

so the (inverse-variance) weight in a fixed-effect meta-analysis is

$$W_i = \frac{1}{\sigma^2/n} = \frac{n}{\sigma^2}$$

and the variance of the summary effect under the fixed-effect model the standard error is given by

$$V_M = \frac{1}{\sum_{i=1}^k W_i} = \frac{1}{k \times n/\sigma^2} = \frac{\sigma^2}{k \times n}.$$

Therefore under the fixed-effect model the (true) standard error of the summary mean is given by

$$SE_M = \sqrt{\frac{\sigma^2}{k \times n}}.$$

Under the random-effects model the weight awarded to each study is

$$W_i^* = \frac{1}{(\sigma^2/n) + \tau^2}$$

and the (true) standard error of the summary mean turns out to be

$$SE_{M^*} = \sqrt{\frac{\sigma^2}{k \times n} + \frac{\tau^2}{k}}.$$

THE NULL HYPOTHESIS

Often, after computing a summary effect, researchers perform a test of the null hypothesis. Under the fixed-effect model the null hypothesis being tested is that there is zero effect in *every study*. Under the random-effects model the null hypothesis being tested is that the *mean effect* is zero. Although some may treat these hypotheses as interchangeable, they are in fact different, and it is imperative to choose the test that is appropriate to the inference a researcher wishes to make.

WHICH MODEL SHOULD WE USE?

The selection of a computational model should be based on our expectation about whether or not the studies share a common effect size and on our goals in performing the analysis.

Fixed effect

It makes sense to use the fixed-effect model if two conditions are met. First, we believe that all the studies included in the analysis are functionally identical. Second, our goal is to compute the common effect size for the identified population, and not to generalize to other populations.

For example, suppose that a pharmaceutical company will use a thousand patients to compare a drug versus placebo. Because the staff can work with only 100 patients at a time, the company will run a series of ten trials with 100 patients in each. The studies are identical in the sense that any variables which can have an impact on the outcome are the same across the ten studies. Specifically, the studies draw patients from a common pool, using the same researchers, dose, measure, and so on (we assume that there is no concern about practice effects for the researchers, nor for the different starting times of the various cohorts). All the studies are expected to share a common effect and so the first condition is met. The goal of the analysis is to see if the drug works in the population from which the patients were drawn (and not to extrapolate to other populations), and so the second condition is met, as well.

In this example the fixed-effect model is a plausible fit for the data and meets the goal of the researchers. It should be clear, however, that this situation is relatively rare. The vast majority of cases will more closely resemble those discussed immediately below.

Random effects

By contrast, when the researcher is accumulating data from a series of studies that had been performed by researchers operating independently, it would be unlikely that all the studies were functionally equivalent. Typically, the subjects or interventions in these studies would have differed in ways that would have impacted on

01 the results, and therefore we should not assume a common effect size. Therefore, in
02 these cases the random-effects model is more easily justified than the fixed-effect
03 model.

04 Additionally, the goal of this analysis is usually to generalize to a range of
05 scenarios. Therefore, if one did make the argument that all the studies used an
06 identical, narrowly defined population, then it would not be possible to extrapolate
07 from this population to others, and the utility of the analysis would be severely limited.

08 09 **A caveat**

10 There is one caveat to the above. If the number of studies is very small, then the
11 estimate of the between-studies variance (τ^2) will have poor precision. While the
12 random-effects model is still the appropriate model, we lack the information needed
13 to apply it correctly. In this case the reviewer may choose among several options,
14 each of them problematic.

15 One option is to report the separate effects and *not* report a summary effect.
16 The hope is that the reader will understand that we cannot draw conclusions
17 about the effect size and its confidence interval. The problem is that some readers
18 will revert to vote counting (see Chapter 28) and possibly reach an erroneous
19 conclusion.

20 Another option is to perform a fixed-effect analysis. This approach would yield a
21 descriptive analysis of the included studies, but would not allow us to make
22 inferences about a wider population. The problem with this approach is that (a)
23 we do want to make inferences about a wider population and (b) readers will make
24 these inferences even if they are not warranted.

25 A third option is to take a Bayesian approach, where the estimate of τ^2 is based on
26 data from outside of the current set of studies. This is probably the best option, but the
27 problem is that relatively few researchers have expertise in Bayesian meta-analysis.
28 Additionally, some researchers have a philosophical objection to this approach.

29 For a more general discussion of this issue see *When does it make sense to*
30 *perform a meta-analysis* in Chapter 40.

31 32 **MODEL SHOULD NOT BE BASED ON THE TEST FOR HETEROGENEITY**

33
34 In the next chapter we will introduce a test of the null hypothesis that the between-
35 studies variance is zero. This test is based on the amount of between-studies
36 variance observed, relative to the amount we would expect if the studies actually
37 shared a common effect size.

38 Some have adopted the practice of starting with a fixed-effect model and then
39 switching to a random-effects model if the test of homogeneity is statistically
40 significant. This practice should be strongly discouraged because the decision to
41 use the random-effects model should be based on our understanding of whether or
42 not all studies share a common effect size, and not on the outcome of a statistical test
43 (especially since the test for heterogeneity often suffers from low power).

01 If the study effect sizes are seen as having been sampled from a *distribution* of
02 effect sizes, then the random-effects model, which reflects this idea, is the logical one
03 to use. If the between-studies variance is substantial (and statistically significant) then
04 the fixed-effect model is inappropriate. However, even if the between-studies var-
05 iance does not meet the criterion for statistical significance (which may be due simply
06 to low power) we should still take account of this variance when assigning weights. If
07 T^2 turns out to be zero, then the random-effects analysis reduces to the fixed-effect
08 analysis, and so there is no cost to using this model.

09 On the other hand, if one has elected to use the fixed-effect model *a priori* but the
10 test of homogeneity is statistically significant, then it would be important to revisit
11 the assumptions that led to the selection of a fixed-effect model.

12 13 **CONCLUDING REMARKS**

14
15 Our discussion of differences between the fixed-model and the random-effects
16 model focused largely on the computation of a summary effect and the confidence
17 intervals for the summary effect. We did not address the implications of the
18 dispersion itself. Under the fixed-effect model we assume that all dispersion in
19 observed effects is due to sampling error, but under the random-effects model we
20 allow that some of that dispersion reflects real differences in effect size across
21 studies. In the chapters that follow we discuss methods to quantify that dispersion
22 and to consider its substantive implications.

23 Although throughout this book we define a fixed-effect meta-analysis as assum-
24 ing that every study has a common true effect size, some have argued that the fixed-
25 effect method is valid without making this assumption. The point estimate of the
26 effect in a fixed-effect meta-analysis is simply a weighted average and does not
27 strictly require the assumption that all studies estimate the same thing. For simpli-
28 city and clarity we adopt a definition of a fixed-effect meta-analysis that does
29 assume homogeneity of effect.

30 31 32 **SUMMARY POINTS**

- 33 • A fixed-effect meta-analysis estimates a single effect that is assumed to be
34 common to every study, while a random-effects meta-analysis estimates the
35 mean of a distribution of effects.
- 36 • Study weights are more balanced under the random-effects model than under the
37 fixed-effect model. Large studies are assigned less relative weight and small
38 studies are assigned more relative weight as compared with the fixed-effect
39 model.
- 40 • The standard error of the summary effect and (it follows) the confidence
41 intervals for the summary effect are wider under the random-effects model
42 than under the fixed-effect model.
- 43

- The selection of a model must be based solely on the question of which model fits the distribution of effect sizes, and takes account of the relevant source(s) of error. When studies are gathered from the published literature, the random-effects model is generally a more plausible match.
- The strategy of starting with a fixed-effect model and then moving to a random-effects model if the test for heterogeneity is significant is a mistake, and should be strongly discouraged.

01
02
03
04
05
06
07
08
09
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43

Criticisms of Meta-Analysis

Introduction

One number cannot summarize a research field

The file drawer problem invalidates meta-analysis

Mixing apples and oranges

Garbage in, garbage out

Important studies are ignored

Meta-analysis can disagree with randomized trials

Meta-analyses are performed poorly

Is a narrative review better?

Concluding remarks

INTRODUCTION

While meta-analysis has been widely embraced by large segments of the research community, this point of view is not universal and people have voiced numerous criticisms of meta-analysis.

Some of these criticisms are worth mentioning for their creative use of metaphor. The first set of Cochrane reviews dealt with studies in neonatology, and one especially creative critic, cited by Mann (1990), called the reviewers *an obstetrical Baader Meinhof gang* (*obstetrical* being a reference to the field of research, and *Baader Meinhof gang* a reference to the terrorist group that operated in Europe during the 1970s and 1980s).

Others were more circumspect in their comments. Eysenck (1978) criticized a meta-analysis as *an exercise in mega-silliness*. Shapiro (1994) published a paper entitled *Meta-Analysis / Shmeta Analysis*. Feinstein (1995) wrote an editorial in which he referred to meta-analysis as 'statistical alchemy for the 21st century'.

01 These critics share not only an affinity for allegory and alliteration but also a
02 common set of concerns about meta-analysis. In this chapter we address the
03 following criticisms that have been leveled at meta-analysis, as follows.

- 04 • One number cannot summarize a research field
- 05 • The file drawer problem invalidates meta-analysis
- 06 • Mixing apples and oranges
- 07 • Garbage in, garbage out
- 08 • Important studies are ignored
- 09 • Meta-analysis can disagree with randomized trials
- 10 • Meta-analyses are performed poorly

11
12 After considering each of these questions in turn, we ask whether a traditional
13 narrative review fares any better than a systematic review on these criticisms. And,
14 we summarize the legitimate criticisms of meta-analysis that need to be considered
15 whenever meta-analysis is applied.

17 ONE NUMBER CANNOT SUMMARIZE A RESEARCH FIELD

18 Criticism

19
20 A common criticism of meta-analysis is that the analysis focuses on the summary
21 effect, and ignores the fact that the treatment effect may vary from study to study. Bailar
22 (1997), for example, writes, 'Any attempt to reduce results to a single value, with
23 confidence bounds, is likely to lead to conclusions that are wrong, perhaps seriously so.'

25 Response

26
27 In fact, the goal of a meta-analysis should be to *synthesize* the effect sizes, and not
28 simply (or necessarily) to report a summary effect. If the effects are consistent, then
29 the analysis shows that the effect is robust across the range of included studies. If
30 there is modest dispersion, then this dispersion should serve to place the mean effect
31 in context. If there is substantial dispersion, then the focus should shift from the
32 summary effect to the dispersion itself. Researchers who report a summary effect
33 and ignore heterogeneity are indeed missing the point of the synthesis.

36 THE FILE DRAWER PROBLEM INVALIDATES META-ANALYSIS

37 Criticism

38
39 While the meta-analysis will yield a mathematically sound synthesis of the studies
40 included in the analysis, if these studies are a biased sample of all possible studies,
41 then the mean effect reported by the meta-analysis will reflect this bias. Several
42 lines of evidence show that studies finding relatively high treatment effects are
43 more likely to be published than studies finding lower treatment effects. The latter,

01 unpublished, research lies dormant in the researchers' filing cabinets, and has led to
02 the use of the term *file drawer problem* for meta-analysis.

04 **Response**

05
06 Since published studies are more likely to be included in a meta-analysis than their
07 unpublished counterparts, there is a legitimate concern that a meta-analysis may
08 overestimate the true effect size.

09 Chapter 30 (entitled *Publication Bias*) explores this question in some detail. In that
10 chapter we discuss methods to assess the likely amount of bias in any given meta-
11 analysis, and to distinguish between analyses that can be considered robust to the
12 impact of publication bias from those where the results should be considered suspect.

13 We must remember that publication bias is a problem for any kind of literature
14 search. The problem exists for the clinician who searches a database to locate
15 primary studies about the utility of a treatment. It exists for persons performing a
16 narrative review. And, it exists for persons performing a meta-analysis. Publication
17 bias has come to be identified with meta-analysis because meta-analysis has the
18 goal of providing a more accurate synthesis than other methods, and so we are
19 concerned with biases that will interfere with this goal. However, it would be a
20 mistake to conclude that this bias is not a problem for the narrative review. There, it
21 is simply easier to ignore.

23 **MIXING APPLES AND ORANGES**

24 **Criticism**

25
26 A common criticism of meta-analysis is that researchers combine different kinds of
27 studies (*apples and oranges*) in the same analysis. The argument is that the
28 summary effect will ignore possibly important differences across studies.

30 **Response**

31
32 The studies that are brought together in a meta-analysis will inevitably differ in their
33 characteristics, and the difficulty is deciding just how similar they need to be. The
34 decision as to which studies should be included is always a judgment, and people
35 will have different opinions on the appropriateness of combining results across
36 studies. Some meta-analysts may make questionable judgments, and some critics
37 may make unreasonable demands on similarity.

38 We need to remember that meta-analyses almost always, by their very nature,
39 address broader questions than individual studies. Hence a meta-analysis may be
40 thought of as asking a question about fruit, for which both apples and oranges (and
41 indeed pears and melons) contribute valuable information. One of the strengths of
42 meta-analysis is that the consistency, and hence generalizability, of findings from
43 one type of study to the next can be assessed formally.

01 Of course, we always need to remember that we are dealing with different kinds of
02 fruit, and to anticipate that effects may vary from one kind to the other. It is a further
03 strength of meta-analysis that these differences, if identified, can be investigated
04 formally. Assume, for example, that a treatment is very effective for patients with
05 acute symptoms but has no effect for patients with chronic symptoms. If we were to
06 combine data from studies that used both types of patients, and conclude that the
07 treatment was modestly effective (on average), this conclusion would not be accurate
08 for either kind of patient. If we were to restrict our attention to studies in only patients
09 with acute symptoms, or only patients with chronic symptoms, we could report how
10 the treatment worked with one type of patient, but could only speculate about how it
11 would have worked with the other type. By contrast, a meta-analysis that includes
12 data for both types of patients may allow us to address this question empirically.

14 GARBAGE IN, GARBAGE OUT

16 Criticism

17 The often-heard metaphor *garbage in, garbage out* refers to the notion that if a
18 meta-analysis includes many low-quality studies, then fundamental errors in the
19 primary studies will be carried over to the meta-analysis, where the errors may be
20 harder to identify.

22 Response

23
24 Rather than thinking of meta-analysis as a process of *garbage in, garbage out* we
25 can think of it as a process of waste management. A systematic review or meta-
26 analysis will always have a set of inclusion criteria and these should include criteria
27 based on the quality of the study. For trials, we may decide to limit the studies to
28 those that use random assignment, or a placebo control. For observational studies
29 we may decide to limit the studies to those where confounders were adequately
30 addressed in the design or analysis. And so on. In fact, it is common in a systematic
31 review to start with a large pool of studies and end with a much smaller set of studies
32 after all inclusion/exclusion criteria are applied.

33 Nevertheless, the studies that do make it as far as a meta-analysis are unlikely to
34 be perfect, and close attention should be paid to the possibility of bias due to study
35 limitations. A meta-analysis of a collection of studies that is each biased in the same
36 direction will suffer from the same bias and have higher precision. In this case,
37 performing a meta-analysis can indeed be more dangerous than not performing one.

38 However, as noted in the response to the previous criticism about *apples and*
39 *oranges*, a strength of meta-analysis is the ability to investigate whether variation in
40 characteristics of studies is related to the size of the effect. Suppose that ten studies
41 used an acceptable method to randomize patients while another ten used a ques-
42 tionable method. In the analysis we can compare the effect size in these two
43 subgroups, and determine whether or not the effect size actually differs between

01 the two. Note that such analyses (those comparing effects in different subgroups)
02 can have very low power so need to be interpreted carefully, especially when there
03 are not many studies within subgroups.

04 **IMPORTANT STUDIES ARE IGNORED**

05 **Criticism**

06 Whereas the *garbage in, garbage out* problem relates to the inclusion of studies that
07 perhaps should not be included, a common complementary criticism is that important
08 studies were left out. The criticism is often leveled by people who are uncomfortable
09 with the findings of a meta-analysis. For example, a meta-analysis to assess the effects
10 of antioxidant supplements (beta-carotene, vitamin A, vitamin C, vitamin E, and
11 selenium) on overall mortality was met with accusations on the web site of the Linus
12 Pauling Institute (Oregon State University) that in this 'flawed analysis of flawed data'
13 the authors looked at 815 human clinical trials of antioxidant supplements, but only 68
14 were included in the meta-analysis.
15
16

17 **Response**

18 We have explained that systematic reviews and meta-analyses require explicit
19 mechanisms for deciding which studies to include and which ones to exclude.
20 These eligibility criteria are determined by a combination of considerations of
21 relevance and considerations of bias, and are typically decided before the search
22 for studies is implemented. Studies should be sufficiently similar to yield results
23 that can be interpreted, and sufficiently free of bias to yield results that can be
24 believed. For both purposes, judgments are required, and not all meta-analysts or
25 readers would reach the same judgments on each occasion. Importantly, in meta-
26 analysis the criteria are transparent and are described as part of the report.
27
28

29 **META-ANALYSIS CAN DISAGREE WITH RANDOMIZED TRIALS**

30 **Criticism**

31 LeLorier *et al.* (1997) published a paper in which they pointed out that meta-
32 analyses sometimes yield different results than large scale randomized trials.
33 Specifically, they located cases in the medical literature where someone had
34 performed a meta-analysis, and someone else subsequently performed a large
35 scale randomized trial that addressed the same question (e.g. *Does the treatment*
36 *work?*). The authors reported that the results of the meta-analysis and the rando-
37 mized trial *matched* (both were statistically significant, or neither was statistically
38 significant) in about 66% of cases, but did not match (one was statistically sig-
39 nificant but the other was not) in the remaining 34%. Since randomized trials are
40 generally accepted as the gold standard they conclude that some 34% of these meta-
41 analyses were wrong, and that meta-analyses in general cannot be trusted.
42
43

Response

There are both technical and conceptual flaws in this criticism. The technical flaws relate to the question of what we mean by *matching*, and the authors' decision to define *matching* as both studies being (or not being) statistically significant. The discussion that follows draws in part on comments by Ioannidis *et al.* (1998), Lelorier *et al.* (1997, 536–543) and others (see further readings at the end of this chapter).

Consider Figure 43.1, which shows a meta-analysis of five randomized controlled trials (RCTs) at the top, and a subsequent large-scale randomized trial at the bottom.

In this fictional example the five studies in the meta-analysis each showed precisely the same effect, an odds ratio of 0.80. The summary effect in the meta-analysis is (it follows) an odds ratio of 0.80. And, the subsequent study showed the same effect, an odds ratio of 0.80.

The only difference between the summary effect in the meta-analysis and the effect in the subsequent study is that the former is reported with greater precision (since it is based on more data) and therefore yields a *p*-value under 0.05. By the LeLorier criterion these two conclusions would be seen as conflicting, when in fact they have the identical effect size.

Additionally, LeLorier concludes that in the face of this conflict the single randomized trial is correct and the meta-analysis is wrong. In fact, though, it is the meta-analysis, which incorporates data from five randomized trials rather than one, that has the more powerful position. (What would happen if we performed a new meta-analysis which incorporated the most recent randomized trial? Would

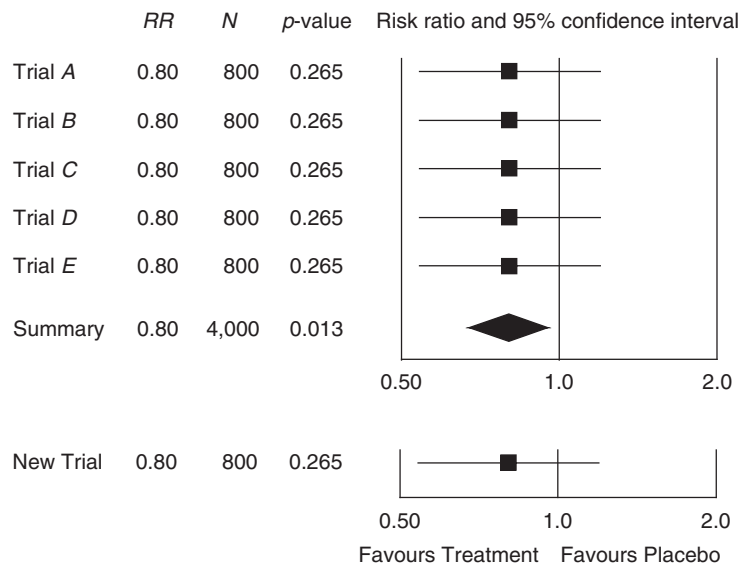


Figure 43.1 Forest plot of five fictional studies and a new trail (consistent effects).

LeLorier now see this new meta-analysis as flawed?) In fact, the real issue is not that a meta-analysis disagrees with a randomized trial, but that randomized trials disagree with each other.

At a meeting of The Cochrane Collaboration in Baltimore (1996), a plenary speaker made the same argument being made by LeLorier *et al.* (that meta-analyses sometimes yield different results than randomized trials) and, like the paper, cited the statistic that roughly a third of meta-analyses fail to match the *comparable* randomized trial. A distinguished member of the audience, Harris Cooper, asked the speaker if he knew what percentage of randomized trials fail to match the next randomized trial on the same topic. It turns out that the percentage is roughly a third.

However, to move on to a more interesting question, let's assume that the results from a meta-analysis and a randomized trial really do differ. Suppose that the meta-analysis yields a risk ratio of 0.67 (with a 95% confidence interval of 0.84 to 0.77) while the new trial yields a risk ratio of 0.91 (0.82 to 1.0). According to the meta-analysis the treatment reduces the risk by at least 23%, while the new trial says that its impact is no more than 18%.

In this case the effect *is different* in the two analyses, but that does not mean that one is wrong and the other is right. Rather, it behooves us to ask why the two results should differ, much as we would if we had two large scale randomized trials with significantly different results. Often, it will turn out that the different analyses either were asking different questions or differed in some important way. A careful examination of the differences in method, patient population, and so on, may help to uncover the source of the difference.

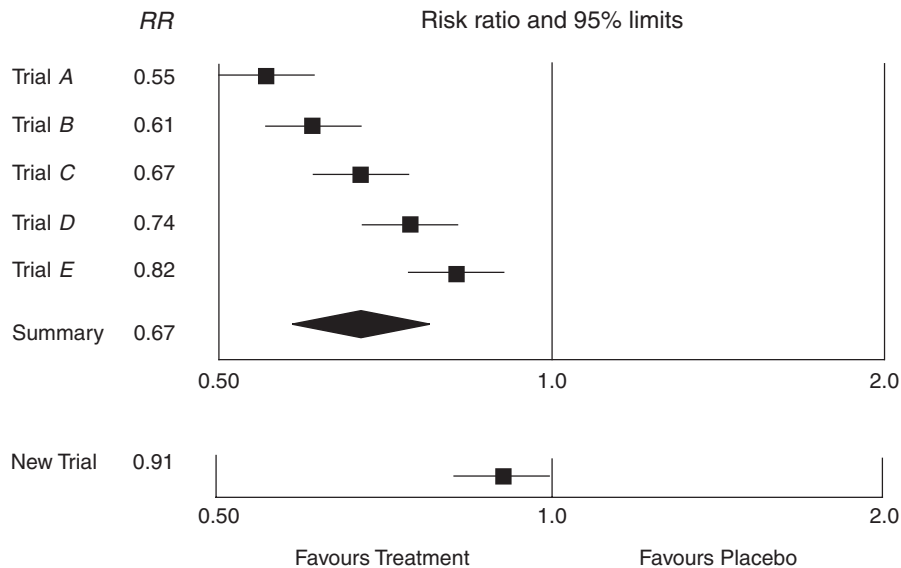


Figure 43.2 Forest plot of five fictional studies and a new trial (heterogeneous effects).

01 Consider the following scenario, depicted in Figure 43.2. A new compound is
02 introduced, which is meant to minimize neurological damage in stroke patients. In
03 1990, the compound is tested in a randomized trial involving patients with a very poor
04 prognosis, and yields a risk ratio of 0.55. Based on these encouraging results, in 1994
05 it is tested in patients with a somewhat better prognosis. Since the patients in this
06 group are more likely to recover without treatment, the impact of the drug is less
07 pronounced, and the risk ratio is 0.61. By 1998 the drug is being tested with all
08 patients, and the risk ratio is 0.82. These are the studies included in the meta-analysis.
09 The new trial is performed using a relatively healthy population and (following the
10 trend seen in the meta-analysis) yields a risk ratio of 0.91.

11 If one were to report a mean effect of 0.67 for the meta-analysis versus 0.91 for the
12 new trial there would indeed be a problem. But, as we have emphasized throughout
13 this volume, the meta-analysis should focus on the dispersion in effects and try to
14 identify the reason for the dispersion. In this example, using either health status or
15 study year as a covariate we can explain the pattern of the effects, and would have
16 predicted that the effect size in the new study would fall where it did.

18 META-ANALYSES ARE PERFORMED POORLY

19 Criticism

20 John C. Bailar, in an editorial for the *New England Journal of Medicine* (Bailar,
21 1997), writes that mistakes such as those outlined in the prior criticisms are common
22 in meta-analysis. He argues that a meta-analysis is inherently so complicated that
23 mistakes by the persons performing the analysis are all but inevitable. He also
24 argues that journal editors are unlikely to uncover all of these mistakes.

27 Response

28 The specific points made by Bailar about problems with meta-analysis are entirely
29 reasonable. He is correct that many meta-analyses contain errors, some of them
30 important ones. His list of potential (and common) problems can serve as a bullet
31 list of mistakes to avoid when performing a meta-analysis.

32 However, the mistakes cited by Bailar are flaws in the application of the
33 method, rather than problems with the method itself. Many primary studies
34 suffer from flaws in the design, analyses, and conclusions. In fact, some
35 serious kinds of problems are endemic in the literature. The response of the
36 research community is to locate these flaws, consider their impact for the
37 study in question, and (hopefully) take steps to avoid similar mistakes in the
38 future. In the case of meta-analysis, as in the case of primary studies, we
39 cannot condemn a method because some people have used that method
40 improperly. As Bob Abelson once remarked in a related context, 'Think of
41 all the things that people abuse. There are college educations. And oboes.'

IS A NARRATIVE REVIEW BETTER?

In his editorial Bailar concludes that, until such time as the quality of meta-analyses is improved, he would prefer to work with the traditional narrative reviews: 'I still prefer conventional narrative reviews of the literature, a type of summary familiar to readers of the countless review articles on important medical issues.'

We disagree with the conclusion that narrative reviews are preferable to systematic reviews, and that meta-analyses should be avoided. The narrative review suffers from every one of the problems cited for the systematic review. The only difference is that, in the narrative review, these problems are less obvious. For example:

- The process of determining which studies to include in the systematic review or meta-analysis is difficult and prone to error. But at least there is a set of criteria for determining which studies to include. If the narrative review also has such criteria, then it is subject to the same kinds of error. If not, then we have no way of knowing how studies are being selected, which only compounds the problem.
- Meta-analyses can be affected by publication bias. But the same biases exist in the material upon which narrative reviews are based. Meta-analysis offers a means to investigate the likelihood of these biases and their potential impact on the results.
- Meta-analyses may be based on low quality primary research. But a good systematic review includes a careful assessment of the included studies with regard to their quality or risk of bias, and meta-analytic methods enable formal examination of the potential impact of these biases. A narrative reviewer may discount a study because of a belief that the results are suspect for some reason. However, a limitation can be found for virtually any study, so in the absence of a systematic quality assessment of every study, a narrative reviewer is free to be suspect about any study's results and to lay the blame on one or more of its limitations.
- The weighting scheme in a meta-analysis may give a lot (or little) weight to specific studies in ways that may appear inappropriate. But in a meta-analysis the weights reflect specific goals (to minimize the variance, or to reflect the range of effects) and the weighting scheme is detailed as part of the report, so a reader is able to agree or disagree with it. By contrast, in the case of a narrative review, the reviewer assigns *weights* to studies based on criteria that he or she does not communicate, and may not even be able to fully articulate. Here, the problem involves not only the relative weights assigned to small or large studies. It extends also to the propensity of one reviewer to focus on effect sizes, and of another to focus on (and possibly be misled by) significance tests.
- Some meta-analyses focus on the summary effect and ignore the pattern of dispersion in the results. To ignore the dispersion is clearly a mistake both in a narrative review and in a meta-analysis. However, meta-analysis provides a full complement of tools to assess the pattern of dispersion, and possibly to explain it as a function of study-level covariates. By contrast, it would be an almost

01 impossible task for a narrative reviewer to accurately assess the pattern of
02 dispersion, or to understand its relationship to other variables.

- 03 • In support of the narrative review, Bailer cites the role of the expert with
04 substantive knowledge of the field, who can identify flaws in specific studies,
05 or the presence of potentially important moderator variables. However, this is not
06 an advantage of the narrative review, since the expert is expected to play the same
07 role in a meta-analysis. Steve Goodman (1991) wrote, 'The best meta-analyses
08 knit clinical insight with quantitative results in a way that enhances both. They
09 should combine the careful thought and synthesis of a good review with the
10 scientific rigor of a good experiment.'

12 CONCLUDING REMARKS

14 Most of the criticisms raised in this chapter point to problems with meta-analysis,
15 and make the implicit argument that the problem would go away if we dispensed
16 with the meta-analysis and performed a narrative review. We have argued that these
17 problems exist also for the narrative review, and that the key advantage of the
18 systematic approach of a meta-analysis is that all steps are clearly described so that
19 the process is transparent.

20 Is meta-analysis so difficult that the method should be abandoned, as some have
21 suggested? Our answer is obviously that it is not. Most of the criticisms raised deal
22 with the application of the method, rather than with the method itself. What we
23 should do is take the valid criticisms seriously and protect against them in planned
24 analyses and by thoughtful interpretation of results.

25 Steven Goodman, in his editorial for *Annals of Internal Medicine* (1991) writes,

26 Regardless of the summary number, meta-analysis should shed light on why trial
27 results differ; raise research and editorial standards by calling attention to the
28 strengths and weaknesses of the body of research in an area; and give the practitioner
29 an objective view of the research literature, unaffected by the sometimes distorting
30 lens of individual experience and personal preference that can affect a less structured
31 review.

35 SUMMARY POINTS

- 36 • Meta-analyses are sometimes criticized for a number of flaws, and critics
37 have argued that narrative reviews provide a better solution.
- 38 • Some of these flaws, such as the idea that we cannot summarize a body of data
39 in a single number, are based on misunderstandings of meta-analysis.
- 40 • Many of the flaws (such as ignoring dispersion in effect sizes) reflect pro-
41 blems in the way that meta-analysis is used, rather than problems in the
42 method itself.
- 43

- Other flaws (such as publication bias) are a problem for meta-analysis. However, the suggestion that these problems do not exist in narrative reviews is wrong. These problems exist for narrative reviews as well, but are simply easier to ignore since those reviews lack a clear structure.

Further Reading

- Bailar, J.C. (1995). The practice of meta-analysis. *J Clin Epidemiol* 48: 149–157.
- Bailar, J.C. (1997). The promise and problems of meta-analysis. *New Engl J Med* 337: 559–561.
- Boden, W.E. (1992). Meta-analysis in clinical trials reporting: has a tool become a weapon? *Am J Cardiol* 69: 681–686.
- Egger, M., & Davey Smith, G. (1998). Bias in location and selection of studies. *BMJ* 316: 61–66.
- Eysenck, H.J. (1978). An exercise in mega-silliness. *Am Psychol* 33: 517.
- Lau, J., Ioannidis, J.P., Terrin, N., Schmid, C.H., & Olkin, I. (2006). The case of the misleading funnel plot. *BMJ* 333: 597–600.
- LeLorier, J., Gregoire, G., Benhaddad, A., Lapierre, J., & Derderian, F. (1997). Discrepancies between meta-analyses and subsequent large randomized, controlled trials. *N Engl J Med* 337: 536–543.
- Responses to Lelorier *et al.*
- Bent, S., Kerlikowske, K., & Grady, D. (1998). *NEJM*, 338(1), 60.
 - Imperiale, T.F. (1998). *NEJM*, 338(1), 61.
 - Ioannidis, J.P., Cappelleri, J.C., & Lau, J. (1998). *NEJM*, 338(1), 59.
 - Khan, S., Williamson, P., & Sutton, R. (1998). *NEJM*, 338(1), 60–61
 - LeLorier, J., & Gregoire, G. (1998). *NEJM*, 338(1), 61–62.
 - Song, F. J., & Sheldon, T. A. (1998). *NEJM*, 338(1), 60.
 - Stewart, L. A., Parmar, M. K., & Tierney, J. F. (1998). *NEJM*, 338(1), 61
- Sharpe, D. (1997) Of apples and oranges, file drawers and garbage: why validity issues in meta-analysis will not go away. *Clin Psychol Rev* 17: 881–901.
- Thompson, S.G & Pocock, S. J. (1991). Can meta-analysis be trusted? *Lancet* 338: 1127–1130.