

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/286122828>

Quasi-Experimental Research Designs

Article · February 2012

DOI: 10.1093/acprof:oso/9780195387384.001.0001

CITATIONS
30

READS
21,730

1 author:



Bruce A. Thyer
Florida State University

390 PUBLICATIONS 4,477 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



I am working on a dozen things. Too many to list. [View project](#)



Autism [View project](#)

Quasi-Experimental Research Designs

POCKET GUIDES TO SOCIAL WORK RESEARCH METHODS

Series Editor

Tony Tripodi, DSW

Professor Emeritus, Ohio State University

*Determining Sample Size:
Balancing Power, Precision, and Practicality*

Patrick Dattalo

Preparing Research Articles

Bruce A. Thyer

Systematic Reviews and Meta-Analysis

Julia H. Littell, Jacqueline Corcoran,
and Vijayan Pillai

Historical Research

Elizabeth Ann Danto

Confirmatory Factor Analysis

Donna Harrington

*Randomized Controlled Trials:
Design and Implementation for
Community-Based Psychosocial*

Interventions

Phyllis Solomon, Mary M. Cavanaugh,
and Jeffrey Drain

Needs Assessment

David Royle, Michele Staton-Tindall,
Karen Badger, and
J. Matthew Webster

Multiple Regression with Discrete

Dependent Variables

John G. Orme and Terri
Combs-Orme

Developing Cross-Cultural Measurement

Thanh V. Tran

*Intervention Research:
Developing Social Programs*

Mark W. Fraser, Jack M. Richman,
Maeda J. Galinsky, and Steven H. Day

*Developing and Validating Rapid
Assessment Instruments*

Neil Abell, David W. Springer, and
Akihito Kamata

*Clinical Data-Mining:
Integrating Practice and Research*

Irwin Epstein

*Strategies to Approximate Random
Sampling and Assignment*

Patrick Dattalo

Analyzing Single System Design Data

William R. Nugent

Survival Analysis

Shenyang Guo

*The Dissertation:
From Beginning to End*

Peter Lyons and Howard J. Doueck

Cross-Cultural Research

Jorge Delva, Paula Allen-Meares, and
Sandra L. Momper

Secondary Data Analysis

Thomas P. Vartanian

Narrative Inquiry

Kathleen Wells

Structural Equation Modeling

Natasha K. Bowen and Shenyang Guo

*Finding and Evaluating Evidence:
Systematic Reviews and
Evidence-Based Practice*

Denise E. Bronson and
Tamara S. Davis

*Policy Creation and Evaluation:
Understanding Welfare Reform in
the United States*

Richard Hofer

Grounded Theory

Julianne S. Oktay

*Systematic Synthesis of
Qualitative Research*

Michael Saini and Aron Shlonsky

Quasi-Experimental Research Designs

Bruce A. Thyer

BRUCE A. THYER

Quasi-Experimental Research Designs

OXFORD
UNIVERSITY PRESS

OXFORD
UNIVERSITY PRESS

Oxford University Press, Inc., publishes works that further
Oxford University's objective of excellence
in research, scholarship, and education.

Oxford New York
Auckland Cape Town Dar es Salaam Hong Kong Karachi
Kuala Lumpur Madrid Melbourne Mexico City Nairobi
New Delhi Shanghai Taipei Toronto

With offices in
Argentina Austria Brazil Chile Czech Republic France Greece
Guatemala Hungary Italy Japan Poland Portugal Singapore
South Korea Switzerland Thailand Turkey Ukraine Vietnam

Copyright © 2012 by Oxford University Press, Inc.

Published by Oxford University Press, Inc.
198 Madison Avenue, New York, New York 10016
www.oup.com

Oxford is a registered trade mark of Oxford University Press

All rights reserved. No part of this publication may be reproduced,
stored in a retrieval system, or transmitted, in any form or by any means,
electronic, mechanical, photocopying, recording, or otherwise,
without the prior permission of Oxford University Press.

Library of Congress Cataloging-in-Publication Data

Thyer, Bruce A.
Quasi-experimental research designs / Bruce A. Thyer.
p. cm. — (Pocket guides to social work research methods)
Includes bibliographical references and index.
ISBN 978-0-19-538738-4 (pbk. : alk. paper)
1. Social sciences—Research—Methodology.
2. Research—Methodology. I. Title.
H62.T49 2012

001.4'34—dc23 2011036983

1 3 5 7 9 8 6 4 2

Printed in the United States of America
on acid-free paper

Contents

Preface vii

1 The Role of Group Research Designs to
Evaluate Social Work Practice 3

2 Pre-experimental Group Research Designs 29

3 Quasi-Experimental Group Research Designs 77

4 Interrupted Time Series Research Designs 107

5 Evaluating and Reporting Quasi-Experimental Studies 127

Glossary 179

References 187

Index 201



Preface

It has been argued that research which focuses on evaluating the outcomes of social work programs, on those services provided to individuals, groups, couples, families, organizations, and communities, or aimed at empirically evaluating the effects of various social welfare programs, are among the most scientifically valuable contributions that can be made by our discipline's researchers. Regardless of whether or not one subscribes to this appraisal, it is clear that the formal research method that is most often used in the evaluation of social work practice is the general type called *quasi-experimental design*. Practice outcomes can, of course, be evaluated using a diversity of approaches, including clinical judgment, narrative case histories, single-case studies, randomized controlled trials, meta-analyses, and systematic reviews. Each of these has its strengths and limitations. However, the adaptability of quasi-experimental designs for use in field settings—those naturalistic environments in which real social work services are provided to real clients, under clinically representative conditions—renders them particularly suitable for use by social work researchers.

In this work, I review the history and background of quasi-experimental designs as used by social workers, and I walk the reader through an increasingly complex array of these designs. I begin with studies of the outcomes obtained by a single group of clients, studies that

are sometimes collectively labeled as pre-experimental designs. A variety of these designs are described, with their strengths, limitations, and practical uses. I next move to a discussion of designs involving evaluating the outcomes of two or more groups of clients, with one group receiving an intervention that is the focus of investigation, and the other(s) receiving either no treatment, standard care or treatment-as-usual, and/or a group receiving an innocuous intervention that serves as a control for nonspecific placebo influences (which are ubiquitous in the human services, sometimes powerfully and positively so). The final category of designs presented are various time-series designs, most often used in policy evaluation studies. Each design is illustrated with a description of its application in one or more previously published articles authored by social workers. The concluding chapter addresses how the data from quasi-experimental designs can be statistically evaluated, reviews some ethical standards and guidelines relating to the protection of human subjects from risk, and describes some contemporary standards that are recommended to be followed when reporting the results of quasi-experimental investigations.

There are, of course, other texts available that cover the methodology of the design and conduct of quasi-experimental studies, but none do so from the particular perspective of the profession of social work. Such a disciplinary focus, I hope, represents the unique contribution of the present volume. I would like to express my gratitude to the four reviewers of this book in manuscript form, and for their many(!) helpful suggestions, which greatly strengthened it. I would also like to thank my collaborators at Oxford University Press, Maura Roessner and Nicholas Liu, whose patience was extensively tested as this volume developed. My spouse, Dr. Laura Myers, as always, was a recurring source of support and inspiration. At one point in the past year, she was busy co-authoring her own social work research textbook and our late-night sessions propped up in bed together, surrounded by laptops, books, and papers, were a source of amusement to our four children.

This book is respectfully dedicated to William Shadish, whose exemplary work in exploring the limits and strengths of quasi-experimental and experimental research designs has inspired a generation of researchers in the human service professions. Over the years, I have followed his continuing oeuvre with admiration, respect, and humility. I was delighted to find him a warm, engaging, and friendly soul when we

attended a conference together at the Lejongdals Slott, near Stockholm, Sweden, in the winter of 2011. He has my thanks for his professional contributions.

Bruce A. Thyer, Ph.D., LCSW, BCBA-D
Tallahassee, FL



Quasi-Experimental Research Designs



1

1

2

3

4

The Role of Group Research Designs to Evaluate Social Work Practice



5 **S**ocial work as a professional discipline has defined itself from its earli-
6 **S**est years as a scientifically grounded field. Initially, its conceptions of
7 “science” were relatively simplistic, consisting of tabulating descriptive
8 information about social problems and of individual clients. The social
9 survey movement of the late 1800s and early 1900s almost defined what
10 was then meant by scientific research within our field, and was very useful
11 in setting the stage for social reform legislation by revealing the extent
12 and seriousness of social problems in urban areas (Gordon, 1973; Bales,
13 1996). To gain a sense of the magnitude of these efforts, recall that Charles
14 Booth’s survey of the *Life and Labour of the People of London* (Booth,
15 1902–1903) required more than a decade to complete and ultimately
16 comprised some 17 volumes! The Pittsburgh Survey (Devine, 1908) is the
17 closest U.S. counterpart to Booth’s massive project, and W. E. B. DuBois’
18 (1899) *The Philadelphia Negro* represents one of the earliest empirical
19 studies of the psychosocial and economic life of African Americans. Jane
20 Addams and her staff at Hull House, in Chicago, devoted considerable
21 effort to the graphic depiction of social conditions through the construc-
22 tion of the Hull House maps (Residents of Hull House, 1895), again a
23 purely descriptive project aimed at illuminating social pathologies and

4 Quasi-Experimental Research Designs

1 suggesting possible ways to alleviate them. The contemporary Hull House
2 museum provides reproductions of these impressive Hull House maps
3 (see [http://www.uic.edu/jaddams/hull/urbanexp/geography/geography.](http://www.uic.edu/jaddams/hull/urbanexp/geography/geography.htm)
4 [htm](http://www.uic.edu/jaddams/hull/urbanexp/geography/geography.htm)). Such descriptive surveys and studies were of immense value, even
5 if the naïve assumption that solutions would somehow become self-
6 evident once social problems were sufficiently identified and quantified
7 turned out to be overly optimistic.

8 Later in the 20th century, social science in general, including social
9 work, became more interested not only in describing social work condi-
10 tions, but in attempting to more accurately understand their causes and
11 correlates. Such research became technically feasible through the devel-
12 opment of more sophisticated methods of statistical description, analysis,
13 and inference. Correlational statistics and inferential tests were developed
14 that helped social and behavioral scientists make more legitimate asser-
15 tions about apparent associations, differences, and changes observed in
16 the heretofore purely descriptive quantitative data they gathered. As sam-
17 pling techniques and statistical methods of analysis improved, so did
18 research methods themselves, to the point at which it became possible to
19 undertake systematic evaluations of the outcomes of social work services.
20 For example, Carl R. Rogers, the founder of person-centered therapy, was
21 employed early on in his career as a social worker, and he conducted a
22 small-scale evaluation of a foster home involving 10 boys. These were tough
23 cases. “In spite of the fact that not one of the whole group was over the age
24 of 11 years of age, sex misconduct of every variety was represented—
25 masturbation, attempted sexual intercourse, filthy language, incest,
26 extreme sexual curiosity, sex perversions of every sort” (Rogers, 1933,
27 p. 21). He was able to present data on the IQ scores of the boys obtained
28 when they initially entered the foster home and again some 3 years later,
29 showing “a slight, but significant increase in mentality” (Rogers, 1933,
30 p. 37). Rogers noted that the effects were small and that it was not possi-
31 ble to conclusively assert that the positive changes observed were *caused*
32 by the foster home’s beneficial effects. Despite these limitations, his study
33 was considered sufficiently exemplary to be worthy of being reprinted in
34 Lowry’s (1939) *Readings in Social Casework: 1920–1938*.

35 Similar small-scale investigations, now known as pre-experimental
36 studies, involving the more or less systematic assessment of clients
37 before they received a social work program and again some time later,
38 provided crude efforts to see if clients really *were* getting better following

1 participation in our discipline's services, began to accrue during the
2 1930s and 1940s. These, in turn, were supplemented by somewhat more
3 sophisticated studies involving naturally occurring control or compari-
4 son groups (such studies became known as quasi-experiments). By being
5 able to *compare* the outcomes of clients receiving special social work ser-
6 vices with those receiving either no such services or some alternative
7 form of intervention, it is possible to have a stronger sense of the actual
8 impact of the special social work services. Later, added to the mix,
9 emerged a small number of true experimental studies, ones in which the
10 control or comparison groups were created using random assignment
11 procedures. This helps to ensure that the groups were essentially equiva-
12 lent on most significant variables, including demographic features, prob-
13 lem severity, the possession of strengths, and any personal or social assets
14 that may impact how a treatment could improve someone's situation.

15 F. Stewart Chapin, then Director of the School of Social Work at the
16 University of Minnesota, wrote favorably on this topic in 1949, in his
17 article titled *The Experimental Method in the Study of Human Relations*.
18 According to Chapin:

19 In the interest of clear thinking about this problem it is helpful to dis-
20 tinguish, first, the trial-and-error "experiments" of social legislation as a
21 means to achieve some desired end; . . . second, the operations of natural
22 social forces that produce an effect; . . . and, third, the use of experimen-
23 tal designs as a method of the study of the first two, in order to determine
24 the degree of success in the attainment of a desired social end, or to mea-
25 sure the effect of some social force. . . (Chapin, 1949, p. 132)

26 There are three general patterns of experimental design in the study of
27 human relations: . . . first, a cross-sectional design in which comparison
28 is made for a given date between an experimental group which receives a
29 social program, and a matched control group denied this program; sec-
30 ond, a projected design in which before and after measurements are made
31 upon an experimental group which received a program over an interval of
32 time, and a matched control group denied this program; and, third, what
33 may be called the *ex post facto* design, in which a present situation is taken
34 as an effect of some assumed and previously operating causal complex of
35 factors, and, depending on the adequacy of accessible records, an experi-
36 mental group and a matched control group are traced back to an earlier

6 Quasi-Experimental Research Designs

1 date when the forces to be measured began functioning upon the experi-
2 mental group but not upon the control group. (Chapin, 1949, p. 133)

3 Chapin then went on to describe examples of using such designs in
4 the evaluation of various social work and welfare interventions (e.g., the
5 Boy Scouts, the federal Works Progress Administration, public housing,
6 juvenile delinquency intervention, the effects of education on income,
7 etc.). As we shall see later in this book, Chapin's description of posttest-
8 only controlled group designs, and of pre- and posttest controlled group
9 studies are distinctions that remain major forms of what are now called
10 quasi-experimental research designs and whose inferential logic remains
11 largely unchanged through to the present. He was optimistic about the
12 potential for utilizing research methods in the appraisal of the results of
13 social work, unabashedly claiming it to be possible to identify causal rela-
14 tionships in human affairs:

15 The experimental method has contributed in large measure to the strik-
16 ing achievements of modern science. This method allows us to analyze
17 our relations of cause and effect more rapidly and clearly than by any
18 other method. It permits verification by many observers. It has sub-
19 stituted for unreasonable prejudice a definite sort of proof that has
20 attained sufficient certainty to justify prediction. . . . Experiment is sim-
21 ply observation under controlled conditions. When observation alone
22 fails to disclose the factors that operate in a given problem, it is necessary
23 for the scientist to resort to experiment. The line between observation
24 and experiment is not a sharp one. Observation tends gradually to take
25 on the character of an experiment. Experiment may be considered to have
26 begun when there is actual human interference with the conditions that
27 determine the phenomenon under observation. (Chapin, 1917, p. 133)

28 The reader may have noted the 42-year gap between these two lengthy
29 quotes by Chapin.

30 In the 1970s, social worker Joel Fischer (1973, 1976) completed a
31 comprehensive review of quasi and experimental studies on the out-
32 comes of what was then labeled "social casework," services largely pro-
33 vided by workers holding master's degrees in social work (MSWs).
34 He found one such study that had been published during the 1940s, two
35 during the 1950s, and 11 that had been published during the 1960s.

1 Fischer found that when all these prior studies were examined, and the
2 original authors' conclusions summarized, there was very little evidence
3 that clients benefited from receiving conventional social casework ser-
4 vices, and considerable evidence showed that a significant percentage
5 were harmed as a result of their receipt of social services. Fischer's review
6 was buttressed by other similar but independent analyses of the same
7 literature, conducted about the same time, which arrived at equally
8 dismal conclusions (e.g., Segal, 1972; Grey & Dermody, 1972). This was
9 very bleak news indeed, and Fischer's conclusions provoked a storm of
10 controversy, with some reactions being largely defensive in nature, others
11 consisting of personalized criticism of Fischer's motivations for conduct-
12 ing these reviews, and some claiming that the effects of social work were
13 simply not amenable to scientific investigations (e.g., Pharis, 1976).
14 Fortunately, this latter view was not widely held.

15 One reaction to the Fischer assessment was a recognition that evalu-
16 ation studies needed to employ more scientifically legitimate outcome
17 measures, finally taking seriously the early recommendations of Mary
18 Macdonald:

19 The first essential, then, for evaluative research on practice is to make
20 explicitly, specific, and concrete the objectives towards which practice
21 is directed. . . . The essence of research is that the findings relate to that
22 *which is observed and not to the individual observer*. This is the criterion
23 of objectivity, or reliability, and it is one to which until recently such
24 evaluative research as we had in social casework has given little or no
25 attention. In research, the burden of proof is on the investigator, and he
26 is expected to show that his results are not a matter of personal whim.
27 One step in this direction has been taken when success is defined in spe-
28 cific and concrete terms. (Macdonald, 1952, p. 136, emphasis added)

29 Other positive reactions to the Fischer report were to focus on more
30 narrowly circumscribed issues and problems addressed using interven-
31 tions that could be operationalized well, thus enabling others to replicate
32 those essential elements of the social work services possibly deemed effec-
33 tive in producing positive change. A further improvement was to adopt
34 research designs of greater scientific credibility, thus permitting a clearer
35 determination of the effects of social work. Within a decade following
36 Fischer (1973), the picture had changed considerably, and for the better.

1 Reid and Hanrahan (1982) published a further review of more recently
2 conducted outcome studies on social work, finding largely positive
3 results, as did a number of other reviewers (Thomlison, 1984).
4 In 1988, Lynn Videka-Sherman undertook the largest effort to date at
5 tracking down outcome studies on social work, and conducted what is
6 called a *meta-analysis* on these studies, a method of integrating the find-
7 ings of disparate studies. Videka-Sherman claimed to find largely posi-
8 tive effects for social work services provided across a wide array of practice
9 domains. Although Videka-Sherman's analysis regrettably included a
10 large number of non-social work studies and included some other sig-
11 nificant mistakes that clouded the conclusions that could be drawn about
12 social work per se, the overall report was seen as important and served to
13 bolster the profession's claim that it was indeed capable of assisting cli-
14 ents in a meaningful manner. Although each of these newer outcome
15 studies can be and were legitimately criticized as overly optimistic and
16 insufficiently critical (e.g., Hogarty, 1989; Epstein, 1990), it is fair to say
17 that, in terms of methodological sophistication and results, the 1980s
18 brought to light an increasing array of evaluation studies that improved
19 the evidentiary foundations and justification for social work services.
20 Subsequent systematic reviews and meta-analyses have reinforced this
21 conclusion (e.g., de Schmidt & Gorey, 1997; Gorey, Thyer, & Pawluck,
22 1998; Grenier & Gorey, 1998). At present, more quasi-experiments and
23 randomized controlled trials of social work services are published in a
24 given year than occurred during entire decades prior to 1980. This reflects
25 the discipline's maturation as a legitimate profession based on credible
26 knowledge derived from high-quality social and behavioral science.

27 THE ROLE OF QUASI-EXPERIMENTAL STUDIES

28 Pre-experimental and quasi-experimental research designs are often
29 used to try to evaluate the effects of a social program, a particular type of
30 psychotherapy or some other form of psychosocial intervention, or the
31 results of public policy. They are also widely used in medicine to evaluate
32 the effects of medications. The term *design* in its broadest sense refers to
33 all the elements that go into creating and conducting a research study,
34 features such as forming a broad research question; creating a specific,
35 directional, and falsifiable hypothesis; deciding upon a unit of analysis

1 (e.g., individuals, groups of persons, organizations, communities, coun-
2 ties, states, countries, etc.); selecting a sample of clients or other units of
3 analysis; choosing one or more outcome measures; developing a way to
4 deliver the intervention in a credible manner; figuring out how to assess
5 client functioning following (and sometimes before) receipt of the inter-
6 vention; analyzing the results; and integrating the findings back into any
7 relevant body of theory that the hypotheses were based upon (recogniz-
8 ing that not all hypotheses are based on an explicit behavioral or social
9 science theory). Although this book discusses each of these elements of an
10 outcome study to some extent, our primary focus will be upon the *logical*
11 *or comparative* aspects of a research project, those features that comprise
12 the most commonly understood meaning of the term “design.”

13 Traditionally, research designs used in outcome studies have been
14 broadly categorized into three types. Those which involve the analysis of a
15 single group of clients have traditionally been called *pre-experimental*
16 *designs*. Those that involve comparing the outcomes of one group receiving
17 a treatment that is the focus of evaluation to one or more groups of clients
18 who receive either nothing or an alternative real treatment, or to a group
19 receiving a placebo-type treatment, have been called *quasi-experimental*
20 *designs*. And the third type, *true experiments*, are characterized by creating
21 different groups (those receiving “real” treatment vs. those receiving noth-
22 ing, alternative treatment, or placebo) by randomly assigning clients (or
23 other units of analysis) to those various treatment conditions. It is with
24 these distinctions in mind that traditional research textbooks have used the
25 terms pre-experimental designs, quasi-experimental designs, and true
26 experimental designs (e.g., Campbell & Stanley, 1963; Cook & Campbell,
27 1979; Rubin & Babbie, 2008; Thyer, 2010a; Yegidis, Weinbach, & Myers,
28 2011). Some authorities have classified pre-experimental designs as quasi-
29 experiments (Shadish, Cook, & Campbell, 2002), and it is with this sense in
30 mind that the title of the present book is derived. However, given the con-
31 tinuing widespread understanding of the traditional distinction between
32 pre- and quasi-experimental designs, these designs will be addressed in
33 separate chapters in this volume.

34 When a researcher has access to a relatively large number of people
35 who have or will receive a particular type of social work intervention, and
36 she wishes to try to figure out what the effects of that intervention may
37 have been, then pre- and quasi-experimental group research designs can
38 be an excellent approach. For many years, authorities in social work have

1 claimed that the design and conduct of outcome studies on social work is
 2 one of the most, if not *the* most valuable type of research project that can
 3 be undertaken, given social work's applied interests (see Chapter 1 in
 4 Royse, Thyer, & Padgett, 2010, which reviews this position). I tend to
 5 agree with this perspective (Harrison & Thyer, 1988).

6 A very large proportion of social work services and programs are
 7 undertaken with no systematic efforts made to evaluate the outcomes of
 8 these services. This is despite language found in the Code of Ethics of the
 9 National Association of Social Workers (2008), as in Standard 5.02 for
 10 dealing with Evaluation and Research:

- 11 (a) Social workers should monitor and evaluate policies, the
 12 implementation of programs, and practice interventions.
 13 (b) Social workers should promote and facilitate evaluation and
 14 research to contribute to the development of knowledge.

15 It is clear that practice evaluation studies are not only the purview of
 16 the academic social worker but of all professional practitioners in our
 17 discipline. Relatedly, social work's accreditation organization, the
 18 Council on Social Work Education, provides guidelines as to what is
 19 required to be taught within the bachelor of social work (BSW) and
 20 MSW programs. Among the policies found in the social work Educational
 21 Policy and Accreditation Standards (Council on Social Word Education,
 22 2008) is the following:

23 **Educational Policy 2.1.6—Engage in research-informed practice and**
 24 **practice-informed research.**

25 Social workers use practice experience to inform research, employ
 26 evidence-based interventions, evaluate their own practice, and use
 27 research findings to improve practice, policy, and social service deliv-
 28 ery. Social workers comprehend quantitative and qualitative research
 29 and understand scientific and ethical approaches to building knowledge.

30 Social workers

- 31 • use practice experience to inform scientific inquiry and
- 32 • use research evidence to inform practice.

33 The content of this book, relating to using quasi-experimental designs
 34 to evaluate the outcomes of social work, is core information needed by all

1 social work students and practitioners to critically evaluate our discipline's
2 research literature. It is obvious that research studies should include prac-
3 titioner participation in order for such evaluations to be effectively
4 designed, executed, and completed. And, competent research methodolo-
5 gists are necessary to help practitioners to design credible evaluations of
6 their services. Perhaps the ideal scenario occurs when an experienced
7 social worker goes on to obtain his or her research-based Ph.D., with a
8 focus on acquiring skills in intervention research in social work (Harrison
9 & Thyer, 1988). Most social workers will not actually undertake empirical
10 research of any kind, and few will apply the types of designs described in
11 this book in the context of evaluation studies. But *all* social workers need
12 to have the ability to digest and understand such studies, regardless of
13 whether they are presented in the form of a journal article, book chapter,
14 books, governmental report, or organizational monograph. This book is
15 intended to facilitate the student's and practitioner's ability to be an effec-
16 tive consumer of the research literature across the human services.

17 **SOME QUESTIONS QUASI-EXPERIMENTAL STUDIES CAN ANSWER**

18 These pre-experimental and quasi-experimental designs can be usefully
19 employed to provide legitimate answers to fundamental questions such
20 as the following:

21 1. *What is the status of clients after they have received a given course of*
22 *treatment?*

23 Treatment can refer to a specific form of psychotherapy, a school or
24 community-based intervention, a new social policy, or the like. And cli-
25 ents can refer to not only individuals, but also to other units of analy-
26 sis, such as couples, families, small groups, organizations, communities,
27 counties, states, or even nations. If a program or intervention really *is*
28 effective, those clients assessed when it is completed should have good
29 outcomes. For example, if 100 families provided a new home via the
30 Habitat for Humanity program were assessed 3 years later, and all 100
31 families continue to reside in the Habitat homes they were provided, this
32 would be a very good outcome. However, if only 30 families were found
33 to have maintained ownership of their homes after 3 years, with 70 being
34 unable to afford the minimal payments, then the Habitat program would

12 Quasi-Experimental Research Designs

1 not be seen as very effective in terms of providing long-term housing for
2 the poor. Recently, I (Thyer, 2010b) reported on the pass rates of MSW
3 graduates taking the licensed clinical social worker (LCSW) examination
4 using this type of simple design. It can also be used to examine, for exam-
5 ple, client satisfaction with social work services, or to evaluate how social
6 workers attending a continuing education seminar viewed the quality of
7 the instruction they received.

8 *2. Do clients improve after receiving a given course of treatment?*

9 This question is an incremental increase in sophistication over the first,
10 in that it requires some sort of formal appraisal of the functioning of a
11 group of clients *prior to* their receipt of social work services. It then sys-
12 tematically compares the posttreatment outcomes with the results of the
13 pretreatment assessment. Both assessments must be conducted in a very
14 similar manner for this type of comparison to be legitimate. If improve-
15 ments are found, and they are meaningful, such a result is consistent with
16 the hypothesis that treatment caused the group of clients to improve.
17 Of course, there are other reasons why the clients could have gotten
18 better, so it is not usually permissible to claim that treatment caused the
19 improvements, only that improvements *happened*. If improvements
20 are not found, and the study was conducted properly, the results would
21 be consistent with the hypothesis that the treatment is not effective. Both
22 results can be useful to know.

23 *3. What is the status of clients who have received a given treatment com-
24 pared to those who did not receive that treatment?*

25 If the answers to questions 1 and 2 are positive, the next increasingly
26 sophisticated question relates to determining if the positive changes can
27 be reasonably attributable to the passage of time alone. This requires
28 comparing the outcomes of the treated clients to a similar group of cli-
29 ents who did not receive the intervention. If treated clients end up better
30 off than untreated ones, such results are consistent with the hypothesis
31 that treatment produces more improvements than no treatment or the
32 passage of time alone. This is good to know since many types of client
33 problems wax and wane in severity over time, either in response to the
34 natural history of the problem itself, in response to environmental
35 changes, or in response to both influences. Intrinsic or environmental
36 influences may be presumed to equally affect clients in the treatment and

1 no-treatment groups. With these influences partially controlled for, the
2 remaining possible source of posttreatment differences may be more
3 credibly presumed to be due to the social work intervention received by
4 the treatment group. Of course, other possible confounds remain, such
5 as placebo influences, social desirability bias in client reports, nonspecific
6 relationship effects, and the like.

7 4. *What is the status of clients who have received a novel treatment com-*
8 *pared to those who received a credible placebo treatment?*

9 If it is shown that treated clients have high levels of functioning, that they
10 actually improved following the receipt of social work services, and that
11 they are better off than clients who did not receive a social work interven-
12 tion, a further question of interest consists of finding out if clients who
13 received a novel or experimental treatment fared better than clients who
14 received an inert or placebo-type therapy. Any given social work interven-
15 tion consists of at least three distinct elements that could possibly be
16 responsible for therapeutic change. One important factor consists of the
17 *therapeutic relationship*, which is known to be powerfully influential.
18 A second important factor consists of the specific treatment techniques
19 that are provided, within the context of the therapeutic relationships.
20 The third important factor consists of the positive expectations for change
21 induced simply by the experience of receiving a credible treatment method
22 delivered by a warm and caring therapist. This latter mechanism of change
23 can be called the *placebo factor*, of which more will be said later. Needless
24 to say, placebo influences are ubiquitous in the delivery of interpersonal
25 services in social work and other health care professions. Truly professional
26 social work services consist of the creation of positive changes in clients
27 induced by relationship factors *and* specific therapies, positive changes that
28 are more powerful than placebo factors alone. One does not need years of
29 graduate training to provide placebo treatments, and one would hope that
30 the profession of social work provides treatments that are considerably
31 more powerful than placebo. One very good way to help determine if this
32 is the case is to compare the outcomes of clients who received a real exper-
33 imental social work treatment with the results of clients who received a
34 placebo or sham therapy. This is not as uncommon as you might think.
35 Over a dozen nomothetic studies of social work intervention have used
36 credible placebo groups in an effort to tease out the specific effects of treat-
37 ment from those induced by relationship or placebo factors.

1 Almost 50 years ago, the distinguished social work researcher
2 Margaret Blenkner (1962) authored a very important article titled
3 “Control-groups and the Placebo Effect in Evaluative Research,” an arti-
4 cle that appeared in our profession’s flagship journal, *Social Work*.
5 Blenkner reviewed the salience of placebo effects in the provision of social
6 work services and of the need for evaluation studies to control for pla-
7 cebo effects by using control groups providing apparently credible but
8 actually innocuous therapies. Although concerns exist about providing
9 placebo therapies to clients in need, circumstances can arise wherein pro-
10 viding a placebo-like therapy is both ethical and methodologically legiti-
11 mate. For example, Wolf and Abell (2003) conducted a controlled study
12 on the effects of a specific form of meditation on client psychosocial
13 functioning. The meditation technique involved providing half the
14 clients with a “real” mantra of Sanskrit words (“hare Krishna, hare
15 Krishna. . .”) whereas half received a fake mantra of meaningless Sanskrit
16 words (“sarva dasa, sarva dasa . . .”). Clients received otherwise identical
17 instructions on how to meditate, and their psychosocial functioning was
18 assessed before and after 4 weeks of meditation practice. Clients who
19 received the real meditation method reported lower levels of stress than
20 did the clients who received placebo medication practice. Afterward,
21 those clients who initially received the fake mantra were debriefed and
22 taught the real technique. Similarly, Pignotti (2005) evaluated the effects
23 of a psychotherapy called thought field therapy (TFT). Thought field
24 therapy involves a supposedly crucial treatment technique involving the
25 client tapping on his or her body in very specific positions and timing
26 sequences. In a test of this therapy, Pignotti arranged for half the clients
27 to receive real TFT, with accurate instructions on where to tap their body.
28 The other half were instructed (without their being aware of it) to tap on
29 randomly chosen parts of their body. The real TFT group experienced
30 modest improvements, but so did the placebo TFT group, thus demon-
31 strating that TFT is essentially a technique that primarily relies on pla-
32 cebo influences to bring about its benefits, not anything specific to the
33 tapping technique. Hyun, Lee, Kang, and Choi (2008) treated smokers
34 with two forms of acupuncture: “real” acupuncture, involving real needle
35 placement according to the acupuncture theory of meridians and invis-
36 ible bodily energies unknown to science; and fake acupuncture, in which
37 needles were placed on sites supposedly unrelated to the theory behind
38 acupuncture. Subjective nicotine craving was measured posttreatment

1 and found to be reduced in all clients but not to differ between the groups
2 receiving real and fake acupuncture. This suggests that any beneficial
3 effects of acupuncture are attributable to placebo factors, not to any
4 technique-specific effects.

5 In a final example, two social workers conducted a study of the effects
6 of eye-movement desensitization and reprocessing therapy (EMDR) on
7 women prisoners with a history of being physically abused. The clients
8 first received a treatment that the researchers were sure would not be
9 helpful—relaxation training (RT), something not known to help persons
10 with posttraumatic stress disorder (PTSD) symptomatology. After a
11 period of receiving RT, all clients then received EMDR. Generally speak-
12 ing, there were few to no differences between the small improvements
13 observed following RT or EMDR. Given that EMDR is claimed by its
14 proponents to produce dramatic improvements above and beyond pla-
15 cebo influences, this small-scale study suggested the opposite: that EMDR
16 has effects similar to those induced by a placebo treatment (Colosetti &
17 Thyer, 2000).

18 Currently, however, placebo-control conditions are not a com-
19 mon feature of social work outcome studies that make use of quasi-
20 experimental designs. Given the seriousness of the conditions that our
21 clients bring us, and the general paucity of evidence that clients are helped
22 at all by what is routinely provided, few intervention research programs
23 have advanced to the stage of needing to control for placebo influences.
24 As our disciplinary outcomes-research endeavors mature over the ensu-
25 ing decades, the use of credible placebo-control groups will become more
26 evident.

27 5. *What is the status of clients who have received a novel treatment com-*
28 *pared to those who received the usual treatment or care?*

29 Treatment innovations occur all the time in social work, and it can be
30 useful to compare the outcomes of clients who receive a new interven-
31 tion, relative to clients who receive standard care or treatment as usual.
32 Early on in social work outcomes research, studies were undertaken
33 comparing the results of standard, longer-term social casework against
34 the results obtained from clients who received time-limited intervention
35 or other forms of briefer therapies (e.g., Reid & Shyne, 1969). Some
36 might even contend that it is ethically important to compare any new
37 treatment against standard accepted care—and to demonstrate that the

1 innovative approach is both safe and effective—prior to the new therapy
2 being adopted on a wide scale.

3 This list of five fundamental questions could be extended a good deal
4 by asking about additional practical issues, such as the durability of any
5 initial effects observed during follow-up periods of various lengths, about
6 possible side effects (good and bad) seen following receipt of treatments,
7 of the costs of care relative to the benefits apparently received, and more.
8 But, for now, we will limit our discussion to these five initial and incre-
9 mentally more complex questions.

10 Notice the cautious language presented in Questions 1–5. We are not
11 making any claims about treatment *causing* any observed outcome or
12 improvement; rather, we are seeking answers to more modest questions.
13 The question *Do clients get better following treatment?* is much easier to
14 answer than is the question *Did the treatment cause the clients to get better?*
15 Questions 1–5 may be answerable using the pre-experimental and quasi-
16 experimental designs presented in this book, because questions regarding
17 the causal effects of intervention usually require more rigorous designs,
18 called *randomized controlled trials* (RCTs), as discussed in Solomon and
19 Draine (2009) and in Cnaan and Tripodi (2010), among many other
20 sources. Nevertheless, we believe that Questions 1–5 are worth investigat-
21 ing. Few social work programs, and few therapies, can provide definitive
22 answers, in a scientifically credible manner, to Questions 1–5. Providing
23 such answers can be a useful preliminary to conducting the more complex
24 and difficult RCTs described elsewhere. For example, if simple investiga-
25 tions using pre-experimental and quasi-experimental designs reveal that
26 clients who received treatment X did not get better, or that X is followed
27 by results no better than standard care or a placebo treatment, then it
28 makes little sense to conduct a more complex and expensive RCT to fur-
29 ther evaluate treatment X. In that sense, the simpler designs described in
30 this chapter can be a useful screening method to distinguish *potentially*
31 effective treatments from ineffective ones. Caution is warranted here in
32 the possible case of a therapy that slows deterioration, but does not
33 enhance clients' absolute levels of functioning. A quasi-experimental
34 study might find that clients are worse following therapy; but, absent
35 proper comparison groups, the researcher might not know that the treated
36 clients were actually better off than if they had not received treatment.

37 It was said earlier that the designs described in this chapter are widely
38 used. How widely? Well, Rubin and Parrish (2007) reviewed every issue

1 of two major journals that tend to report more outcome studies than do
2 other social work journals, *Research on Social Work Practice* and *Social*
3 *Work Research*, published during a 6-year period (2000–2005). They cat-
4 egorized the articles according to the type of design used in empirical
5 outcome studies. They found a total of 28 quasi-experimental studies of
6 social work practice and 21 pre-experimental studies (using the one-
7 group pretest–posttest design described below), making a total of 49
8 published studies using the methods described in this book. In contrast,
9 they found only 16 studies that used a randomized experimental design,
10 nine used single-system research designs, 11 were correlational investiga-
11 tions, and one was a qualitative article. In a related study, Holosko (2010)
12 examined 3 years' (2005, 2006, and 2007) worth of articles published in
13 three major social work journals (*Research on Social Work Practice*, *Social*
14 *Work Research*, *Journal of Social Service Research*), with a focus on the
15 designs used in evaluating the outcomes of practice. He found a total of
16 five randomized controlled trials, but 53 studies making use of pre- and
17 quasi-experimental research designs. Clearly, these latter approaches
18 are the designs most frequently used to evaluate the outcomes of social
19 work practice, at least as reflected in these three selective and highly cited
20 journals.

21 **PURPOSES OF QUASI-EXPERIMENTAL STUDIES**

22 Some contend that the highest or most sophisticated type of knowl-
23 edge is *causal knowledge*, information that allows us to accurately predict
24 what will happen to people who receive a particular social work interven-
25 tion. As a result, less sophisticated research is sometimes minimized,
26 seen as not worth doing, or as relatively unimportant. In this section,
27 my hope is to disabuse the reader of such notions by describing some
28 ways in which quasi-experimental designs can be useful in social work
29 research.

30 **Initial Screening of Treatments**

31 The design and conduct of an intervention study can be a massive
32 undertaking, expensive in terms of time, money, and other resources.

1 Many treatments, sadly, will ultimately prove to not be genuinely helpful.
2 It is possible through small-scale quasi-experiments to identify interven-
3 tions that clearly *do not work*. This information can be of considerable
4 value—clients can subsequently avoid being exposed to ineffective inter-
5 ventions, and researchers can move on to devote their attentions to more
6 promising lines of treatment. An example from the field of medicine can
7 help convey this point. Suppose it is predicted that a new strain of influ-
8 enza will emerge next year, and the Snape Potions Company has pre-
9 pared a special vaccine that is supposed to protect those immunized with
10 it from contracting this new strain of flu. It would be relatively simple to
11 recruit a sample of 100 healthy volunteers, immunize all of them, then
12 expose them to the new flu virus. If 100% of the sample, all 100 patients,
13 subsequently came down with the new type of flu, this would be pretty
14 convincing evidence that the new vaccine was not worth pursuing. With
15 this simple study used as a preliminary screening test, the company could
16 avoid the expense of a much larger-scale study involving many hundreds
17 of patients, various control groups, and the like.

18 We could envision something similar used in the human services.
19 Suppose a novel psychotherapy called rectification therapy (a term
20 I made up for the purposes of this book) was claimed to be very highly
21 effective at preventing relapse (defined as subsequent attempts to kill
22 oneself) among persons with major depression who attempted suicide
23 for the first time. To test this prediction, we could provide rectifica-
24 tion therapy to a consecutive series of depressed patients who had
25 recently made their first suicide attempt and follow-up on their incidence
26 of reattempting suicide over the next 12 months. If 100% were deter-
27 mined to have tried to kill themselves again, then rectification therapy
28 would rather convincingly have been shown not to be very effective in
29 achieving its intended goal of preventing all further suicide attempts.
30 With proper comparison groups, it might be possible to demonstrate
31 whether clients treated with rectification therapy had *fewer* suicide
32 attempts than untreated clients, as well. It can be a mistake to undertake,
33 as an *initial* appraisal of the effectiveness of a new treatment, a very com-
34 plex design with hundreds of clients. Most therapies will, in the fullness
35 of time, turn out to be not very useful. This finding can be determined
36 with very simple studies. Do small-scale studies first as a screen, and only
37 pursue more ambitious ones if the intervention passes the preliminary

1 trials presented by the simpler studies and show at least *some* promise as
2 useful.

3 **Testing and Advancing Theory via Corroboration**
4 **or Falsification of Hypotheses**

5 Another useful function of quasi-experiments is to provide preliminary
6 tests of hypotheses or answers to questions that are not about practice
7 outcomes, but that address other potentially important issues, say, related
8 to the causes or etiology of selected psychosocial problems. Take a clinical
9 observation of Sigmund Freud, related to the etiology of agoraphobia:

10 In the case of agoraphobia etc., we often find the recollection of an
11 anxiety attack and what the patient fears is the reoccurrence of such an
12 attack, under the special circumstances in which he believes he cannot
13 escape. (Freud, 1962/1894, p. 81)

14 Based on his work with a number of agoraphobic patients, Freud
15 came to believe that the agoraphobia stemmed from their experiencing a
16 panic attack. Notice that he said this for persons with agoraphobia, not
17 persons who are depressed, schizophrenic, or otherwise afflicted. Thus,
18 this is a specific and easily falsifiable hypothesis. You could test it via a
19 highly controlled quasi-experiment by systematically looking at several
20 groups of clients (agoraphobics, depressed, schizophrenic, etc.) and
21 assessing them via a reliable and valid history-taking of their experiences
22 with panic attacks. If such experiences were very prevalent among those
23 with agoraphobia and *not* among those with other disorders, then Freud's
24 hypothesis could be said to be corroborated or supported (we rarely say
25 "confirmed" in the behavioral sciences). Such a study would be a com-
26 plex undertaking. However, you could do a simpler test. Ask a number of
27 persons with agoraphobia to complete a simple questionnaire about their
28 experience with panic attacks. If the incidence is high, then you have
29 provisional or preliminary support for Freud's hypothesis and thus some
30 justification for undertaking a larger-scale investigation. But if, contrary
31 to his hypothesis, you found that very few persons with agoraphobia
32 reported such a history, then you might decide that this is a line of etio-
33 logical research that will likely prove to be a dead end, and thus opt to do

1 something else. In such an instance, this very simple research design is
2 quite valuable.

3 **Developing Generalizable Knowledge**

4 In many areas of psychosocial and health research, it is simply not possi-
5 ble to conduct true experiments, experimental studies that afford the best
6 opportunity to develop true knowledge about the causes of problems and
7 the real effects of interventions. Sometimes these difficulties are logistical.
8 It may not be possible to gain access to a sufficiently large number of
9 subjects to make a large-scale study feasible. People presenting with some
10 issues, such as sex offenders, sex workers, prisoners, members of visible
11 minorities, substance abusers, and the like may be reluctant to voluntarily
12 participate in any kind of research project. And, perhaps most commonly,
13 it may not be practical to develop and implement procedures to randomly
14 assign clients to experimental versus standard versus placebo versus
15 no-treatment conditions. Ethical considerations may preclude random
16 assignment methods, or one's institutional review board may not approve
17 of randomly assigning clients in need to control or comparison condi-
18 tions. In such instances, true experiments may simply not be possible,
19 and one must, perforce, rely on an array of research designs of lesser
20 potential validity. This need not doom one's quest for developing causal
21 knowledge, however. Take the case of smoking and lung cancer.

22 No one has ever designed a randomized controlled study wherein
23 young children were assigned to a condition requiring them to smoke
24 cigarettes from an early age into late adulthood, while others were strictly
25 prohibited from ever smoking, and then having researchers look at the
26 incidence of lung cancer in the two groups. Such a horrible study would
27 be a very good way, scientifically, to see if smoking causes lung cancer; but
28 fortunately, absent such research, there are other lines of evidence we can
29 use to investigate possible associations. One could retrospectively look at
30 the smoking histories of people who do and do not have lung cancer. If a
31 smoking history is disproportionally present in the backgrounds of the
32 lung cancer patients, this correlational evidence points in the direction of
33 concluding that smoking causes cancer. It *points* to this conclusion but
34 does not prove it to be true. One could look at the lung cancer rates among
35 groups of people with high rates of smoking (e.g., the poor, or the French)
36 versus those with lower rates of smoking (the well-to-do, or Mormons),

1 and see if the incidence of lung cancer is higher among those who smoke
2 more. One can look at lung cancer rates in countries with high and low
3 rates of smoking and see if country-wide lung cancer rates systematically
4 vary across these countries. One can examine the incidence of newly diag-
5 nosed lung cancer and see if it varies as the prevalence of smoking changes.
6 For example, in the United States, new cases of lung cancer have declined,
7 roughly proportionately with the declines in the numbers of people
8 smoking. And one can look at animal research, examining the emergence
9 of lung cancer in laboratory animals exposed to tobacco smoke versus
10 clean air for long periods of time. If (as it so happens), over time, across
11 research groups, and across countries, the nonexperimental research con-
12 sistently points in the direction favoring the hypothesis that smoking
13 causes lung cancer, and absent credible counterfactual evidence, the sci-
14 entific community eventually concludes that there is a true causal asso-
15 ciation, and we thus have warning labels on cigarettes telling us that
16 smoking causes lung cancer.

17 The lesson here is that an array of quasi-experimental, correlational,
18 and epidemiological studies have the potential to provide our field with
19 relatively plausible causal knowledge about the true effects of certain
20 psychosocial interventions. If quasi-experimental study after quasi-
21 experimental study points to the same conclusion (e.g., rectification
22 therapy helps people with problem X) then, even in the absence of true
23 experimental evidence, the field can provisionally accept this hypothesis,
24 always being ready, of course, to revisit this conclusion as more evidence
25 accumulates. This illustrates another manner in which quasi-experiments
26 can be useful—in the development of generalizable knowledge.

27 **Obtaining Pilot Data in Support of Research Grant Applications**

28 In many circles, the acme of academic success is getting a large-scale federal
29 research grant funded. Many research proposals are prepared in response
30 to a federal request for proposals, for areas of research the government
31 particularly wishes addressed. Others are unsolicited applications. In either
32 case, grant applications proposing to conduct a large-scale randomized
33 true experiment can be considerably strengthened by including informa-
34 tion from previously conducted quasi-experiments—pilot studies, if
35 you will—preliminary to a more rigorous investigation. One of social
36 work's most successful recipients of large-scale federal research funding is

22 Quasi-Experimental Research Designs

- 1 Dr. Gail Steketee, Dean of the Boston University School of Social Work.
 2 Here is what she and her colleague Scott Geron have to say on this topic:

3 Pilot data are important because they demonstrate the applicant's exper-
 4 tise in a target area and serve as a basis from which the proposed research
 5 is built. Pilot data are essential for obtaining most federal funds and
 6 show the investigator's capacity to complete the study . . . the investiga-
 7 tor should highlight the results from pilot studies that illustrate the need
 8 to conduct the proposed research, including the relevance of the findings
 9 to specific hypotheses, the proposed sample size and methodology, and
 10 the likelihood that study hypotheses will be supported. The investigator
 11 seeking larger-scale funding will need more substantial pilot data illus-
 12 trating good effects in the predicted direction. Key points to consider in
 13 describing pilot studies include the following:

14 Complete pilot studies before submitting federal grant proposals.

15 Refer only to pilot studies . . . that clearly demonstrate one's technical
 16 skills and expertise in the proposed research area.

17 Note how the pilot data are promising but insufficient and that, there-
 18 fore, more data are needed. (Geron & Steketee, 2010, p. 626)

19 Rubin and Babbie (2008, p. 262) echo this advice:

20 [I]f you seek funding for a more ambitious experiment or quasi-exper-
 21 iment, your credibility to potential funding sources will be enhanced if
 22 you can include in your proposal for funding evidence that you were able
 23 to successfully carry out a pilot study and that its results were promising.

24 **Quasi-Experimental Studies As a Teaching Tool**

25 The design and conduct of outcome evaluations in social work and other
 26 human service disciplines is a sophisticated skill, and like all sophisti-
 27 cated skills, they are unlikely to spring forth fully formed, like Athena
 28 from the forehead of Zeus. It is more likely that advanced skill devel-
 29 opment will be based on learning preliminary skills by doing simpler
 30 tasks, completed a number of times to the point of mastery, prior to

1 undertaking more technically difficult projects. In my career as an aca-
2 demic, I have frequently worked with master's- and doctoral-level stu-
3 dents on completing simple pre- or quasi-experiments as a preliminary
4 to undertaking more ambitious ones. Dorothy Carrillo, for example,
5 conducted a pre-experimental outcome study on training students in
6 interviewing skills (Carrillo, Gallant, & Thyer, 1993) and subsequently
7 completed a more sophisticated quasi-experiment as her Ph.D. disserta-
8 tion project on the same topic (Carrillo & Thyer, 1994). Betsy Vonk was
9 involved in conducting two quasi-experimental studies of relatively
10 simple design (Thyer, Vonk, & Tandy, 1996; Vonk, Zucrow, & Thyer,
11 1996), prior to undertaking her more complex dissertation project (Vonk
12 & Thyer, 1999). Similarly, Patrick Bordnick worked on two small-scale
13 pre-experimental outcome studies (Capp, Thyer, & Bordnick, 1997;
14 Crolley, Roys, Thyer, & Bordnick, 1998) to help acquire the skills he
15 needed to do a larger-scale RCT of various inpatient therapies for cocaine
16 addicts (Bordnick, Elkins, Orr, Walters, & Thyer, 2004). It can be a mis-
17 take for an inexperienced researcher to undertake an exceedingly ambi-
18 tious outcome study without having acquired the necessary preliminary
19 experience and skills to successfully pull off the larger project. Working
20 on pre- and quasi-experimental outcome studies serves the dual func-
21 tions of modestly contributing to disciplinary knowledge and in helping
22 the graduate student to develop advanced skills in intervention research.

23 I hope that this section has persuaded the reader that the design and
24 conduct of quasi-experimental outcome studies in social work can be of
25 significant value. They can provide an initial screening of the possible
26 effectiveness of interventions. If the results are positive, then the inter-
27 vention *may* be effective and worth further investigation. If the results are
28 negative, then the intervention is pretty certain to *not* be effective, and
29 you may shortcut a potentially futile line of inquiry. Quasi-experiments
30 testing particular treatments, especially if a series of such studies reach
31 conclusions that consistently point in a similar direction, may yield con-
32 clusions you can be pretty confident in. As Rubin and Babbie (2008,
33 p. 255) note “Despite the lack of random assignment, well-designed qua-
34 si-experiments can have a high degree of internal validity.” Studies by
35 research methodologist William Shadish and others have examined the
36 conclusions reached via high-quality quasi-experiments compared to the
37 same interventions evaluated using randomized controlled experiments.

1 He has found “substantial cause for optimism that conditions do exist
2 under which nonrandomized experiments can yield accurate answers”
3 (Shadish, 2011, p. xxx; see also Shadish & Ragsdale, 1996; Shadish, Clark,
4 & Steiner, 2008; Shadish, Galindo, Wong, Steiner, & Cook, 2011). This is
5 not to say that quasi-experiments are always essentially equivalent to true
6 experiments in their rigor at arriving at valid conclusions, but neither
7 should they be cavalierly dismissed as inadequate for evaluation pur-
8 poses. There are also simple pragmatic considerations, as set forth by
9 Ottenbacher:

10 Cronbach (1983) suggested that outcome studies should be guided by
11 attempts to construct designs that meet situational needs, rather than
12 focusing strictly on the requirements of an idealized true experiment.
13 Creating the best research design, in this view, involves multiple con-
14 siderations, including the purpose of the investigation, the specific set-
15 ting, and the available resources. Cronbach (1983) argues that there is no
16 single ideal standard for designs in clinical or applied environments. Any
17 design is an interplay of resources, possibilities, creativity, and personal
18 judgments. (Ottenbacher, 1997, p. 233)

19 These quasi-experimental designs can also serve a useful role in
20 providing pilot data to be included in research grant applications, in
21 addition to being published in their own right. Grant applications are
22 considerably strengthened via the inclusion of solid pilot data and may
23 pave the way to receive the funding necessary to conduct stronger evalu-
24 ation studies. The final function mentioned was the pedagogical role of
25 participation in quasi-experimental research projects for doctoral stu-
26 dents and others new to the evaluation research field. Practical participa-
27 tion in such projects is often much more valuable than reading about
28 research designs in a textbook.

29 THE THEORETICALLY NEUTRAL ASPECT OF QUASI-EXPERIMENTAL DESIGNS

30 One of the strengths of quasi-experimental designs is that they are suffi-
31 ciently versatile as to be useful in the evaluation of virtually any psycho-
32 social or medical intervention, regardless of the theoretical basis of
33 that treatment. Whether an intervention is derived from social learning

1 theory or psychodynamic principles, transactional analysis or hypnosis,
2 the strengths perspective or a personal deficit orientation, so long as
3 a treatment is hypothesized to produce real improvements in clients'
4 lives, quasi-experimental designs may be usefully employed to empiri-
5 cally ascertain if those improvements occurred or not, and in some cases,
6 to permit tentative causal inferences as to the effects of treatment. It is in
7 this sense that these designs are atheoretical.

8 Quasi-experimental designs are, of course, based upon certain philo-
9 sophical assumptions about the nature of the world (ontology) and how we
10 come to arrive at legitimate knowledge of that world (epistemology), but
11 these assumptions are those shared by conventional scientific inquiry as a
12 whole. Among these foundational positions are empiricism, realism, opera-
13 tionalism, scientific skepticism, naturalism, determinism, parsimony, prag-
14 matism, and the like. Although these fundamental philosophical assumptions
15 may be enjoyably debated by philosophers, such discussions are to some
16 extent fruitless since the issues are not capable of being satisfactorily resolved
17 (which is why they continue to be endlessly argued). Those who accept the
18 philosophical assumptions of mainstream science will find the application
19 of quasi-experimental designs to be a very useful tool to answer questions
20 and to test hypotheses. Those who do not accept these assumptions may
21 find quasi-experimental designs to be an unsatisfactory approach to creat-
22 ing knowledge. It certainly seems reasonable that social scientists have the
23 latitude to adopt any set of philosophical assumptions and research method-
24 ologies they deem appropriate. The ultimate validation of the usefulness of
25 these competing approaches to research will reside in their ability to produce
26 knowledge that helps prevent and solve the serious psychosocial problems
27 that social work clients bring to members of our discipline.

28 Radical feminist and qualitative researcher Liane Davis wryly noted
29 that when she was asked by the National Association of Social Workers'
30 National Committee on Women's Issues to examine gender disparities
31 within the profession of social work itself, she

32 [O]btained a large data-set and was cranking out statistic after statis-
33 tic on her office computer. Using this quantitative research method
34 I was once again demonstrating that female social workers earn less
35 than male social workers, even when controlling for important relevant
36 variables. . . . Clearly this is a task that can only be accomplished with
37 quantitative methodology. (Davis, 1994, p. 73)

1 Apart from quantitative statistics, Davis also used a quasi-experimental
 2 design, with gender as an independent variable. This research anecdote
 3 illustrates the principle that one should adopt the research method that
 4 will provide the best, most accurate, and credible answer to one's ques-
 5 tions. If a self-described radical feminist and assertive advocate of qualita-
 6 tive research methodologies like Liane Davis can comfortably make use of
 7 quasi-experimental designs, what better proof do we have of the ubiqui-
 8 tous value of these latter approaches?

9 NOMENCLATURE, SYMBOLS, AND ASSUMPTIONS IN DESCRIBING 10 QUASI-EXPERIMENTS

11 Both pre-experimental and quasi-experimental designs use some simple
 12 nomenclature and symbols to provide an outline or sketch of what was
 13 done. The typical symbols and their meaning are described below:

- 14 • O means a period of time in which clients were assessed
 15 (or *Observed*).
- 16 • O_1 means the first time clients were assessed, O_2 the second time,
 17 O_3 the third, and so forth. The measures used in the assessment
 18 of research participants are often called *dependent variables* in
 19 behavioral and social science, but in this book we will call them
 20 outcome measures, since we are focusing on evaluating the
 21 results of intervention.
- 22 • X usually means a novel treatment, the service that is the
 23 primary focus of an interventive study. In much social and
 24 behavioral science research, an intervention such as X may be
 25 called the *independent variable*, that which is manipulated
 26 (e.g., some clients get a therapy, and others do not). While keeping
 27 this usage in mind, this book will usually refer to these conditions
 28 more simply as “treatments” since most of the illustrations will
 29 involve evaluating social work practices or programs. Obviously,
 30 when the factor under investigation varies and we examine its
 31 presumed effects (which is not done in the context of an outcome
 32 study), the phrase “independent variable” is more accurate.
- 33 • Y, Z, or other letters mean other conditions or therapies
 34 received by a group of clients. Y might mean, for a given study,

- 1 treatment as usual (also known as TAU), Z might refer to a
2 placebo treatment, and so forth. Such designations have no
3 standard usage (e.g., Y = treatment as usual), so the meaning
4 of these symbols may vary.
- 5 • X_1 X_2 means the same intervention given on different occasions.
 - 6 • R means that the group was composed using random assignment
7 methods.
 - 8 • n means the size of a group or sample.
 - 9 • N means the size of the population from which n was obtained.

10 The designs used in this chapter are founded on the traditions of
11 what is called *nomothetic research*, that is, research using large numbers
12 of people. In the case of outcome studies of social work practice, it is
13 more respectful to refer to the people involved in our research as clients,
14 since they are, after all, often real-life social work clients. This is in con-
15 trast to most other behavioral and social work science research projects
16 that use people to test hypotheses derived from theories solely for the
17 purposes of knowledge development, to advance knowledge for its own
18 sake. In the latter instances, the people being studied are usually called
19 *subjects* or perhaps, more recently, research *participants*.

20 To properly assess the results of one's observations of large numbers
21 of clients, social work evaluators usually use one or more of several meth-
22 ods commonly called *inferential statistics*, statistical tests which, when
23 used properly, help us in making correct inferences about the status of
24 clients after treatment, whether or not a single group of clients appreciably
25 changed following treatment, or whether two or more groups of clients
26 differ at a given point in time (e.g., after receiving an intervention). Such
27 outcome measures are rated or scored in such a way as to produce aggre-
28 gated information expressed in the form of arithmetic averages, or mean
29 scores and their associated standard deviations. These measures are usu-
30 ally analyzed by a class of inferential statistics called *parametric inferential*
31 *statistics*, statistical tests based in part on the assumption that the large
32 amounts of data involved would approximate a normal or bell-shaped
33 curve, when plotted on a graph. Parametric tests can be used, for example,
34 to see if the mean scores of a single, large group of clients have changed,
35 posttreatment, compared to their pretreatment status. They can also be
36 used to see if, posttreatment, the clients who received Treatment X sig-
37 nificantly differ from the clients who received Treatment Y, and the like.

1 Some outcome measures involve simple yes-or-no categorizations,
2 or may be expressed in terms of numbers, frequency, or percentages.
3 In such cases, the inferential statistics involved in the analysis of the data
4 are often *nonparametric tests*, ones *not* based on the assumption that the
5 data are roughly normally distributed. For example, if at pretreatment
6 100% of the clients met the diagnostic criteria for panic disorder, and
7 after participating in a treatment for anxiety, posttreatment only 60%
8 (or 90%, or 80%, or 30%, etc.) were so diagnosed, a simple nonparamet-
9 ric test could tell you if this was a statistically significant (e.g., reliable, or
10 not likely due to chance) change. One's choice of an inferential statistic
11 should be driven by the type of data obtained and the research design
12 employed, and *not* by a researcher's desire to use a familiar or novel
13 method of statistical analysis. Let the tools fit the job. Do not arrange the
14 job to be assessed by a particular tool. There will be further discussion on
15 the use of inferential statistics in the final chapter of this book.

16 SUMMARY

17 Group research designs have long been used in the evaluation of social
18 work programs and policies, and to produce more basic scientific knowl-
19 edge. One type of group design, known as quasi-experiments, possesses
20 particular advantages when conducting intervention research in real-life
21 agency and other practice settings. Quasi-experimental designs are capa-
22 ble of answering some very important and fundamental questions about
23 how clients fare after receiving social work services. All too often, practi-
24 tioners and agencies do not possess systematic information on client out-
25 comes and follow-up status. Quasi-experimental studies can provide this
26 data. These studies are also useful for screening out ineffective treatments,
27 identifying potentially effective therapies, testing theories and producing
28 generalizable knowledge, as a teaching tool for graduate students, and in
29 producing pilot data to accompany grant applications.

Pre-experimental Research Designs



This chapter reviews the design and conduct of the simplest of nomothetic (involving relatively large numbers of clients) quasi-experimental research designs. These designs involve looking at only *one group* of people, those who received a given social work intervention that is the focus of the study, referred to here and in later chapters as the *treatment group*. The prerequisites to conduct this type of study are straightforward. You need clients, of course, ideally folks experiencing a similar problem, and they should have received a similar program of social work treatment. Without *both* of these two features, you have a jambalaya of ingredients that form an indigestible recipe, and it is very difficult to draw any legitimate conclusions about the effects of treatment when you have very diverse client problems addressed with an array of different interventions. Clients can, and usually do, present with many disparate problems, sometimes concurrently.

For the purposes of evaluation research, one of the lessons learned from Fischer's (1973, 1976) assessment is that intervention should focus on a fairly narrow range of problems. In mental health, it might be clients who meet the diagnostic criteria for a particular disorder, say major depression or obsessive-compulsive disorder (OCD). With the very poor, it might be unemployment only. With substance abusers, it might be individuals who primarily abuse one particular type of drug, say cocaine. By restricting your initial evaluation efforts to a fairly narrowly defined clientele, you actually

1 enhance the likelihood of finding a positive effect. Take mental health as an
2 example. Few psychotherapies or medications can be expected to help with
3 *every type* of disorder. Rather, therapies are often tailored for particular
4 conditions. Thus, we have, for example, the treatment called exposure
5 therapy and response prevention useful in helping folks who meet the cri-
6 teria for OCD; we have assertive community treatment (ACT) for indi-
7 viduals with schizophrenia, Alcoholics' Anonymous (AA) for alcoholics,
8 and social worker Myrna Weissman's interpersonal psychotherapy (IPT)
9 or Aaron Beck's cognitive therapy for persons with depression. We would
10 not expect ACT to be helpful for someone with OCD or for AA to help
11 people with schizophrenia (unless they also abused alcohol). Prozac not-
12 withstanding, we have yet to develop any true panacea treatments. Imagine
13 having a widely varying mix of clients whose primary diagnosis involves
14 dozens of mental disorders, applying a single intervention to them, and
15 expecting to see overall positive results. It is extremely unlikely. But, if we
16 can do a study of the effects of AA alone on persons only suffering from
17 alcohol abuse, the chances of seeing a clear effect are greatly enhanced, *if*
18 the treatment is genuinely useful.

19 In some ways, the situation is like that of the chemist who wishes
20 to control for extraneous effects by using only very pure chemicals com-
21 bined at conditions of standard temperature and pressure. Keep your
22 research picture as pure as possible—successfully isolate an effect—
23 replicate (e.g., reproduce) that effect, *then* and only then, introduce greater
24 variability into the picture. If you can demonstrate that rectification ther-
25 apy seems to help the pure alcoholic, then do further studies of people
26 who abuse alcohol, as well as another drug, say cocaine. If you find the
27 therapy is still beneficial with this more complex clinical picture, consider
28 adding another complicating element or variable, say marijuana. And so
29 on, gradually approximating research clients who resemble more and
30 more the “typical” substance-abusing client seen by social workers in
31 agency-based practice.

32 Psychologists in particular are known for conducting such “pure
33 studies” on various types of psychotherapies. These studies utilize care-
34 fully chosen clients presenting with only one major problem, carefully
35 trained psychotherapists with sterling clinical credentials receiving expert
36 supervision while they treat clients during the study, therapy conducted
37 in controlled surroundings, and the like. Such studies are often critiqued,
38 and legitimately so, for not reflecting representative clinical conditions;

1 that is, clients with multiple and complex problems, clinicians with lim-
 2 ited training and less than adequate supervision, shabby offices, limited
 3 time, etc. The answer to such critiques is not to repudiate undertaking
 4 therapy research but to gradually conduct outcome studies under condi-
 5 tions that increasingly reflect everyday clinical realities, to see if the orig-
 6 inal positive results obtained under “ideal” circumstances hold up under
 7 less than perfect situations. In this way, such replication studies provide
 8 us with greater confidence regarding the usefulness of a new treatment
 9 under clinically representative conditions and can be more effectively
 10 promoted in mainstream practice.

11 Such work is being undertaken. Stewart and Chambless (2009)
 12 recently published a review of cognitive-behavioral therapies (CBTs) for
 13 adults with anxiety disorders. These treatments were initially developed
 14 under very tightly controlled conditions using carefully screened clients
 15 and supervised therapists. Such experiments are called *efficacy studies*
 16 and are basically intended to see if an intervention works under *ideal*
 17 circumstances. Once such efficacious treatments are identified, they can
 18 then be put into practice under clinically realistic conditions and their
 19 effects reevaluated—these studies are then labeled *effectiveness studies*.
 20 In this case, Stewart and Chambless found that the CBTs found useful in
 21 helping adults with disabling anxiety disorders in tightly controlled effi-
 22 cacy studies were also very effective in everyday practice, not just in uni-
 23 versity and hospital clinics, but in public agency practice. This conclusion
 24 is the results of many years of patient research, involving hundreds of
 25 clinicians and researchers, and thousands of clients, but the end result is
 26 very rewarding in determining that we do indeed have some highly effec-
 27 tive therapies to help seriously handicapped clients (see also Shadish,
 28 Matt, Navarro, and Phillips, 2000).

29 All right, let us review some of the simpler group outcome studies
 30 and see how they may be of value.

31 ONE GROUP POSTTREATMENT-ONLY DESIGN

32 Take a look at the following design, and see if you can figure out what it
 33 means:

34 $X - O_1$

1 If you said something along the lines of “A group of social work cli-
 2 ents received a treatment, labeled X, and afterward they were assessed on
 3 some outcome measure,” you would be exactly right. This design is called
 4 the *one group post-treatment-only design* and can be exceedingly useful in
 5 attempting to answer Question 1 from Chapter 1: *What is the status of*
 6 *clients, after they have received a given course of treatment?* Noting the sim-
 7 plicity of this design, you may be tempted to be rather skeptical as to its
 8 value. Let me see if I can persuade you otherwise.

9 Supposed a group of middle school students are provided a compre-
 10 hensive drug education program aimed at deterring drug use. These stu-
 11 dents were able to be followed-up 5 years later, and it was found that 90%
 12 of them admitted to regularly using drugs. What would this tell you about
 13 the effectiveness of the drug education program? Obvious, it is not *highly*
 14 effective. Suppose another group of middle school students received a
 15 school-based abstinence-oriented sex education curriculum, and 5 years
 16 later, 90% anonymously admitted to having regular sexual intercourse
 17 outside of marriage? Would this be useful (if disappointing) information
 18 to know, especially for a school system contemplating adopting the same
 19 abstinence-oriented sex education program? If you are a master’s degree
 20 social work (MSW) student, you may run across advertisements for per-
 21 sons or firms selling social work licensing examination preparation pro-
 22 grams, workshops, manuals, or online tutorials said to enhance your
 23 likelihood of passing your state’s licensing test. Would it make a differ-
 24 ence to you to learn that 95% of MSWs who completed a given proprie-
 25 tary licensed clinical social worker (LCSW) test preparation program
 26 passed the licensing test? And that another program’s “graduates” only
 27 had a pass rate of 60%? Which preparation program would you want to
 28 pay for and complete? Perhaps the $X - O_1$ design has some potential value
 29 after all? Here are some published examples of this type of simple design.

30 Do Children on Medicaid Receive Required Preventive Screening Services?

31 Medicaid, the federal- and state-funded medical insurance program for
 32 the poor, requires early and periodic screening, diagnostic, and treatment
 33 examinations for children and youth under the age of 21. The theory is
 34 that these preventive tests will detect health problems early on and, in
 35 doing so, will result in improved health for poor children. These required
 36 evaluations are in the areas of a comprehensive health and developmental

1 history, an unclothed physical examination, immunizations, laboratory
2 tests, and health education. Staff within the Department of Health and
3 Human Services reviewed the medical records for 345 children receiving
4 Medicaid in nine different states to assess the extent to which these young
5 persons were receiving required screenings. Fully 75% of the children did
6 not receive all required medical, vision, and hearing screenings, and 41%
7 did not receive *any* required medical screenings. When young people did
8 receive such screenings, they were often incomplete (Levinson, 2010).
9 This finding has rather important implications. First, taxpayers are paying
10 for health services from which poor people obtain little benefit, primarily
11 because of their lack of contact with health care providers. The issue is
12 not cost, since these services are free to the patients. It also bears on the
13 potential for planned national health care to significantly impact the
14 access of the poor to physicians and other health care providers.
15 Apparently, simply making free health care available does not ensure
16 ready access, at least not for poor children. This is a very important find-
17 ing and was brought to the attention of policy-makers by a simple post-
18 test-only group design.

19 **Do Families Follow-Through with Referrals?**

20 MSW student intern Wendy Pabian was placed in a large diagnostic and
21 evaluation agency serving clients with serious developmental disabilities
22 and their families. The major purpose of the agency was to not only pro-
23 vide a solid diagnostic assessment in heretofore difficult-to-diagnose
24 cases, but also to provide an array of medical and psychosocial treatment
25 recommendations to the families, so that their children with develop-
26 mental disabilities could be provided the services most likely to facilitate
27 their growth and development to the greatest possible extent. As her fac-
28 ulty liaison, in conversations with Wendy, I asked her if the agency's cli-
29 ents were actually following through on the treatment recommendations
30 so laboriously and expensively provided by the agency's interdisciplinary
31 treatment teams. She said she did not know. I asked her to check it out at
32 the agency itself, and a week's inquiries confirmed her impression—no
33 one at the agency really knew if families were obtaining the recommended
34 services. I suggested that doing a post-service assessment of this issue
35 would be a valuable MSW intern project that would be useful to the
36 agency. She agreed, as did her on-site field instructor and agency staff.

1 Wendy obtained a list of clients and families seen during a given time
2 frame some months earlier, and then contacted them by phone, asking
3 the parents about their seeking or obtaining each of the services recom-
4 mended by Wendy's agency. Happily, it turned out that most *were*
5 obtaining most recommended services, suggesting that the agency was
6 indeed providing a useful service to clients with developmental disabili-
7 ties and their families. Wendy prepared an internal report for the agency,
8 as well as a more formal journal article that was subsequently published
9 (Pabian, Thyer, Straka, & Boyle, 2000).

10 **Are State Medicaid Application Enrollment Forms Readable?**

11 Medicaid is a state- and federally funded program that provides health
12 insurance to poor individuals, primarily among families whose annual
13 income is below the federal poverty level, and to poor pregnant women.
14 For persons to receive Medicaid benefits, they must first apply for them via
15 a written or computer-based application form. Individuals who are poor
16 disproportionately suffer from illiteracy or low-literacy, and/or represent
17 persons whose primary language is not English. The readability of the appli-
18 cation materials used to enroll persons in Medicaid can be a significant
19 barrier for such individuals to qualify for this benefit to which they may be
20 legitimately entitled. If the reading level is too high, the application attempt
21 may be incomplete or abandoned entirely, resulting in a lack of benefits.
22 In an attempt to determine if this may be a problem for the poor, Wilson,
23 Wallace, and DeVoe (2009) obtained Internet-based Medicaid enrollment
24 applications from 49 states (excluding Kentucky, which was unavailable
25 online) and the District of Columbia. The readability of these forms was
26 assessed using a standardized format widely used in the measurement
27 of this construct (readability) as it pertains to health-oriented materials.
28 The readability of the Medicaid applications ranged from the 11th to the
29 18th grade level. The authors believed these levels to be excessively high,
30 exceeding the reading levels recommended for patient education in health
31 literature and the actual reading levels of the average American.

32 Language excerpts from the actual Medicaid forms were presented
33 along with examples of how this same content could be rewritten to a
34 6th-grade level, thus permitting the poor to more readily comprehend
35 the language used in these forms. A very similar posttest-only design was
36 used by Zite, Philipson, and Wallace (2007) to evaluate the readability of

1 the Consent to Sterilization form used within the Medicaid system; this
2 form is used when a person requests that he or she be sterilized as a vol-
3 untary method of birth control. Obviously, it is very important that the
4 application forms and consent materials used in welfare eligibility deter-
5 minations and requests for voluntary sterilization be easy to understand
6 by the members of the population for whom they are intended. This
7 simple design is an appropriate—indeed excellent—approach to see if
8 the poor are being discouraged from applying for and receiving birth
9 control via application materials using an excessively high reading level.

10 **Do Adults Raised as Children in Orphanages Fare Well in Life?**

11 In another example, social worker Laura Myers collaborated with a tra-
12 ditional orphanage located in Florida to conduct a follow-up study of the
13 psychosocial well-being of the adults who had been raised in the orphan-
14 age many years ago. The orphanage maintained a mailing list of all of
15 its “alumni,” and Laura crafted a mailed written survey combined with
16 some standardized and previously published measures of life satisfaction
17 and quality of life to send to these alums. The research question generally
18 dealt with how these individuals fared, later on in life, as assessed on
19 educational, financial, familial, and social variables. There are actually
20 very few studies that have been conducted on this topic, even though
21 traditional orphanages have always played a major role in the country’s
22 child welfare system of care. The orphanage mailed out the surveys and
23 an explanatory cover letter, which were returned directly to the orphan-
24 age by the alums who completed it. After personally identifying informa-
25 tion was removed, the orphanage staff forwarded the completed surveys
26 on to Laura, who analyzed the results. Retrospective appraisals of the
27 care the alums received at the orphanage were quite positive, and they
28 were, overall, doing quite well as adults in terms of education, income,
29 life satisfaction, and other indicators. These results also yielded a respect-
30 able journal publication (Myers & Rittner, 2001) and have important
31 implications for foster care and adoption services.

32 **Are Consumers of Mental Health Services Satisfied with Their Treatment?**

33 Social worker Jan Ligon was a senior administrator with the state
34 mental health services program in Atlanta, Georgia. He helped design

1 and conduct a consumer satisfaction study for the clients of a crisis
2 and stabilization program providing services to clients with a serious
3 mental health or substance abuse program and their families. He used a
4 previously published, reliable, and valid measure of client satisfaction,
5 and surveyed some 54 clients and 29 family members in the program,
6 finding generally high levels of satisfaction, which is of course a good
7 thing (see Ligon & Thyer, 2000).

8 **Does a Coalition-building Program Promote Integrated Services?**

9 MSW student Grace Smith was undertaking her internship at an agency
10 that provided workshops called *coalition-building forums*, to service pro-
11 viders in the fields of developmental disabilities (DD), aging services, and
12 mental health and substance abuse treatment—service providers
13 who typically had little contact with each other in their professional roles.
14 The purpose of these forums was to promote person- and family-
15 centered care for older persons with developmental disabilities and their
16 families. Information was provided on the services available through the
17 aging, mental health, and DD networks, and a presentation was made by
18 an older consumer with a developmental disability (or a family member).
19 A networking session among the service providers and small-group
20 activities related to integrating services across the three service programs
21 were also provided. The purpose of these day-long meetings was to pro-
22 mote networking, collaboration, and increased political activities among
23 the service providers in the months following the networking forums,
24 so that services could be more effectively integrated with this high-risk
25 population.

26 Ms. Smith contacted a sample of about 20% (n = 64 people) of all
27 those who participated in these networking sessions several months later
28 and asked them ten standardized questions regarding the possible impact
29 of the forums on their daily practice. About half (47%) indicated that
30 they had made professional contact with people they had met or heard
31 about through attending the networking sessions. In general, respon-
32 dents indicated improvements in their knowledge, awareness, and atti-
33 tudes relating to older persons with developmental disabilities but
34 relatively few (15%) indicated that they had engaged in lasting collab-
35 orative activities with service providers in other areas to serve such per-
36 sons. This “bottom-line” result was disappointing and suggested that the

1 networking sessions were not an effective mechanism to promote the
2 long-term integration of services, which was the actual aim of the pro-
3 gram (see Smith, Thyer, Clements, & Kropf, 1997). If you were the man-
4 ager or administrator of the agency providing these networking forums,
5 would you find it useful to learn whether or not your participants were
6 undertaking the activities your sessions were intended to promote?
7 Of course you would. Lacking this information, you might go on con-
8 tinuing to deliver the same old format, month after month, and have
9 little incentive to try to improve your ability to really promote better
10 integrated services.

11 **Do Individual Developmental Accounts Help the Poor**
12 **Achieve Financial Goals?**

13 On a more macro scale, we can look at evaluations of individual develop-
14 ment accounts (IDAs). IDAs are a widely used mechanism to provide
15 incentives for the poor to save toward achieving one or more of several
16 very limited goals—save for a house down payment, save for college, save
17 to start a new business, or save to pay for significant home improve-
18 ments. IDAs work by providing financial education and a match (often
19 2:1 or 3:1; i.e., for every dollar you place in your IDA account, it is even-
20 tually matched by \$2 or \$3), which poor individuals are able to set aside
21 toward one of these goals. There are over 500 IDAs programs in the
22 United States. Do IDAs help the poor achieve these goals? How could we
23 find out? Well, one simple way would be to create a group of poor per-
24 sons who enrolled in an IDA program and examine their attainment of
25 one of these targeted goals after a reasonable period of time had elapsed,
26 say 5 years. If, for 1,000 participants saving for a home down payment it
27 was found, 5 years later, that 800 had been able to buy their own home
28 using their IDA savings, this would likely be seen as a good outcome.
29 If only 100 had done so, then the outlook for the effectiveness of IDAs
30 would be less positive. See Richards and Thyer (2011) for a review of
31 the evidence on IDA effectiveness, most of which uses relatively simple
32 quasi-experimental designs.

33 This $X - O_1$ design has many merits. It is a great way to initially eval-
34 uate programs that have not had the advantages of formal pretreatment
35 assessments; it can be used to help screen out obviously ineffective treat-
36 ments; and positive results can encourage further, more sophisticated

1 appraisals of an intervention initially seen as promising through this ini-
 2 tial, simple, evaluation. However, this design is seriously limited in that,
 3 without any systematic assessment of the clients' functioning *before* they
 4 received the intervention, it is logically difficult to make any legitimate
 5 inferences regarding whether the group of clients receiving treatment
 6 actually had *changed*. And, without any kind of comparison group that
 7 did not receive intervention X, it is difficult to make any conclusions as
 8 to how the group that received X might have changed if they had not
 9 received X. Shadish, Cook and Campbell (2002, p. 107) do note:

10 However, the design has merit in rare cases in which much specific back-
 11 ground knowledge exists about how the dependent variable behaves.
 12 For example, background knowledge of calculus is very low and stable in
 13 the average high school population in the United States. So, if students
 14 pass a calculus test at levels substantially above chance after taking a cal-
 15 culus course, this effect is likely due to the course. . . . But for valid,
 16 descriptive causal inferences to result, the effect must be large enough to
 17 stand out clearly, and either the possible alternative causes must be
 18 known and be clearly implausible or there should be no known alterna-
 19 tives that could operate in the study context.

20 This posttest-only design has, in prior years, been labeled the “one-
 21 shot case study.” This language is no longer recommended because of
 22 the confusion the term engenders with the similarly named qualitative
 23 method widely used in the psychotherapy literature, the narratively pre-
 24 sented “case study” of an individual client (e.g., Brandell & Varkas, 2010).
 25 Simple though it is, there are a number of methodological refinements
 26 that can be added to the basic design $X - O_1$ that can strengthen it as an
 27 evaluation method. Some of these techniques are described below.

28 **WAYS TO STRENGTHEN THE POSTTEST-ONLY GROUP DESIGN**

29 **Use Outcome Measures Known to be Reliable and Valid**

30 All outcome measures are not alike. Some have strong psychometric
 31 properties—they are internally consistent and have high test–retest reli-
 32 ability. Their face and content validity is evident, their internal factor

1 structure is robust, and they have strong concurrent and predictive
2 validity. They are easy to complete, score, and interpret, and are low in
3 cost and readily available. Use such measures whenever possible, in lieu
4 of surveys or rating scales especially developed for a particular study and
5 not previously demonstrated to be reliable and valid. Outcome measures
6 must also be sensitive, in that they are capable of detecting small but
7 meaningful differences and changes.

8 **Use Several Different Outcome Measures**

9 A client's level of functioning, strengths, deficits, or psychopathology
10 may be capable of being concurrently measured in several different ways.
11 In the measurement of depression, for example, the Beck Depression
12 Scale (BDI) is intended to be completed by depressed clients themselves
13 and is widely considered to be among the best measures of depression
14 available. The Hamilton Rating Scale also is used to assess depression,
15 but it is intended to be completed by a health care professional or care-
16 giver, thus providing another perspective on a client's depression. A third
17 potential way to assess depression would be to ask the client to complete
18 a sleep log, indicating the times at which he or she fell asleep, woke up,
19 and total duration of sleep each night. A posttest-only study with three
20 somewhat distinct measures of depression, each legitimate in its own
21 right yet tapping into different aspects of depressive phenomenology,
22 can yield a stronger assessment of a client's depressive status than may be
23 obtained by using only one measure. In research, this approach is known
24 as *triangulation*—using different, perhaps imperfect, measures to more
25 adequately capture some outcome variable.

26 **Use Larger Samples of Clients**

27 In nomothetic research, size matters (20 is better than 10, 50 is better
28 than 20). The larger the sample size, the more persons you have to draw
29 conclusions from and the more likely it is that your sample reflects some
30 larger population of interest. Now, to be sure, using a randomly selected
31 sample, one wherein every person in a population has the same likeli-
32 hood of being selected, is the most solid way to ensure having a sample
33 that accurately represents a larger population of interest. This ideal is
34 often not feasible however, especially in research settings such as social

1 service agencies. If your agency sees 300 new cases a year, you can
 2 see how a sample of 100 of these clients, albeit nonrandomly selected,
 3 will yield a more representative sample of your 300 cases than would
 4 having only 20 nonrandomly selected individuals. Also, as we shall see in
 5 Chapter 5, using larger sample sizes increases what is known as *statistical*
 6 *power*, the ability to detect true differences or changes when they really
 7 exist. And this is a good thing. However, be aware that simply having a
 8 larger sample size alone is no guarantee that the sample is more represen-
 9 tative of the larger population of interest. A very large sample of persons
 10 chosen in some biased manner can still be unrepresentative, even if the
 11 bias is unintentional. For example, suppose you wished to sample “tele-
 12 phone users,” and you used the local white pages to select your sample.
 13 Well, no matter how big a group you contacted, you would still largely
 14 leave unsampled cell phone users (who are not listed in the white pages)
 15 and those who do not have a landline registered in their names.

16 Use Multiple Posttreatment Assessments

17 Assessing client functioning once after receipt of services is a good start.
 18 But to do this again after, say, 3 months have elapsed from the termina-
 19 tion of treatment is even better, and 6 months better still. Providing ser-
 20 vices that produce a strong effect immediately after treatment is great,
 21 but if these improvements evaporate after a month or two, we will be
 22 much less excited about the effectiveness of our program. The only way
 23 to make this determination is to conduct suitable follow-up assessments.
 24 If gains are maintained months or years down the road, this is of far
 25 greater import to clients and social work services.

26 The posttest-only design with repeated measures can be diagrammed
 27 as follows:

$$28 \quad X - O_1 - O_2 - O_3$$

29 By adding one or more additional posttreatment assessments, this design
 30 is strengthened by permitting an appraisal of how well any initial effects
 31 (e.g., improvements) have been maintained, or if any other effects emerge
 32 over time. It may be that, immediately posttreatment, the clients' status
 33 was not very good, perhaps leading to the conclusion that treatment

1 could not have possibly had any beneficial effects. But a further post-
2 treatment assessment, perhaps weeks or months after the first one, could
3 disclose that folks were doing very well indeed. The delay in observing
4 any effect complicates inferences about the possible effects of treatment,
5 since the closer in time a presumed effect occurs after exposure to a pre-
6 sumed cause (e.g., treatment), the stronger is the logical warrant to link
7 the apparent results to the treatment.

8 The elements of theory and plausibility also need to be taken into
9 account in trying to make a causal inference. Someone who takes aspirin
10 for a headache expects an effect (e.g., pain relief) in an hour or two,
11 not days later. However, persons with an earache will not display any
12 improvements immediately after taking an antibiotic medication;
13 improvements can be predicted to begin only a day or two later. Similarly,
14 some interventions may also be predicted (ideally in advance) on the
15 basis of clinical observation, prior research, or even of theory, not to
16 produce immediate changes, and only to have improvements become
17 evident some time following receipt of treatment. For example, certain
18 psychotherapies that aim at symptomatic improvement on the basis of
19 the client developing “insight” could reasonably be expected to not yield
20 any improvements for the first month or so. If this is the case, then the
21 social work researcher can deliberately arrange for posttreatment assess-
22 ments to take place only after the interval of time supposedly needed for
23 improvements to become evident.

24 Other treatments could be predicted to have an additive effect,
25 wherein treatment benefits accrue over time, even after intervention has
26 been discontinued. Say that treatment X consists of intensive tutoring in
27 reading provided to low-income elementary school students. One might
28 expect a certain level of reading proficiency when reading ability is mea-
29 sured immediately after the tutoring program is discontinued, say at
30 time 1, or $X - O_1$. It would be reasonable to anticipate that at time 2, say
31 3 months later, reading proficiency would have gotten even better, or
32 $X - O_1 - O_2$, and so forth, with the group of children’s reading scores
33 progressively improving as they gain practice in reading.

34 Use These Designs Prospectively

35 These designs are best used *prospectively*; that, is planned for and with
36 data gathered after the study has been planned. This is opposed to using

1 them *retrospectively*, a method that looks at data after the fact, without
2 having had any plans to use the data for evaluation services when the
3 data were initially gathered.

4 Look for Generalized, Positive Outcomes or Potential Negative Ones

5 It is common in therapy outcome studies to include one or more mea-
6 sures related to a client's primary issue, such as depression, marital dis-
7 cord, abuse of a given substance, or domestic violence. This is, of course,
8 important. If the intervention may be predicted to produce other broad-
9 er-ranging effects, these too might be measured. For example, in the
10 instance of an intervention designed to reduce spousal battering, obtain-
11 ing credible measure of the episodes of abuse is crucial, but including
12 additional measures, say, of *quality of life* or of *life satisfaction*, could also
13 be important. One would hope that a reduction in marital violence would
14 also improve quality of life or life satisfaction, even though these may not
15 have been among the client's primary complaints when initially seeking
16 treatment. If a therapy is effective in alleviating the dysfunctional mood
17 swings associated with bipolar disorder, but the client's level of life satis-
18 faction declines, this is obviously not as satisfactory an outcome as one
19 with a treatment that reduces mood swings and promotes life satisfac-
20 tion. Some therapies can have negative side-effects (even non-drug-
21 related interventions), and screening for these negative side effects is also
22 a valuable adjunct to assessing changes in focal problems.

23 You can probably think of many other ways in which the basic $X - O_1$
24 design could provide useful information. For example, if you are a BSW
25 or MSW student, do you have the opportunity to complete a course eval-
26 uation at the end of each class? Some universities make this information
27 publicly available, for each professor. Suppose you learn that Dr. Thyer
28 consistently earns high ratings on his students' course evaluations?
29 Or that he usually earns very low ratings? Would this information be of
30 any use to you in deciding whether or not to take Thyer's classes? Student
31 course evaluations, client satisfaction studies, follow-ups on prevention
32 and other service programs, studies on the natural history of disorders—
33 these are all types of inquiry in which the $X - O_1$ design can provide
34 very useful information. Do not dismiss them as scientifically useless.
35 If you have a simple question (e.g., *What is the status of clients after they*
36 *received a social work intervention?*), these simple designs can help answer it.

- 1 But if you have a more complicated question, you may need more com-
 2 plicated studies, such as those described below.

3 THE ONE-GROUP PRETEST–POSTTEST DESIGN

- 4 If you would like to answer the second question mentioned Chapter 1,
 5 *Do clients improve after receiving a given course of treatment?*, you can add
 6 a pretreatment assessment to the basic $X - O_1$ design and end up with an
 7 approach that can be diagrammed as follows:

$$8 \qquad O_1 - X - O_2$$

9 Clients are assessed twice, and their responses are aggregated (often using
 10 an average score for all clients) at the pretest and again after they com-
 11 plete treatment, with this second assessment called the posttest. Not sur-
 12 prisingly, this design is called the *one-group pretest–posttest design*.
 13 As before, you need to make use of one or more reliable and valid out-
 14 come measures that assess some aspects of client functioning. These can
 15 be diagnostic measures or an assessment of a problem, a deficit, or a
 16 strength. Clients are assessed, subsequently participate in a social work
 17 treatment program or intervention, and are then assessed again, in the
 18 same manner as during the pretest phase of the study. You can then com-
 19 pare their aggregated posttest scores or measures with their pretest scores
 20 and see if things have changed, overall, for that group of clients. Ideally,
 21 they have, and for the better. Either way, you can at least partially answer
 22 Question 2: “Yes, they are better following Treatment X,” or “No they are
 23 not better after Treatment X,” or “They are worse off, following Treatment
 24 X,” or, least satisfactorily, “The data are too unclear to say if they have
 25 changed.” It is not usually enough to *simply look* at the aggregated scores
 26 pre- and posttreatment when assessing a sample of clients. You usually
 27 need to “test” the changes using one or more inferential statistical tests to
 28 rule out random fluctuations in the data as being responsible for any
 29 noted differences. Such tests will be briefly discussed in Chapter 5.

30 This design is often feasible, even in conventional agency-based set-
 31 tings without many resources, and can be undertaken by individual social
 32 workers who lack advanced degrees and research training (although

1 these can obviously be of help!). These designs can be done *retrospec-*
2 *tively*, in circumstances wherein clients complete some sort of pretreat-
3 ment and posttreatment assessments as a part of routine agency or
4 clinical practice, and the data are later accessible to a given social worker
5 or research team. Simply (!) obtain agency permission to extract the
6 desired information, comply with oversight by any pertinent institu-
7 tional review board, place the data in a spreadsheet or statistical package
8 database, and look at the results. There are many opportunities to con-
9 duct such retrospective studies, as agencies are often sitting on a gold
10 mine of data gathered in this manner, data which often just sit there,
11 unanalyzed, unused, and not contributing to the expansion of knowl-
12 edge so desperately needed by the field.

13 As this book is being written, I am working with a doctoral student in
14 social work who is in charge of a substance abuse treatment program at a
15 state women's prison. Prisoners who are known to have a history of drug
16 abuse are assigned to this residential drug abuse program (RDAP) within
17 the prison. Although the RDAP is a widely used intervention sponsored by
18 the Federal Bureau of Prisons to treat substance abusers, little appears
19 to be known as to whether it is truly effective at deterring substance abuse
20 (see http://www.bop.gov/inmate_programs/substanceabuse_faqs.jsp). We
21 could find no published outcome studies evaluating RDAPs. As a part of
22 the RDAP, prisoners complete several outcome measures pertaining to
23 knowledge about illegal drugs and their health consequences, as well as
24 their attitudes about drug abuse. And they do this again after several
25 months of participation. This information is available to my student
26 within her work role via the prison records for several hundred RDAP
27 participants, and it has not been analyzed. Such data can be characterized
28 as "low-hanging fruit," easy to harvest, and I am urging my doctoral
29 student to gain access and analyze these data as her Ph.D. dissertation
30 research project. In this instance, the average scores on, say, an attitudi-
31 nal measure related to substance abuse could be assessed when the pris-
32 oners are initially enrolled in the RDAP program—the O_1 assessment.
33 Then, the enrollees begin the RDAP and remain in it for 3 months, after
34 which they complete the same attitudinal measure, or O_2 . This results in
35 the complete design being diagrammed as $O_1 - X - O_2$. This *would be*
36 adequate to test the predictive hypothesis that "RDAP participants
37 will display statistically significant improvements related to attitudes
38 about drug abuse, following 3 months of participation in the program."

1 If the average score on a measure of attitudes at time 2, or O_2 , were sig-
2 nificantly improved relative to the average score for the entire group at
3 time 1, or O_1 , then we could legitimately claim that the hypothesis was
4 supported or corroborated. It would be very rare for one to be able to say
5 that the hypothesis was *confirmed*, as this is generally too strong a posi-
6 tion to take with results obtained from nonexperimental research designs.
7 Given the limited empirical evidence that the RDAP program really
8 works, a relatively simple pretest–posttest study of this nature would be
9 a useful addition to the professional literature on the drug rehabilitation
10 of incarcerated offenders.

11 These designs may also be *prospectively* undertaken, planned in
12 advance by an individual or research team to improve or refine the evalu-
13 ation. Sometimes prospectively designed studies can make use of more
14 valid outcome measures than those routinely used by agencies, or steps
15 can be taken to ensure that the treatment is provided to the clients in a
16 manner that promotes adherence to the best practices appropriate to the
17 model of treatment being evaluated. A prospective study may opt to use
18 independent evaluators to assess clients pre- and posttreatment, to try to
19 partially control for the bias engendered by having the therapists who
20 delivered the treatment be the ones who assess whether the clients bene-
21 fited from the treatment. Or, steps can be taken to ensure that the post-
22 tests are really administered in a manner equivalent to the pretests,
23 to partially control for possible changes in how clients were assessed,
24 pre- and posttreatment. For these reasons, prospectively designed out-
25 come studies are usually seen as more rigorous tests of a treatment's
26 effectiveness than are retrospectively designed ones, which are limited to
27 existing data that may not have been gathered with formal research pur-
28 poses in mind. But, having said that, retrospectively undertaken investi-
29 gations are also very worthwhile in their own right. Here are some
30 examples of social workers using the $O_1 - X - O_2$ design.

31 **Evaluating Group Therapy for Children After Homicide**

32 Traumatic grief reactions are common among children who witness homi-
33 cide and other forms of violence, and clinical social workers and other
34 mental health professionals are often called upon to provide therapeutic
35 services to them. Obviously, we wish to deliver care that is genuinely helpful,
36 without adverse side effects, that produces lasting change and is low in cost.

1 There are not many treatments like this, and social workers delivering ser-
2 vices that lack a strong evidence base are understandably interested in evalu-
3 ating the programs they do deliver, to help assure that children are not being
4 harmed (see Lilienfeld, 2007; Barlow, 2010) and perhaps are receiving ben-
5 efit. Such were the circumstances faced by mental health workers in the city
6 of New Orleans, who developed a “Loss and Survival Team” to provide
7 services to homicide victims and children exposed to violence. Each child
8 received a semi-structured group therapy program consisting of eight to ten
9 sessions, and completed a previously published, reliable and valid posttrau-
10 matic stress disorder (PTSD) assessment at the beginning of group therapy
11 and again when it ended. Over time (October 1997–December 2001), some
12 21 groups were completed in ten different public schools, involving a total
13 of 102 African American children. All the group therapists were MSWs
14 or MSW interns. Symptoms of PTSD dropped statistically significantly for
15 the 102 children as a whole, and various subgroup analyses were also per-
16 formed (boys vs. girls, younger vs. older, etc.). This was a very nice illustra-
17 tion of using the one-group pretest–posttest design, as well as of the efforts
18 of social workers providing clinical services to children systematically evalu-
19 ating outcomes following treatment. Question 2 was answered nicely, and
20 the answer was an agreeable “Yes!” The authors appropriately discuss the
21 limitations of their study in their conclusions and also present suggestions
22 to enhance future research with this clinical group and problem situation
23 (see Salloum, 2008).

24 **Evaluating Abstinence-based Sex Education**

25 This design was also used to evaluate an abstinence-oriented empower-
26 ment program for public school youth. The problem area was teenage
27 pregnancy, and the issue was the possible effectiveness or ineffectiveness
28 of the widely used approach called *abstinence education*, aimed at deter-
29 ring the teenage initiation of sexual intercourse. A total of 130 public
30 school children in grades 5 through 9 participated in 18 eight-week group
31 interventions. The psychoeducational program was theoretically based
32 and designed to be consistent with recommended “best practices” in this
33 area. Outcome measures included a previously published and widely used
34 measure of children’s self-esteem, a knowledge measure of coping with
35 peer pressure, a measure of intention to abstain from sex, and a measure
36 of parent–adolescent communication, all given pre- and posttreatment.

1 All measures showed improvements that were statistically significant,
2 and the authors discussed the strengths and limits of their study (see Abel
3 & Greco, 2008). Abstinence-only sex education is a controversial approach
4 to sex education, and it reflects favorably on the providers of this
5 program that they were willing to subject it to this type of evaluation.
6 The basic elements of the pretest–posttest design in the above example
7 were augmented by using several, not just one, reliable and valid outcome
8 measures. Your conclusions are obviously strengthened if you use several
9 measures of what you are trying to change and all these measures consis-
10 tently point, posttreatment, to changes in the desired direction.

11 **Evaluating Cognitive Behavior Therapy for Depression**

12 Chen, Jordan, and Thompson (2006) used this design to evaluate the
13 possible outcomes of CBT on depression among 30 clients receiving
14 intensive outpatient services at a psychiatric hospital in a Texas city.
15 Clients completed standardized self-report measures of depression and
16 of their problem-solving ability at the beginning of outpatient treatment
17 and again when the daily group therapy program (each session lasting
18 about 2.5 hours per day!) was concluded. Posttreatment depression
19 scores significantly improved, compared to the group's pretreatment
20 mean scores, when pretreatment problem-solving ability scores were
21 controlled for (this is a methodological twist we need not get into here).
22 This is a modest study, but one that can and should be used as an instruc-
23 tive example by social workers seeking guidance on how they can inte-
24 grate simple research methods to evaluate their own practice.

25 **Evaluating School Social Work**

26 Like the posttest-only design, the pretest–posttest single-group design
27 can be improved by using more than one reliable and valid outcome
28 measure and by conducting repeated posttreatment assessments. This
29 was done by Diehl and Frey (2008), in the evaluation of a community–
30 school model of social work practice. The study involved 12 schools
31 located in one school district in the Midwest. The intervention consisted
32 of referring youth with behavior problems to the local school social
33 worker, who provided a standardized system of case management and
34 direct intervention involving individual counseling and home visits,

1 sometimes with group and family therapy sessions added. Upon referral,
 2 the child completed a standardized problem behavior scale, and the par-
 3 ents and teachers completed forms measuring their behavioral concerns
 4 about the child. The program evaluation occurred for all kids referred
 5 from August 2000 to December 2002, a total of 154 youth. Assessments
 6 were completed not only at intake but also 3 and 6 months following
 7 treatment, resulting in the following design:

$$8 \quad O_1 - X - O_2 - O_3$$

9 There were statistically significant improvements in the kids' behavior at
 10 3 months (O_2) and 6 months (O_3) compared to the children's assessment
 11 when they were initially referred (O_1). This is good thing, overall. However,
 12 there were some problems. Data were not available on all 154 kids at the
 13 3-month follow-up assessment. In fact, it was only available for about half
 14 of them, and for only 29% at 6 months. Can you see how this problem of
 15 *client attrition* complicates making any conclusions about the possible
 16 effectiveness of the school social work program? We will return to this
 17 issue of attrition, also known as *mortality*, later on in this book.

18 Evaluating Internet-based Treatment for PTSD

19 Knaevelsrud and Maercker (2007) conducted a randomized controlled
 20 trial of Internet-based treatment for 96 participants with PTSD. About
 21 half ($n = 41$) were randomly assigned to a 5-week, ten-session cognitive-
 22 behavioral writing program that included exposure work, social shar-
 23 ing, and cognitive reappraisal. The remaining participants were initially
 24 assigned to a waiting list control condition. Three distinct measures were
 25 given at the pretest: a measure of PTSD symptomatology (the Impact of
 26 Event Scale), a measure of mood (the Brief Symptom Inventory), and a
 27 measure of health (physical and mental). After the treated clients com-
 28 pleted their program, they and the waiting list control participants were
 29 reassessed using the same measures, and it was found, basically, that the
 30 treated group had marked improvements on all measures except physical
 31 health, whereas the waiting list control group did not change much.
 32 These results clearly favored the Internet-based treatment. At this point,
 33 the initial phase of the project was completed, and the wait-listed clients

1 were provided the Internet-based therapy. Knaevelsrud and Maercker
2 (2010) later went back and recontacted the original group of 41 clients
3 who were initially treated and asked them to take the assessment mea-
4 sures again at 3 months posttreatment and for a fourth time 18 months
5 posttreatment. A total of 34 of the original 41 clients were assessed at
6 18 months (an attrition of 7 out of 41, or 17%). This later aspect of their
7 study could be diagrammed as follows:

$$8 \quad O_1 - X - O_2 - O_3 - O_4$$

9 with the posttreatment assessments corresponding to the initially treated
10 group being evaluated immediately after treatment (O_2), 3 months after
11 treatment (O_3), and 18 months after treatment (O_4). Because these fol-
12 low-up data are reported for only the one group, the initial randomized
13 controlled trial had been broken down to a quasi-experimental pretest-
14 posttest study with repeated follow-up assessment. It is still a good study,
15 especially since multiple outcome measures were used and the results
16 were consistently favorable during all posttreatment measures, even
17 18 months after treatment was discontinued. Many studies only under-
18 take a single posttreatment evaluation of client functioning, usually
19 immediately after treatment, so having three evaluations, with the rela-
20 tively lengthy follow-up period of 18 months, is a real strength. Providing
21 services for persons with PTSD via the Internet appears to provide a
22 promising approach for extending the availability of therapy beyond the
23 confines of the consulting room, which would be helpful.

24 **Post-Psychiatric Hospitalization Follow-up of Adolescents**

25 Welner, Welner, and Fishman (1979) conducted a follow-up study of the
26 disposition of 77 adolescents (average age at treatment was 16 years) who
27 had been psychiatrically hospitalized. Follow-up was conducted some
28 8–10 *years* after discharge. Subgroup analyses were completed, looking
29 at rates of suicide and rehospitalization among patients with various ini-
30 tial diagnoses. These found that adolescents with a diagnosis of bipolar
31 illness fared particularly poorly, as did those with adolescent-onset
32 schizophrenia. This study, whose third author was a social worker, was
33 not so much an evaluation of the effects of treatment as a study in the

1 prognosis of serious mental illness among the young. Still, it is a good
 2 example of answering Question 1, and took place over a remarkably long
 3 time frame. It appeared in one of the premier psychiatric journals.

4 There are some other ways to enhance the usefulness of the pretest–
 5 posttest single-group design. One is to use more than one pretreatment
 6 assessment, as in:

$$7 \quad O_1 - O_2 - X - O_3$$

8 Having more than one pretest provides more credible evidence regarding
 9 the clients' problem status prior to intervention, just as having more than
 10 one posttest enhances our sense of the client's long-term status after
 11 treatment. However, most researchers using the pretest–posttest design
 12 only make use of one pretest and one posttest assessment, not several. If
 13 you employ a whole series of pre- and posttests, your design begins to
 14 segue into another conceptually similar evaluation method called the
 15 *time series design*, which will be covered in Chapter 4.

16 Another refinement in trying to more legitimately figure out the
 17 effects of a given treatment is to add additional elements to your basic
 18 $O_1 - X - O_2$ design through use of a subsequent period of time during
 19 which treatment is *removed* and the client's status assessed again. This
 20 can be diagrammed as:

$$21 \quad O_1 - X_{\text{introduced}} - O_2 - X_{\text{removed}} - O_3$$

22 In terms of logical inference, this design modification only makes sense if
 23 you have reason to believe that the effects of X are likely to be temporary.
 24 Here is an example. A group of kids (say $n > 20$) diagnosed with hyperac-
 25 tivity disorder is assessed using a reliable and valid measure of hyperactive
 26 behavior. They are then treated early the next day with a short-lived medi-
 27 cation intended to reduce hyperactivity; in the middle of the day, they are
 28 reassessed as before. The following day they do *not* receive the medicine for
 29 hyperactivity, and in the middle of that day, they are assessed a third time.
 30 The logic of this design is that if the drug really improves behavior, this will
 31 be evident at the second assessment, O_2 . If positive effects are indeed
 32 observed at O_2 for the group of kids as a whole, this is tentative evidence
 33 that the drug had its intended effect. If it is true that the drug caused

1 the kids to improve, and the drug is deliberately withheld the next day, at
2 the third assessment, their behavior should return to close to that
3 observed at O_1 . If such outcomes are forthcoming at O_3 , you have stronger
4 logical grounds for concluding that X (the drug) was causally responsible
5 for the changes in comportment. This is pretty nifty in terms of trying to
6 make causal inferences, but not a good result clinically, of course.

7 There are many psychosocial interventions you can conceive of that
8 would have immediate but short-lived effects and could be amenable to
9 being studied in this manner. For example, many classrooms use some
10 sort of point system to reward good behavior and deter misbehavior.
11 Do these really work as intended? The question is not a moot one. Point
12 systems take time to devise, implement, and run. If they work well, they
13 may be a blessing. But if they are not really useful, then why go to all the
14 bother? So, imagine the following study. A class of kids is assessed during
15 one day with no particular program in place to encourage on-task aca-
16 demic behavior. The next day, a token or point system is explained and put
17 into place, and at the end of the day, rewards are provided to the kids who
18 performed well. Behavior is assessed all during this day as in the previous
19 day. The third day, the point system is not in operation, and behavior is
20 assessed again. If you saw improvements during the day when the point
21 system was in place, improvements over the first day's behavior, and a res-
22 toration to the original level of functioning on the third day, most teachers
23 (and social workers) would be pretty convinced that the point program
24 was effective in promoting on-task behavior and reducing misbehavior.
25 You may be able to think of other psychosocial interventions that could be
26 evaluated using this type of design. For example, the use of a token econ-
27 omy on the behavior of psychiatric inpatients, the effects of daily atten-
28 dance at Alcoholics Anonymous meetings on one's craving to drink, the
29 effects of a day treatment program on the mood of people with clinical
30 depression, and more. Basically, any intervention whose effects are expected
31 to be immediate but *temporary* can be evaluated using this approach.

32 **SUMMARY ON STRENGTHENING THE ONE-GROUP** 33 **PRETEST–POSTTEST DESIGN**

34 All those methods listed above for use in strengthening the posttest-only
35 design apply here. In addition, you may do the following.

1 Use Multiple Pretreatment Assessments

2 Client problems may wax and wane in response to the ebb and flow of
 3 naturally occurring life events, via their own natural history (e.g., bipolar
 4 disorder, depression), or in response to biological changes occurring
 5 within clients (e.g., changes in diet, sleeping, exercise, medication use,
 6 use of food supplements and herbal products, menstrual cycle, etc.). This
 7 can complicate making any inferences about meaningful changes follow-
 8 ing receipt of social work services. Using two or more preassessment
 9 periods permits a more informed appraisal of true changes, relative to
 10 measuring a group at a single point in time. If problems are seen to be
 11 tending upward or downward during these pretreatment assessments, or
 12 if they are relatively stable, one is better able to detect any true changes by
 13 comparing the several pretreatment values with the posttreatment one
 14 (or several posttreatment evaluations, which is even stronger). Clearly a

$$15 \qquad \qquad \qquad O_1 - X - O_2$$

16 design is greatly improved upon by using the

$$17 \qquad \qquad \qquad O_1 - O_2 - O_3 - X - O_4 - O_5 - O_6$$

18 For example, suppose in the first example immediately above (the
 19 $O_1 - X - O_2$ Design), the average pretest score was 60 and the average post-
 20 test score was 40 (the meaning of these numbers is unimportant here),
 21 with higher scores meaning greater problems. This would look pretty
 22 good, with problems decreasing from 60 to 40 points following treatment.
 23 It would indeed look good if this was all the information you had. But
 24 suppose you used the second design, and the three average pretest scores
 25 were 100, 80, and 60, with the three average posttest scores being 40, 20,
 26 and 0. Would having this additional information change how you viewed
 27 the originally presented average pretest score of 60 and posttest score of
 28 40? Most likely it would, since you could see that the three pretreatment
 29 scores were tending downward, and the posttest scores simply reflected
 30 this pretreatment trend extended over time. In other words, in the second
 31 instance, it looks as if treatment had *no effect*. This possibility could only
 32 be ascertained by using multiple pretests and posttests.

1 Use a Removal or Withdrawal Phase

2 During a removal or withdrawal phase, clients' are reassessed after the
3 treatment is deliberately removed. This approach is viable when the
4 effects of treatment are short-lived. Assess the group, treat the group,
5 reassess. If you see a positive change, terrific. Then, withhold the treat-
6 ment and reassess. If you see a relapse when treatment is not provided,
7 your ability to infer that the original positive changes were the results of
8 treatment is enhanced. Not perfect, just enhanced. And, of course, it is
9 not desirable to terminate a study at this point, but you may, if feasible,
10 reinstate the original treatment condition on an enduring basis.

11 Use These Designs Prospectively

12 Use these designs prospectively; that is, in a way that is planned for and
13 with data gathered in advance, as opposed to using them *retrospectively*—
14 looking at data after the fact, without having had any plans to use the
15 data for evaluation services when the data were initially gathered.

16 These designs are widely used across the human services and health
17 professions, and regularly appear in some of the world's most prestigious
18 journals. It is ill-informed to simply dismiss them as unworthy of consid-
19 eration because of their frequently low internal validity.

20 Some additional social worker examples of using these designs can be
21 found:

- 22 • Raskin, Johnson, and Rondestvedt (1973) evaluated the
23 pretest–posttest changes that occurred among ten patients
24 who suffered from chronic anxiety, and who were treated with
25 muscular relaxation-based biofeedback. Outcome measures
26 included a measure of anxious mood, sleep difficulties,
27 headaches, and sense of relaxation. This simple study appeared
28 in a leading journal, the *Archives of General Psychiatry*.
- 29 • A novel form of comprehensive care for troubled youth has
30 been developed, known as *wrap-around services*. Wrap-around
31 is a family-centered and strengths-based program involving
32 interdisciplinary team treatment with elements of educational,
33 mental health, child welfare, and juvenile justice services being
34 provided in a coordinated manner. Many hundreds of millions
35 of dollars have gone into the federal funding for wrap-around,

1 with relatively little evidence that it really helps kids. Copp,
2 Bordnick, Traylor, and Thyer (2007) conducted a pretest–
3 posttest evaluation of the first 15 youth who received wrap-
4 around in a pilot program initiated in Georgia, for whom
5 complete baseline and 6-month follow-up information was
6 available. No significant changes were found on any of the
7 outcome measures. This was disappointing. As the authors
8 pointed out, these negative results may have been due to the
9 relatively small sample size and to some concerns that the
10 newly initiated wrap-around model may not have been
11 implemented as recommended.

- 12 • Holly Bragg Capp, a MSW intern, arranged to have newly
13 admitted inpatients on a psychiatric unit complete a reliable
14 and valid measure of psychiatric symptoms, and for them to
15 complete the same measure again, just prior to their discharge.
16 Over the course of one semester, she was able to obtain pretest–
17 posttest data on 78 consecutively admitted patients. After
18 examining the data, she was able to conclude that the patients’
19 reported symptoms improved both statistically and clinically.
20 Very few programs can answer the question “Did our patients
21 improve?” with empirical data, and Holly’s study was a nice
22 step in providing a preliminary answer to this question.
23 Can she assert that they improved because of the inpatient
24 program? No. But still it is useful to undertake simple studies
25 like this. Certainly, the program’s administrators were happy
26 with the result! See Capp, Thyer, and Bordnick (1997) for a
27 full report of this project.

28 Some non–social work examples of using the pretest–posttest design
29 include the following:

- 30 • Spinelli (1987) used a pretest–posttest design with 13 participants
31 to evaluate possible effects of interpersonal psychotherapy on a
32 group of clients with which it had not been previously used. This
33 study was published in one of the world’s leading psychiatric
34 journals, the *American Journal of Psychiatry*.
- 35 • Whitt-Glover, Hogan, Lang, and Heil (2008) evaluated a faith-
36 based physical activity program to promote exercise among

- 1 sedentary African Americans. This appeared in a well-respected
2 public health journal.
- 3 • Schneider et al. (2006) evaluated the outcomes of a large-scale
4 community-based program aimed at promoting the
5 consumption of fresh fruits and vegetables among school
6 children. This appeared in the prestigious *Journal of the*
7 *American Medical Association*.
 - 8 • Novak, Cusick, and Lowe (2007) evaluated the impact of an
9 occupational therapy program provided in the homes of
10 Australian children with cerebral palsy.
 - 11 • Keller et al. (2009), psychologists, evaluated a statewide suicide
12 prevention program delivered to 416 providers of child welfare
13 services.

14 Clearly, these designs have usefulness in evaluating a wide array of
15 social care and health services, and, if well-conducted and circumspect in
16 drawing causal inferences, are capable of being published in very presti-
17 gious professional journals. They should not be cavalierly dismissed as
18 useless or unimportant.

19 SOME LIMITS OF THE PRE-EXPERIMENTAL DESIGNS

20 You will probably have noticed that even the best of the pre-experimen-
21 tal designs are generally only capable of answering Questions 1 and 2
22 from those introduced in Chapter 1. They are not very good at answering
23 the other three questions. This is because Questions 3–5 all involve trying
24 to compare outcomes from a group of clients who received an interven-
25 tion of interest with the outcomes of clients who received no treatment,
26 standard treatment, or placebo treatment. Because the pre-experimental
27 designs have no other groups to evaluate treated clients against, any
28 comparisons between groups are not possible. It is this feature, the lack
29 of any comparison or control groups, which has traditionally separated
30 the pre-experimental designs from the so-called quasi-experimental
31 designs. The latter designs *do* involve comparison groups. But the quasi-
32 experimental designs do not rise to the level of sophistication of true
33 experiments because they lack a further refinement; namely, the deliber-
34 ate random assignment of clients to various conditions or groups.

1 You can use the bulleted points below to help distinguish among these
2 three types of nomothetic research designs, as they are often traditionally
3 construed:

- 4 • *Pre-experimental designs*. Looks at a group of clients
5 posttreatment only, *or* compares posttreatment outcomes with
6 pretreatment observations, obtained from the same group.
- 7 • *Quasi-experimental designs*. Compares a treated group against
8 other groups of clients who received no treatment, standard
9 treatment, another treatment, or placebo treatment.
- 10 • *Experimental designs*. Includes the elements of pre- and quasi-
11 experimental designs, *plus* the added feature of creating these
12 groups through random-assignment methods.

13 Do keep in mind that, in this book, pre-experimental and the tradi-
14 tional quasi-experimental designs are collectively referred to as quasi-
15 experimental designs. A specific type of experimental design used in
16 more tightly controlled social work intervention research, the *random-*
17 *ized controlled trial* (RCT), is covered in another volume in this series
18 (Solomon, Cavanaugh, & Draine, 2009), as well as in a number of other
19 similarly excellent books (Shadish et al., 2002; Nezu & Nezu, 2008), and
20 is often covered in separate chapters in general social work research text-
21 books (e.g., Cnaan & Tripodi, 2010; Pignotti & Thyer, 2009). The virtue
22 of RCTs resides in their stronger potential ability to permit true causal
23 inferences, to be able to say with some degree of confidence that a given
24 effect was the result of a given treatment. This is obviously important
25 because if we cannot sort out legitimately effective interventions from
26 ineffective ones, those that fail to produce positive effects above and
27 beyond those induced by placebo effects, or those that are harmful to
28 clients, in some ways the rationale for the existence of the profession of
29 social work is called into question. Quasi-experimental designs are an
30 attempt to build upon the methods of the pre-experimental designs by
31 adding some type of control groups or conditions, for the purposes of
32 strengthening our ability to test causal hypotheses of a more complicated
33 nature than those implied in Questions 1 and 2.

34 For example, take something fairly new to our field, called *narrative*
35 *therapy*. The specifics of it are not important. It would be a nice thing to
36 show that clients who received narrative therapy were very satisfied with

1 their treatment, or were functioning very well. But this does not prove
2 that narrative therapy is genuinely effective. It would be an improve-
3 ment to show that clients who received narrative therapy were better off
4 after narrative therapy than they were before receiving this treatment.
5 But even this level of evidence is insufficient to prove that narrative ther-
6 apy is genuinely effective. However, if we could show that clients who
7 received narrative therapy were truly better off afterward, compared to a
8 group of clients who got no treatment at all, this would greatly add to the
9 credibility of narrative therapy as an effective intervention. And, if we
10 could show that narrative therapy produces better results than treatment
11 as usual, or compared to a credible placebo treatment, this would be
12 better still. And this is the purpose of quasi-experiments, to conduct such
13 studies to provide stronger evidence of the effectiveness (or ineffectiv-
14 eness) of various social work interventions. Ideally, intervention research-
15 ers maintain a studiously neutral stance by not attempting to prove that
16 something does or does not “work.” Rather, they aim to empirically
17 determine a treatment’s effectiveness, whatever the outcomes. Seeking
18 after “truth” is the foremost agenda, not supporting one’s preferences.

19 The ability to make legitimate causal inferences requires more than
20 simply demonstrating change. One must have some confidence that
21 other factors potentially responsible for these changes have been effec-
22 tively ruled out, and it is for this purpose that quasi-experiments make
23 use of various comparison or control groups. Next, we will review some
24 of the more common reasons, apart from Treatment X, which might
25 be responsible (at least in part) for client improvement. These alternative
26 explanations are collectively called *threats to internal validity*. We will
27 then examine how various quasi-experimental designs try to control for
28 these confounding effects.

29 **SOME THREATS TO INTERNAL VALIDITY**

30 **Placebo Effects**

31 Placebo effects can be defined as:

32 Any therapy or component of therapy (or that component of any ther-
33 apy) that is intentionally or knowingly used for its nonspecific, psycho-
34 logical, or psychophysiological effect, or for its presumed therapeutic

58 Quasi-Experimental Research Designs

1 effect on a patient, symptom, or other illness but is without specific
2 activity for the condition being treated. (Bausell, 2007, p. 29)

3 or, similarly,

4 Changes in a dependent variable that are caused by the power of sugges-
5 tion among the participants in an experimental group that they are
6 receiving something special that is expected to help them. These changes
7 would not occur if they received the experimental intervention without
8 that awareness. (Rubin & Babbie, 2008, p. 642)

9 Placebo effects can be assumed to be present to some degree when-
10 ever a client believes he or she is receiving a credible treatment or inter-
11 vention of some sort. The credibility of an intervention may involve not
12 only the features of the treatment itself, but also the appearance, manner,
13 and confidence of the social worker providing the therapy; the physical
14 surroundings associated with treatment (a well-appointed office in a nice
15 part of town vs. shabby digs on the wrong side of the tracks); the legiti-
16 mate and perhaps not-so-legitimate degrees, credentials, and certifica-
17 tions possessed by the social worker; the therapist's reputation; and more.
18 For two different treatments to be legitimately compared in terms of
19 their effectiveness, both interventions must possess equivalent credibil-
20 ity. If treatment X seems highly credible, and treatment Y much less so,
21 a study comparing X and Y is automatically biased in favor of X from the
22 onset of the study. Hence, well-designed outcome studies attempt to
23 ensure that the placebo-engendering features surrounding X and Y are
24 roughly the same. It is sad but true that if therapy X is provided by
25 Dr. George Clooney, and therapy Y by Dr. Quasimodo, clients may
26 respond as much to *who* provided the treatment, as to the essential fea-
27 tures of the treatment itself.

28 Regression to the Mean

29 It is not uncommon for clients to seek treatment, even to participate in
30 clinical trials of psychosocial intervention, when their problems are at their
31 peak. Many of the psychosocial and health disorders for which social work
32 clients seek assistance have a natural tendency to wax and wane on their
33 own. This is obviously true for conditions such as moderate depression

1 and bipolar disorder, but also true for spousal abuse, alcoholism, and
2 schizophrenia. This ebb and flow of symptom severity can be said to fluctuate
3 around a general trend or average (mean) level. Thus, clients who
4 enroll in treatment studies when their problems are particularly bad may
5 experience a lessening of symptom severity over the ensuing weeks or
6 months, an apparent amelioration that has little to do with the actual
7 effects of treatment but more to do with the natural history or progression
8 of their difficulties. This reversion back to the more general level of severity
9 is called *regression to the mean*. In these instances, it is natural for such
10 individuals to often ascribe their improvements to their participation in a
11 treatment and, if actual improvements are measured, it is equally natural
12 for therapists, even those participating in research projects, to similarly
13 attribute these improvements to treatment, rather than to natural (tempo-
14 rary) remission. This however, commits the logical fallacy of *post hoc ergo*
15 *propter hoc*, Latin for “after this, therefore because of this,” by reasoning
16 along the lines of “Since that event *followed* this one, that event must have
17 been *caused* by this one.”

18 **Maturation**

19 *Maturation* refers to developmental changes that occur over time in some
20 client populations. Young children, adolescents, and the elderly are client
21 groups in which maturational changes are particularly salient. In outcome
22 studies transpiring over long periods—months, perhaps even years—
23 clients may change rather dramatically for reasons that have nothing to
24 do with their receipt of therapy. Children may make striking advances in
25 social or cognitive development in a surprisingly brief period, and the
26 elderly may experience rapid changes in the opposite direction, toward
27 more impairment, cognitive abilities, or senescence.

28 **Passage of Time**

29 A number of the specific threats to internal validity mentioned thus
30 far involve the element of time—time is required to elapse before fac-
31 tors such as maturation, regression to the mean, concurrent historical
32 influences, mortality, etc. become potentially operative. There is another
33 element in which time poses a threat, and it refers to simple changes
34 induced by one’s experience with a condition or disorder. Over time,

1 the morale-eroding influence of a problem may become more or less
2 impactful. Think of common phrases of everyday speech such as *time*
3 *heals all wounds*, *the tincture of time*, or *one day at a time*. With time, one's
4 craving for cigarettes may decline after stopping smoking. The pain of
5 being newly unemployed may diminish, or, conversely, the misery
6 of living with severe obsessive-compulsive disorder may lead to thoughts
7 of suicide. One may adapt to living with a verbally abusive spouse by
8 withdrawing, and so reduce the psychological anguish of feeling unloved.
9 These are all changes that can operate in the context of an evaluation
10 study and give rise to improvements or deterioration that may superfi-
11 cially look like the results of an intervention but are actually a function of
12 a far simpler accounting—the mere passage of time alone. It is not bio-
13 logically mediated maturation, nor the statistical artifact of regression,
14 just time alone.

15 **Mortality/Attrition**

16 Although the term *mortality* usually connotes death, in social work inter-
17 vention research it more often refers to the problem of clients dropping
18 out of treatment. If you begin a simple pretest–posttest study with a
19 client sample of 100 and 3 months later, following intensive treatment,
20 you only have 70 clients continuing to participate in the study, the threat
21 to internal validity called mortality may be present. Suppose you had an
22 outcome measure that consisted of scores on a rapid assessment instru-
23 ment, and you looked at the group's mean pretest score and compared it
24 with the group's mean posttest score. The first mean was based on the
25 original sample of 100 clients, whereas the second was based only on the
26 remaining 70 clients still participating in your study. It is possible that
27 the 30 folks who dropped out differed in meaningful ways from the 70
28 “survivors.” Perhaps the 30 had more severe problems? If so, the mean
29 for the remaining 70 will display improvements at the posttest assess-
30 ment *not* because of the effects of participating in therapy, but because
31 the mean score for the remaining clients is no longer dragged down
32 by the scores of the more seriously impaired ones who dropped out.
33 But mortality may not be due to symptom severity. More prosaic prob-
34 lems may be operative, such as a lack of access to reliable transportation
35 among the clients who dropped out, or a lack of adequate child care
36 needed for clients to attend clinic appointments. It can be difficult to

1 know why people dropped out of a study. Even if you compare demo-
2 graphics and symptom severity of the dropouts versus survivors and find
3 no difference, there may be undetected or unassessed variables at work
4 that are responsible for mortality. Thus, this would make it unwise to
5 claim that, since the clients who dropped out were similar to those who
6 survived, you can safely ignore your study's high mortality as a potential
7 threat to internal validity. A practical example of this problem occurred
8 in the study evaluating school social work services conducted by Diehl
9 and Frey (2008) described above.

10 **Differential Attrition**

11 Imagine you are comparing the effects of two treatments, X and Y.
12 X makes many demands on clients and is psychologically stressful,
13 a treatment such as primal scream therapy (PS, invented by a social
14 worker, Arthur Janov). Treatment Y, on the other hand, is much less
15 upsetting, say, narrative therapy (NT, invented by a social worker!).
16 In fact, Y is so reinforcing that the clients look forward to their sessions.
17 Say you began your study with two groups receiving these two treat-
18 ments, with 50 clients in each group. After some months, 20 of the 50
19 folks receiving PS therapy had dropped out, whereas only three of those
20 enjoying NT stopped participating. At the end of the study, instead of
21 comparing 50 clients to 50 clients, you compared 30 clients to 47 and
22 found that, on average, those who received NT were dramatically better
23 off compared to those who got PS therapy. Could you legitimately con-
24 clude that NT is a more effective treatment than PS? Not really, since the
25 confounds associated with mortality noted above are now differentially
26 present between your two groups, and these complicate your ability to
27 make causal inferences.

28 **Concurrent History**

29 The threat known as *concurrent history* refers to impactful events taking
30 place in clients' lives outside the context of their participation in a treat-
31 ment outcome study. Sometimes these can be very conspicuous things,
32 such as 9–11, Hurricane Katrina, or the election of a very popular presi-
33 dent. Such events can broadly affect the mood and well-being of research
34 participants, and these effects may be reflected in your posttreatment

1 outcome measures and complicate your ability to make any inferences
2 about the effects of treatment per se. Macdonald noted this threat early
3 on, urging us to “take cognizance of the possibility that the patients might
4 have improved if left untreated. . . . We should never lose sight of the
5 possibility that the improvement of the patient is not related to our
6 efforts” (Macdonald, 1952, p. 137).

7 **Nonblind Assessments**

8 If outcome measures used in a treatment study involve the judgments of
9 others, such as therapists’ ratings of their own clients or the ratings of
10 independent observers, and these raters/judges are aware of which phase
11 of a study (pre- vs. posttreatment) is being rated or which treatment con-
12 dition the client received (experimental therapy, treatment as usual, pla-
13 cebo treatment), it can introduce an element of bias into their appraisals.
14 Suppose, as a part of a study on treating depressed clients, the therapists
15 themselves were asked to determine if the clients met the diagnostic
16 criteria for major depression before they began therapy and then again,
17 after the clients’ completion of a treatment program administered by the
18 therapists themselves. You can see that using the treating therapists to
19 determine if the clients still met the criteria for major depression after they
20 were treated would incur a strong tendency for the therapists to be biased
21 in favor of detecting clinical improvements. After all, who would want to
22 judge their own clients as *unimproved* following treatment? A more rigor-
23 ous approach would be to have specially trained diagnosticians assess the
24 clients; these raters would be uninvested in the outcomes of the study,
25 have no allegiance to any particular type of therapy being evaluated, and
26 not be aware of whether the client was beginning treatment or had com-
27 pleted therapy. Knowing that you are assessing someone at the *end* of a
28 treatment outcome study may also create a bias toward a more favorable
29 assessment. Having these independent diagnosticians be unconnected
30 with the actual treatment would help ensure their independence, as they
31 would have less of an investment in detecting improvements.

32 Also, if you have clients receiving different treatments, with one
33 group receiving a novel therapy and a second receiving treatment as
34 usual, if at the time of the posttreatment evaluations the assessors are
35 aware of which treatment the clients had received, this may bias their
36 judgments. Thus, really well designed outcome studies have assessors

1 who do not know if the client they are evaluating is at the beginning or
 2 end of a treatment program, nor do they know which treatment condi-
 3 tion the client received. An even more rigorous technique is to ask, after
 4 the assessments have been completed, the assessors to *guess* if the client
 5 was pre-treatment or posttreatment, or if posttreatment, which treat-
 6 ment they had received (novel treatment, treatment as usual, placebo
 7 treatment, etc.). If the assessors can guess no better than chance, then the
 8 “blindness” of the ratings has been maintained. If they do guess better
 9 than chance, then blindness has been compromised, and the study should
 10 not be reported as having made use of truly blind evaluators.

11 **Multiple Treatment Interference**

12 Sometimes—indeed, most times—psychosocial treatments contain mul-
 13 tiple elements. Most therapies, except perhaps bibliotherapies or Internet-
 14 or computer-based self-help programs, contain important client–therapist
 15 relationship elements, features separate from the specific *techniques* that
 16 may be the manifest focus of treatment. Narrative therapy involves
 17 the telling of stories, solution-focused treatment is based on envisioning
 18 life without a given problem or issue, dynamic psychotherapy involves
 19 recounting childhood experiences, behavior therapy involves arrang-
 20 ing for certain desired behaviors to be reinforced, and so on. And within
 21 each therapeutic model may be found many different components.
 22 A program of behavior analysis could involve reinforcement, modeling,
 23 shaping, manipulation of antecedent stimuli, and the like. Dynamic psy-
 24 chotherapy could involve recounting childhood memories, the interpre-
 25 tation of dreams, the analysis of resistance, and free association. And
 26 sometimes clients receive different treatments at the same time, as in
 27 interpersonal psychotherapy combined with receiving an antidepressant
 28 medication. This receipt of multiple or combined treatments makes it
 29 problematic to assert that only *one* element of a therapeutic regiment is
 30 causally responsible for any observed improvements. In such instances,
 31 one may at best conclude that a given *combination* of treatments was fol-
 32 lowed by certain changes, but one cannot legitimately assert that only one
 33 or more elements of that combination were responsible for these improve-
 34 ments. Such a design could be diagrammed along the following lines:

35 $O_1 - (X + Y) - O_2$

1 with (X + Y) reflecting the fact that clients got two discrete interventions.
2 It is commonly understood, however, that X implies all of the particular
3 therapeutic elements that go into X, given the complexities of social work
4 intervention. Once X has been shown to be genuinely helpful, subsequent
5 studies may be used to isolate the active ingredients of X, compared to
6 the inactive one.

7 Assessment/Treatment Interactions

8 In a treatment study, clients may be formally evaluated prior to begin-
9 ning therapy, then receive therapy, and then be reevaluated. It is both
10 possible and plausible that the act of being assessed *combined with* the
11 subsequent receipt of treatment exerts an effect different from that to be
12 obtained if a client received the treatment alone. For example, if a client,
13 as a part of a weight reduction program was asked, prior to treatment,
14 to weigh himself for a number of days prior to beginning formal treat-
15 ment, and then at the conclusion of the program was weighed again, the
16 act of consciously weighing himself could, *by itself*, motivate him to begin
17 some efforts at losing weight. The outcome of people receiving this com-
18 bination of assessment and treatment would be different relative to a
19 group of clients who got the formal treatment alone, without a formal
20 preliminary period of self-monitoring, and were weighed only at the
21 conclusion of formal treatment. Such self-monitoring effects are not
22 uncommon in the world of therapy, and many times these efforts produce
23 modest changes irrespective of other formal treatments. In treatment out-
24 come studies of people with specific phobias (to say dogs, cats, birds,
25 snakes, etc.), one assessment method is called a *behavioral approach test*
26 (BAT), wherein the phobic person is asked to approach the restrained,
27 feared animal as closely as possible, with the closest distance obtained being
28 one measure of the severity of fear (more phobic persons approach less
29 closely). Repeated BATs themselves can produce mild improvements, even
30 in the absence of formal therapy. Thus, if you designed a study involving a
31 preliminary BAT, then provided a completely useless therapy, followed by
32 a posttreatment BAT, you might see small reductions in avoidance. It
33 would be very tempting to conclude that “treatment” produced these
34 gains, when in reality it was simply the act of being assessed by a BAT.

35 It is also known that individuals taking certain standardized tests, such
36 as the Graduate Record Examination (GRE), tend to improve their scores

1 the second time they take the test. Some commercial firms provide test
2 preparation programs for the GRE and similar exams, a component of
3 which involves taking practice tests very similar to the real examinations.
4 They then report the pass rate or scores of those who took the commercial
5 test preparation program and compare those scores to national test scores
6 (usually favoring the people who enrolled in the test prep course). A more
7 legitimate comparison is to compare the test scores of those who took the
8 test preparation class with the scores of those who simply practiced taking
9 the test an equivalent number of times on their own, not those taking the
10 real test for the first time. Such an evaluation would likely show a much
11 narrower gap in pass rates between those who completed a prep class and
12 those who simply practiced taking the test an equivalent amount.

13 **Instrument Change**

14 The term *instrument* refers to the methods used to assess client func-
15 tioning. This may involve a wide variety of approaches, including client
16 self-reports; ratings of others such as caregivers, teachers, or family
17 members; direct behavioral observations; formal or informal interviews;
18 the client's completion of various tests, rating scales, or rapid assessment
19 instruments, and the like. To compare the posttest scores of one group of
20 clients with their pretest scores requires that the assessment methods or
21 instruments used be similarly conducted, ideally identically, at the two
22 points in time, pre- and posttreatment. If the assessment method differs
23 in some meaningful manner, then you cannot legitimately compare the
24 two sets of scores and make any valid conclusion regarding real changes
25 in client functioning. Similarly, if you wish to compare assessments of
26 two groups of clients, at pretest, at posttest, or both, then not only must
27 the assessments at pretest and posttest be similarly conducted, but they
28 must also be similarly provided *between* the two groups.

29 Instrument change may not be much of a threat when clients do
30 something simple, such as complete an easy-to-read scale or self-rating
31 form. In this case, they can be given the form, read the simple directions,
32 and fill out the form; then the forms are collected and later analyzed.
33 But when more complex assessments are used, the threat of changes in
34 instrumentation may arise. For example, let's say that at the beginning of a
35 study a clinician is asked to conduct a structured diagnostic interview to
36 determine if a client meets the criteria for a particular psychiatric condition

1 such as PTSD. The clinician interviews a large number of clients (say, >30)
2 in the context of a pretest–posttest study. At the beginning of the study, she
3 interviews potential client #1 at pretest, determines that he does meet the
4 criteria for PTSD, and the client is duly enrolled in the study. She then
5 interviews client #2, finds that this person also meets the criteria for PTSD
6 and is eligible to participate in the study, and so on. Eventually, a sufficient
7 number of participants are enrolled and therapy begins. A few months
8 later, therapy is terminated and posttreatment interviews are conducted by
9 the same clinician, who again determines whether or not the client meets
10 the Diagnostic and Statistical Manual (DSM) criteria for PTSD. At the
11 beginning of the study, 100% of the participants were judged to meet the
12 criteria for PTSD; at the end, only 30% were, thus giving an apparent out-
13 come that 70% of the clients no longer “had” PTSD and raising the possi-
14 bility that the treatment is a highly effective cure for this serious condition!

15 One possible confound or threat to internal validity here is the
16 clinician-diagnostician’s changes in diagnostic skill over the course of the
17 study. It is very likely that her interviewing skills and ability to apply the
18 DSM criteria with the first five clients assessed pretreatment were consid-
19 erably different from her skills applied to the last five clients assessed
20 posttreatment. Simply put, she may have experienced a considerable
21 improvement in her diagnostic acumen during the course of applying
22 the structured clinical interview over 60 clients. Practice may make one
23 perfect. Or, she might have gotten bored with the repetitiveness of the
24 process and began to cut some corners by leaving out some questions,
25 resulting in a less accurate diagnostic determination. In other words, the
26 “instrument” used to assess clients (the diagnostician) changed over the
27 course of the pretreatment and posttreatment assessment process, and
28 the possibility exists that it is this confounding factor that resulted in the
29 apparent decline in the numbers of clients meeting the criteria for PTSD,
30 *not* the curative powers of the presumptive therapy.

31 How can this confound be controlled for? One approach is to train
32 evaluators to some criterion prior to really using their appraisals in a
33 study. Don’t use someone new to a complex assessment method—use
34 highly experienced people. Another approach, when using human raters,
35 interviewers, or observers, is to use *two* independent evaluators and cal-
36 culate the extent to which they agree. For example, in the PTSD study
37 just described, the potential research participant could be evaluated by
38 two clinicians acting independently of each other. Each determines if

1 potential client #1 did or did not meet the DSM criteria for PTSD, and
2 only those clients who were determined by *both* raters acting indepen-
3 dently are enrolled in the study. And, when therapy was completed, it
4 required *both* raters to judge the client as no longer meeting the DSM
5 criteria for PTSD in order for the client to be tabulated in the “cured”
6 category. This raising of the methodological bar makes for better science,
7 but it is costly in terms of time and money.

8 In a study on changes in clinical interviewing skills found among
9 MSW students taking a social work methods course, Carrillo, Gallant, and
10 Thyer (1993) used two independent raters to judge the MSW students use
11 of facilitation, questioning and clarification, and support and empathy
12 during structured interviews with a simulated client. The raters’ scores
13 had very high reliabilities, and it was found that, by the end of the class,
14 the MSW students’ scores on these skills had improved. This sounds good,
15 but something problematic happened in the course of the study. The sim-
16 ulated client used in the mock interviews at the beginning of the term
17 became unavailable, and a different person had to be recruited to portray
18 a client at the end of the term, for the posttraining mock interview.
19 The improvements observed at the end of the class *may* have reflected not
20 only genuine enhancements in the students’ interviewing skills, but also
21 changes in the ease with which the different simulated client used at the
22 end of the class could be interviewed. It is possible that this second, differ-
23 ent person was simply easier to interview compared with the person por-
24 traying the mock client at the beginning of the class. In other words, the
25 instrument (in this case, the simulated client) used to assess client func-
26 tioning (in this case, the MSW students) was significantly changed, making
27 it very difficult to strongly claim that taking the class really improved
28 interviewing skills. Of course, that is what the authors would like to have
29 believed, but they recognized that this change in the person serving as the
30 simulated client introduced the possibility that their students’ improve-
31 ments were the results of changes in instrumentation, and so noted this in
32 the published article. This is honest reporting, but simply recognizing the
33 problem as a possible confound does not adequately control for it.

34 **Differing Treatment Credibility**

35 Some therapies seem, on their face, to be highly credible and indeed
36 make sense. The gradual real-life desensitization of phobic persons to

1 their feared object, animal, or situation is one example of a common-sense–
2 based therapy. The degree of credibility people perceive with respect to
3 the treatment they receive in the context of a treatment outcome study
4 can affect clients' responses to the treatment. If, in the name of therapy,
5 you are asked to do something that makes no sense, or may even seem
6 silly, the crucial element of the positive placebo effect is reduced and less
7 improvements may be forthcoming. If you are exposed to a treatment
8 that seems highly credible, then placebo effects are maximized and
9 greater improvements may be elicited, relative to those caused by less
10 credible therapies. Really good psychotherapy outcome studies ask cli-
11 ents at the beginning and end of treatment how much they expected
12 to benefit from the treatment they were going to, or had, received.
13 Or, clients are asked how “believable” the treatment seems/seemed to
14 them. To make a legitimate comparison of the true effects of treatment
15 X versus treatment Y, X and Y should have equivalent credibility or
16 believability on the part of the clients. Otherwise, you are not comparing
17 the true effects of X versus Y, but rather, for example, a seemingly legiti-
18 mate treatment X versus a silly-seeming treatment Y, and this is not a
19 genuine comparison of the real effects of X and Y, relative to each other,
20 independent of credibility issues.

21 The costs of treatment, oddly enough, may affect how effective ser-
22 vices may be. It has long been believed that clients who pay money for
23 psychotherapy services tend to benefit more than do clients who receive
24 free services, a belief stemming from both psychoanalytic and cogni-
25 tive dissonance theory (Shipton & Spain, 1981; Wood, 1982, Yoken &
26 Berman, 1984; Herron & Sitkowski, 1986). And, paying higher fees
27 induces a greater benefit than does paying lower fees. Recently, Waber,
28 Shiv, Carmon, and Ariely (2008) exposed healthy volunteers to a painful
29 task, receiving a series of electric shocks to the wrist of increasing inten-
30 sity, and had them rate the pain induced by each individual shock received.
31 All participants were given a pill before the task and, when given the pill,
32 half of the subjects were informed that it cost of \$2.50 and the other half
33 were told that it cost 10 cents. The pill was supposedly an analgesic (pain
34 reliever) but in reality it was a placebo. The researchers examined the pain
35 intensity ratings of the two groups and found that those subjects who
36 received the pill supposedly costing \$2.50 reported significantly lower
37 pain ratings than did those who received the pill costing 10 cents. In other
38 words, the person's perceptions of the cost of a treatment apparently

- 1 affected how well the treatment performed. Similar influences may effect
- 2 the outcomes of psychosocial interventions as well.

3 Selection Bias

4 There are two ways in which selection bias may complicate the results of
5 your evaluation study. In the first, you may conduct a study on a group
6 of clients that possess characteristics that render them particularly
7 susceptible to either positive or negative responses to a given interven-
8 tion, leading to a false conclusion about the treatment's presumptive
9 positive (or negative) effects. If your sample consists of particularly high-
10 functioning persons, the results may not be generalizable or relevant to
11 the larger population of people with a particular problem. Hence, your
12 conclusion that therapy X helps clients with that problem may be incor-
13 rect. For example, if a program that was intended to help the unemployed
14 find work only enrolled persons with a college degree into an interven-
15 tion study, and it was later found that a very high proportion did indeed
16 obtain good jobs, this would not mean that the program was an effective
17 method for helping the vast majority of persons (who are less well edu-
18 cated) obtain work. Conversely, if a program of assertive community
19 treatment (ACT), a complex intervention requiring intense daily moni-
20 toring by a treatment team of the functioning of persons with chronic
21 mental illness, only was tested on persons meeting the criteria for schizo-
22 phrenia of the paranoid subtype, the intervention may meet with great
23 resistance, given the particularities of that form of schizophrenia.
24 The negative result may lead to a conclusion that ACT does not work,
25 when in reality it does not work with this particular type of client, but
26 may well be very effective for persons with other forms of schizophrenia.

27 Selection bias also refers to the possibility that two or more groups in
28 a quasi-experimental study may have unrecognized preexisting treat-
29 ment differences that only became evident following treatment, with
30 these differences then being incorrectly ascribed to the differing impacts
31 of the treatment conditions. Suppose a social work researcher wanted to
32 determine if clients who received group therapy at her agency improved
33 as much as those who received individual treatment. In this hypothetical
34 example, the normal agency practice is for an intake clinician to conduct
35 all initial evaluations of new clients and to suggest treatment options
36 (in this case group or individual therapy), with most clients following the

1 recommended course of action. Clients begin treatment, say for 3 months
 2 of *either* individual or group therapy, and are then formally assessed
 3 using one or more reliable and valid outcome measure(s). This design
 4 could be diagrammed as:

5 $X - O_1$

6 $Y - O_1$

7 With X indicating clients who received individual therapy, Y being those
 8 clients who received group therapy, and O_1 the posttherapy assessment
 9 scores.

10 The social work researcher gathers the data and finds out that clients
 11 who received group therapy had much higher scores on functioning com-
 12 pared with those clients who received individual treatment. This could
 13 lead an unsophisticated researcher to the conclusion that group therapy
 14 was more effective than individual care. But such a conclusion would
 15 not take into account the possibility of selection bias in the composition
 16 of the two groups. It is possible that the intake clinician unknowingly
 17 (or perhaps knowingly) assigned more seriously disturbed clients to
 18 receive individual treatment, believing that this was somehow more intense
 19 therapy. Thus, at the beginning of treatment, the two groups of clients—
 20 those getting group therapy and those getting individual treatment—were
 21 *already different* in that the individual therapy clients were more seriously
 22 impaired. Thus, when the posttreatment evaluations were made, what was
 23 revealed were these *preexisting* difference in the two groups of clients, not
 24 the results of group therapy being more effective at helping people.

25 How can this threat be controlled for? One common way is to con-
 26 duct *pretreatment* assessments of clients, so that their scores on relevant
 27 measures of functioning, strengths, or psychopathology can be ascer-
 28 tained to determine if they are equivalent or not. Showing that they are
 29 similar goes some way toward reducing the possibility that selection bias
 30 resulted in nonequivalent groups from the beginning. And, apart from
 31 formal assessment of functioning, it is also very useful to compare the two
 32 (or more) groups of clients on a variety of important demographic vari-
 33 ables, with age, race, socioeconomic status, and gender being among the
 34 most important in this regard. Showing that your groups are equivalent

1 demographically helps alleviate the possibility that any posttreatment dif-
2 ferences may be due to some inherent differences in the groups due to
3 their demographic composition. Men and women, whites and blacks,
4 young and old, rich and poor, may all react differently to receiving certain
5 social work services. If one group is composed predominantly of men and
6 the other of women, posttreatment differences may reflect how the two
7 genders react to the treatment differentially (maybe men do not get as
8 much benefit from a given intervention). By being able to show that the
9 groups are pretty much the same on important demographic variables,
10 this threat can be partially controlled for. But only partially, because there
11 may be some unrecognized feature that distinguishes the two (or more)
12 groups and that is not measured but that nevertheless exerts an impact on
13 response to treatment. Such latent differences may never be able to be
14 detected, yet they may influence outcomes. This is why the methodologi-
15 cal refinement of creating treatment and control groups on the basis of
16 random assignment is considered by research methodologists to be one
17 of the best safeguards against groups not being equivalent on virtually all
18 relevant factors. More on this later.

19 **Diffusion/Contamination of Treatments**

20 This threat refers to a breakdown in the essential features supposedly dis-
21 tinguishing one treatment condition from others in a quasi-experiment.
22 Suppose you have two groups of clients, about half of whom are assigned
23 to receive group therapy alone and the other half individual counseling
24 during the normal process of agency operations. A social work researcher
25 wishes to see if the two groups of clients had differing outcomes. If they
26 did, then one approach to treatment might be seen as more effective than
27 another. Diffusion or contamination of treatments could occur if it
28 transpired that a number of persons supposedly receiving only group
29 therapy were later found to have sought out and obtained individual
30 counseling on their own, perhaps through a local church, independent of
31 the services your agency provided. This would compromise the research-
32 er's ability to compare the outcomes of group versus individual counsel-
33 ing since, in reality, you are comparing group therapy *plus* individual
34 counseling versus individual counseling alone, not individual versus
35 group therapy alone. Sometimes clients assigned to one treatment condi-
36 tion encounter and share information and experiences they have had

1 with clients assigned to another treatment or to a control condition. This
2 may inadvertently contaminate the “purity” of your supposedly differing
3 treatments.

4 **Therapist Bias/Allegiance Effects**

5 Individuals who evaluate social work services are often heavily personally
6 invested in the services they provide. If you have centered your profes-
7 sional life around a given form of service or therapy, be it psychodynamic
8 treatment, behavior therapy, solution-focused brief treatment, narrative
9 therapy, group work, family treatment, cognitive therapy, multisystemic
10 treatment, or other, it is understandable that your efforts at evaluating
11 these services could be biased by your preexisting investment in your
12 favored approach. This is a simple fact of life and not intended as a per-
13 sonal criticism of therapy researchers. This is why, when the founder
14 of eye-movement desensitization and reprocessing (EMDR) Francine
15 Shapiro published astonishingly positive outcome studies on EMDR;
16 when the founder of cognitive therapy, Aaron Beck, did the same with
17 cognitive therapy; or when the promoter of facilitated communication,
18 Douglas Biklen, claimed that seriously developmentally disabled people
19 could type with a high level of fluency, the scientific community wanted
20 to see if these initial reports of success undertaken by a new treatments’
21 advocates could be successfully replicated by independent researchers
22 lacking the personal and perhaps financial investments in these novel
23 approaches. This is not to imply that a therapy’s enthusiasts are some-
24 how lacking in honesty, but merely to recognize that all of us have our
25 biases and preferences and that, when potential conflicts are obvious,
26 more stringent standards of evidence may be called for. Sometimes, ini-
27 tially promising treatment results cannot be replicated by independent
28 researchers, giving rise to concerns that the initial results were a fluke or
29 perhaps contaminated by misguided zeal. At other times, independent
30 researchers repeatedly corroborate the effectiveness of new therapies,
31 which is the best possible outcome.

32 A particularly egregious threat pertaining to allegiance effects is the
33 rising practice known as *ghost authorship*. Ghost authorship occurs when
34 scientific writers who are employees of a corporation (e.g., a pharmaceu-
35 tical company, a tobacco firm, etc.) actually write up a complete journal
36 article, and the completed article is then offered to various respected

1 authorities in that particular field for their editing or reviewing. In return,
2 their name is appended to the article as an author or co-author; some-
3 times, the industry ghost writers who actually wrote the study do not
4 appear anywhere on the list of authors, or even in footnotes. This gives
5 rise to an independent-appearing paper apparently originating from
6 respected authors at a prestigious academic institution. However, this is
7 basically fraudulent in terms of authorship, and presents possibly ques-
8 tionable scientific findings, since the results were filtered through an
9 industry with a significant financial investment in the product being
10 evaluated. For example, one influential study that found that tobacco
11 advertising bans had little effect in reducing tobacco consumption was
12 “authored” by a respected marketing professor. In reality, the study was
13 ghost-written by tobacco associations (see Davis, 2008), who paid the
14 professor to present the study at conferences and before the U.S. Congress.
15 The ghost-authorship and payments were not disclosed as conflicts of
16 interest at the time. Similar conflicts have emerged in the reporting of
17 drug trials (see DeAngelis & Fontana, 2008). These are serious threats to
18 the integrity of research reporting and have led to more stringent edito-
19 rial policies relating to the mandatory disclosing of possible conflicts
20 of interest. Many journals, including *Research on Social Work Practice*,
21 now require the authors of accepted papers to describe all such potential
22 conflicts of interests, including factors such as receipt of grant funding,
23 consulting fees, or providing paid trainings in the treatments under
24 investigation.

25 **Lack of Treatment Fidelity**

26 *Treatment fidelity* refers to the extent to which services are delivered to
27 clients in the manner in which they were intended. Important aspects to
28 fidelity include the adherence of the therapists to the proper services
29 model, as well as the competence of the service providers. Unfortunately,
30 good adherence to prescribed therapies delivered by an incompetent
31 therapist fails to provide a fair test of a given service. Similarly, a highly
32 competent therapist who inadvertently blurs treatment techniques and
33 fails to adhere to assigned treatment protocols also compromises a study’s
34 treatment fidelity. Treatment fidelity can be enhanced through various
35 prospective methods. These include using well-proceduralized treatment
36 manuals, if these are appropriate to the clinical situation and available;

1 developing these if they are not; utilizing as therapists only those indi-
 2 viduals who are appropriately credentialed (e.g., licensed or certified,
 3 if appropriate) and skilled in the particular treatment method they are
 4 being asked to deliver (many persons possess generic credentials as a psy-
 5 chotherapist, but that does not mean they are necessarily competent in
 6 providing very specific forms of treatment); building into the delivery of
 7 services careful, regular supervision by a supervisor experienced in the
 8 treatment being provided; and providing some practice sessions with live
 9 supervision, prior to unleashing “treatment” therapists on real clients
 10 enrolled in the study.

11 Some other methods include audio- or videotaping treatment ses-
 12 sions and having these immediately reviewed by supervisors to ensure
 13 that therapists are delivering therapy as planned and taking corrective
 14 supervisory actions if it is not, and using any of a number of existing
 15 measures of therapist adherence or treatment fidelity (see Nezu & Nezu,
 16 2008, pp. 263–281 for a review of some of these).

17 **Concluding Remarks on Threats to Internal Validity**

18 This review of some common threats to internal validity need not be
 19 overwhelming. If you would like a more light-hearted approach to trying
 20 to understand how these factors can complicate your ability to conclude
 21 that a given treatment is effective, go to the *YouTube* website and look up
 22 some of the videos of the comedic magician team of Penn and Teller
 23 (e.g., <http://www.youtube.com/watch?v=MzjoKhBkLYg>) who demon-
 24 strate the power of placebo therapies and alternative medicine treatments
 25 by setting up bogus health clinics in which fake doctors offer to treat
 26 people with magnets, kazoo music, snail mucus (!), and toilet plungers.
 27 These videos are really hilarious.

28 Keep in mind that one needs to control for *plausible* controlling fac-
 29 tors, not every conceivable potential confound. For example, if you use a
 30 no-treatment control group or a treatment-as-usual comparison group,
 31 it is a useful practice to demonstrate that your two groups are statistically
 32 equivalent at pretreatment on important demographic and clinical fac-
 33 tors, on variables such as age, race, gender, education, socioeconomic
 34 background, severity of clinical symptomatology, etc. But you need not
 35 examine essentially irrelevant factors such as clients’ astrological signs,
 36 phases of the moon when treatment is administered, and the like.

1 The designs described in the next chapter illustrate how social work
2 researchers attempt to introduce various controls for some of these legit-
3 imate threats by using quasi-experimental methods. All the threats to
4 internal validity mentioned in this chapter as potentially impacting the
5 interpretation of results from studies using pre-experimental research
6 designs also apply to the quasi-experimental designs presented in the
7 next chapter.

8 SUMMARY

9 This chapter has provided an overview of the logical and design features of
10 those approaches to evaluation that have been traditionally labeled as pre-
11 experimental, in that they involve the assessment of only a single group of
12 clients. These designs may assess clients only after they have received an
13 intervention, or they may assess client functioning before and after expo-
14 sure to a treatment. A variety of ways were presented to potentially
15 strengthen these pre-experimental designs, and numerous examples were
16 described, taken from the published social work journal literature, describ-
17 ing how each of these designs has been used to evaluate the outcomes of
18 social work services. Also provided was a review of various threats to
19 internal validity, factors that can complicate one's ability to conclude that
20 a given treatment *caused* any apparently positive outcomes.



3

Quasi-Experimental Group Designs



5 **I**n Chapter 1, a series of five questions was presented, questions which
 6 required an increasingly more stringent level of evidence in order for
 7 them to be credibly answered. Question 1, which asks, *What is the status of*
 8 *clients after they have received a given course of treatment?* can be answered
 9 using the $X - O_1$ design, and Question 2, *Do clients improve after receiving*
 10 *a given course of treatment?* can be addressed with the $O_1 - X - O_2$ design,
 11 with O referring to a point in time when clients are systematically assessed
 12 using some reliable and valid outcome measure and X depicting the
 13 individuals' receipt of some form of social work intervention. However,
 14 Question 3, *What is the status of clients who have received a given treatment*
 15 *compared to those who did not receive that treatment?* cannot usually be
 16 answered with a pre-experimental design because the question involves
 17 not comparing the same group posttreatment with its pretreatment level
 18 of functioning, but instead requires looking at the results for a group of
 19 clients versus a similar group who did not receive the treatment. Thus,
 20 some sort of comparison group is required. As noted before, it is this ele-
 21 ment of having some sort of comparison or control group that has tradi-
 22 tionally distinguished the quasi-experiments from the pre-experiments.
 23 The terms *control group* and *comparison group* have somewhat different
 24 meanings. A control group is one that does not receive any kind of formal
 treatment at all. Generally, the sole contact that members of a control group

1 may have as a result of their participation in this type of study is via the pro-
 2 cess of being assessed. A comparison group refers to individuals who receive
 3 some sort of alternative treatment. This may be treatment as usual (TAU),
 4 which is some other legitimate intervention used when trying to see if one
 5 treatment is superior to another or to a placebo treatment. A *placebo control*
 6 *group* occurs when clients receive a seemingly credible treatment, which in
 7 reality the researchers presume to be ineffective in terms of helping improve
 8 the problem or situation being presented by the clients. Sometimes, placebos
 9 are otherwise legitimate practices, like relaxation therapy or hypnosis, which
 10 may have their place in the effective treatment of some problems, but not in
 11 others. Relaxation training may have a role in helping people with general-
 12 ized anxiety disorder, but when used to treat persons with schizophrenia,
 13 a condition for which it has no discernible value, it would be considered
 14 a placebo therapy. At other times, a placebo may be a frankly bogus or fake
 15 treatment, a condition that is either deceptive or simply known to be useless.
 16 For example, audiotapes or CD recordings containing supposedly sublimi-
 17 nal messages have been marketed to help people lose weight or stop smok-
 18 ing, but extensive research has shown that they do not work; other placebos
 19 may include sham needle placement in acupuncture studies, eye movements
 20 conducted in ways that the treatment theory suggests should not have any
 21 effect, administering homeopathic pills that actually contain no trace of an
 22 active ingredient, and the like. If presented in a believable manner, all these
 23 can make a useful placebo treatment for comparison purposes. The term
 24 *experimental group* usually refers to those clients who received the actual,
 25 legitimate, or novel treatment under formal investigation.

26 The next sections describe some common variations of more sophis-
 27 ticated control and comparison group quasi-experimental designs.

28 **DESIGNS WITH CONTROL GROUPS AND POSTTESTS ONLY**

29 **The Posttest-only No-treatment Control Group Design**

30 Moving incrementally, the simplest of the quasi-experimental designs
 31 can be diagrammed as follows:

Group 1 Received Treatment	$X - O_1$
Group 2 Did Not Receive Treatment	O_1

1 Some researchers might pose a formal predictive hypothesis, such as
2 *Clients who received social work treatment X will have higher functioning*
3 *than will clients who did not receive treatment X*, and use this design to see
4 if this hypothesis is supported or not. The inferential logic is simple:
5 If the group who received X differs significantly from those who did not
6 receive X, this would support the hypothesis that X produces certain
7 effects, above and beyond receiving nothing.

8 Keep in mind that, for the purposes of our discussion, X can literally
9 be *anything* within the scope of social work practice—individual therapy,
10 group therapy, marital or couples counseling, a community-wide inter-
11 vention, a local or state law, or even a national welfare policy. To give you
12 an idea of how this design is used at a macro level, one might compare
13 some outcome measure across two similar states, one with a certain law
14 in effect and the other lacking that law. In this case, the differing laws
15 represent our X or intervention or independent variable. Florida, for
16 example, does not require motorcycle riders to wear safety helmets.
17 Other states do require this. One could compare fatality rates for motor-
18 cycle accident victims across states, comparing those with and those
19 without a protective helmet law, to try to investigate the contention that
20 wearing helmets saves lives. Closer to the field of social work, some states
21 require that abortion providers ensure that the parents of minors be
22 notified prior to performing an abortion on that minor. Some stakehold-
23 ers have questioned if these parental notification laws deter minors from
24 seeking abortions, and they have used this posttreatment no-treatment
25 control group design to look at abortion rates across states with and
26 without a parental notification law.

27 Ideally, the period of time during which O occurs is roughly the same
28 for the two groups. For example, if the treatment group is provided a
29 social work service in July, with intervention lasting 3 months, then the
30 group would be reassessed immediately posttreatment, near the end of
31 October (when all clients had completed treatment). And, also around
32 the end of October, the group of people who had been identified as not
33 having received treatment (those in the no-treatment control group)
34 would be similarly assessed. Conducting these O assessments at about the
35 same time helps control for events in the external world that can influ-
36 ence client functioning, events unrelated to response to treatment or the
37 waxing and waning of symptomatology or problem severity. Imagine if
38 you conducted a study like this and your treatment group was assessed

1 around September 1, 2001, and the nontreatment group was assessed in
2 mid-September, after the terrorist attack on September 11, 2001? Clearly,
3 your study might have been seriously compromised. Similarly, if your
4 study was in Louisiana and your O evaluations for the two groups
5 occurred on either side of Hurricane Katrina or the Gulf oil spill, it would
6 be very difficult to make any legitimate comparisons in such instances.
7 Here are some examples of using this type of research design.

8 *Evaluating Virginity Pledges*

9 This design, which is sometimes called the *static-group comparison design*,
10 was used by Rosenbaum (2008) in her analysis of the effectiveness of
11 *virginity pledges* in terms of deterring the initiation of sex and of their
12 influence on the use of birth control. This is certainly an important issue,
13 since the U.S. government spends (as of the time of this writing) over
14 \$200 million annually on abstinence promotion programs, some of
15 which involve virginity pledges. Abstinence-only sex education in gen-
16 eral, and virginity pledges in particular, are of uncertain effectiveness in
17 preventing the initiation of sexual activity, pregnancy, or sexually trans-
18 mitted diseases. In this study, derived from a large-scale national survey,
19 289 adolescents aged 15 years or older at the time they voluntarily com-
20 pleted virginity pledges were compared 5 years after taking their pledge
21 with 645 nonpledgers, teenagers generally equivalent demographically to
22 those who took the pledge to abstain from sex until marriage. Note that,
23 ethically and practically, adolescents could not be randomly assigned to
24 undertake a personal pledge to abstain from sexual intercourse, so this
25 quasi-experimental design was an excellent method to initially evaluate
26 such interventions. The results?

27 Five years after the pledge, 82% of pledgers denied ever having pledged.
28 Pledgers and matched nonpledgers did not differ on premarital sex, sex-
29 ually transmitted diseases, and anal and oral sex variables. Pledgers had
30 0.1 fewer past-year partners but did not differ in lifetime sexual partners
31 and age of first sex. Fewer pledgers than matched nonpledgers used birth
32 control and condoms in the past year and birth control at last sex.
33 (Rosenbaum, 2008, p. e110)

34 These results will be disappointing to the advocates of abstinence-
35 only sex education and virginity pledges, and again point to the possible

1 dangers of public policy (funding certain types of programs) getting
2 ahead of the evidentiary curve. But the picture is admittedly mixed.
3 For example, Martino et al. (2008), also using this design, found a pro-
4 tective effect for virginity pledges, comparing 12- to 17-year-olds who
5 took a virginity pledge with those who had not, some 3 years later. About
6 42% of the nonpledgers had initiated intercourse at 3-year follow-up,
7 compared to only 34% of the pledgers. The effect was modest but real.

8 Please note that the above examples, and others presented in this
9 book, are being used to accurately illustrate the use of selected quasi-
10 experimental researcher designs. They should not be interpreted to reflect
11 comprehensive conclusions based on systematic reviews of all the relevant
12 research evidence dealing with the possible effectiveness or ineffectiveness
13 of the interventions being discussed. Although the summaries of the stud-
14 ies are accurately presented, any conclusions drawn from an individual
15 study should *not* be presumed to reflect the value of the presented treat-
16 ments as assessed by a comprehensive review of all relevant research.

17 *Evaluating Foundation Master's of Social Work Skills Training*

18 Social worker Dorothy Carrillo's Ph.D. dissertation used this design in a
19 way that took advantage of a naturally occurring situation. Dorothy was
20 assigned to instruct a class of second-year master's degree in social work
21 (MSW) students devoted to the topic of teaching direct practice skills.
22 In the natural course of events, Dorothy's class consisted of two types of
23 students, 15 had earned the bachelor of social work (BSW) degree and
24 were enrolled in the advanced standing program, wherein they could
25 exempt the first-year foundation course in direct practice. An additional
26 23 students were in the traditional 2-year MSW program and had com-
27 pleted a first-year foundation class in direct practice skills. Thus, without
28 any manipulation on Ms. Carrillo's part, she had two different types of
29 students taking her class; some had completed the first-year direct prac-
30 tice class as a part of their MSW program, and the others had not had this
31 course as a part of their MSW curriculum. As holders of the BSW, they
32 were exempted from taking this course earlier.

33 The Council on Social Work Education permits MSW programs to
34 exempt BSW students from taking certain foundation MSW courses,
35 such as in direct practice, using as its rationale that selected BSW classes
36 are redundant with first-year graduate MSW training. This was a hereto-
37 fore untested assumption, and Ms. Carrillo was in a position to test it.

1 At the beginning of the semester, the standard practice in this class
2 was for students to interview a simulated client who role-played being an
3 elderly person. The students were given some background information
4 and asked to obtain information about this pseudo-client's life history
5 and present circumstances. These interviews were videotaped. In effect,
6 Ms. Carrillo had a ready-made posttest-only no-treatment control group
7 design dropped in her lap. She arranged for these videotapes, which were
8 made very early on in the semester, to be reliably coded in terms of
9 the student's use of selected core interviewing skills related to concepts
10 such as facilitation, questioning/clarification, and support/empathy. One
11 coder (a social work doctoral student) rated all 38 tapes and a reliability
12 coder (another social work doctoral student) independently rated 13
13 (34%) of them as well. Inter-rater agreement on the use of the three
14 selected interviewing skills ranged from 86% to 92%, an acceptably high
15 level of agreement, indicating that the ratings were really of what the
16 videotaped interviewer was doing, as opposed to the rater's purely sub-
17 jective impressions. The study's purpose was to evaluate the following
18 null hypothesis: Advanced standing and 2-year MSW program students
19 will display equivalent use of selected foundation interviewing skills
20 when assessed at a comparable point in their curriculum.

21 After the semester was over and the tapes were blindly rated, the code
22 was broken and Ms. Carrillo could see whether the two groups of students
23 indeed displayed similar interviewing skills or if one group was better
24 than the other. The results (fortunately for social work educational policy)
25 were consistent with the CSWE's policy of granting advanced standing to
26 BSW students, exempting them from foundation MSW courses, since the
27 BSW students and second-year MSW students *did* display similar levels
28 on the three selected skills (see Carrillo & Thyer, 1994).

29 *Increasing Access to Dental Care for Medicaid Preschool Children*

30 It is important that children receive regular dental care. Many children,
31 especially children from poor families, fail to receive such care. The state
32 of Washington developed an Access to Baby and Child Dentistry (ABCD)
33 program offering extended benefits to participating Medicaid-enrolled
34 children and higher fees to dentists seeing such children. Participation in
35 the ABCD program was voluntary, and over the course of time, some
36 families receiving Medicaid enrolled their children in this program and
37 some did not. After the ABCD program had been in place for a year,

1 researchers contacted 282 parents of children aged 13–36 months, with
 2 about half the families being enrolled in the ABCD program and about
 3 half not being enrolled. Fully 43% of the children whose parents had
 4 enrolled them in the ABCS program had visited a dentist within 1 year of
 5 their being signed up, whereas only about 12% of the Medicaid children
 6 who had not been enrolled in the ABCD program had visited a dentist.
 7 In other words, enrollment in the program seemed to result in a child
 8 being 5.3 times as likely to have had at least one dental visit compared to
 9 a child not in the program. Also, the parents of ABCD-enrolled children
 10 reported that their kids were less fearful of the dentist. The program
 11 seemed very successful in promoting early dental care for children from
 12 poor families (Grembowski & Milgrom, 2000).

13 *Promoting Positive Attitudes Toward Computer Use Among Master's*
 14 *of Social Work Students*

15 The use of technology, such as computers and the Internet, has dramati-
 16 cally transformed social work education and practice. Like any new
 17 innovation, incorporating technology into our field met with critics, and
 18 it only moved ahead via fits and starts. Over a decade ago, some social
 19 work programs developed classes aimed at enhancing MSW students'
 20 technical skills and promoting positive attitudes related to computer use
 21 in practice. One such program was developed at the School of Social
 22 Work at Bar Ilan University in Israel (Monnickendam & Elliot, 1997).
 23 Faculty there offered direct practice students a new course devoted to
 24 computer literacy. Some direct students took the computer literacy class
 25 and some did not. Students in the administrative track did not take it
 26 either. At the end of the school term, all MSW students completed mea-
 27 sures related to attitudes toward using computers in the human services.
 28 The design can be diagrammed as follows:

N = 34 Direct Practice Students	X – O ₁
N = 30 Direct Practice Students	O ₁
N = 32 Administrative Track Students	O ₁

29 The results were an interesting mix. Direct practice students who took
 30 the class (depicted as X in the design) had more positive attitudes about

1 computer use than did direct practice students who did not take the com-
2 puter literacy class. This supports the hypothesis that taking such a course
3 can indeed promote more positive attitudes about using computers.
4 However, the administrative track students who *did not* take the com-
5 puter literacy class had attitudes that were even more positive about com-
6 puters than did the direct practice students who did take the class. This
7 complicates things a bit. Plausibly, it could be contended that adminis-
8 trative track students had preexisting positive attitudes about technology,
9 that the direct practice students were more touchy-feely in their approach
10 to social work, and some presorting into direct practice or administrative
11 tracks occurred partly on this basis before they even took the computer
12 course. This is another nice example of taking advantage of a naturally
13 occurring situation to imbed a quasi-experimental research design into
14 one's evaluation work. After learning more about the designs described
15 in this book, you will be surprised to find how many professional situa-
16 tions arise that lend themselves to being evaluated using them.

17 *Moral Development of Sex Offenders*

18 Social worker Frederick Buttell is one of the country's foremost research-
19 ers in the area of treating sex offenders who are men. Dr. Buttell adopted
20 this design to evaluate the levels of moral reasoning evidenced by
21 convicted sex offenders ordered into treatment. Basically, 72 male sex
22 offenders completed a previously published, supposedly reliable and
23 valid measure of moral reasoning, an instrument called the Defining
24 Issues Test (DIT), a measure that presents one with a series of written
25 moral dilemmas and asks the reader to evaluate a list of questions he or
26 she might consider when making a decision about what to do in the
27 depicted situation. The DIT presents six moral dilemma scenarios and
28 results in an overall score said to be related to "principled morality, with
29 higher scores reflective of higher levels of moral reasoning." Buttell's
30 (2002) study was *not* a treatment outcome study; rather, it was some-
31 thing called a *cross-sectional survey*. In this instance, X reflected the back-
32 ground of being a sex offender, and O the scores on the DIT. The group
33 lacking X was comprised of normative data on the DIT reported by the
34 developer of this test. Dr. Buttell was thus able to compare the levels of
35 moral reasoning exhibited by convicted sex offenders with the levels of
36 moral reasoning evidenced by males nationwide. It was found that the
37 sex offenders (mean age of 38 years) had significantly lower levels of

1 moral reasoning than did graduate students, college students, adults in
2 general, and high school students, but scored similarly to junior high
3 school students (as assessed against normative data for the DIT, not
4 actual groups of clients). Buttell speculates that this low level of moral
5 reasoning, focusing more on self-interest than on empathy, may be
6 linked to being a sex offender and to the high rates of repeat offenses and
7 recidivism displayed by this group. This example also illustrates how
8 quasi-experimental designs can be used in other forms of research besides
9 outcome studies, with cross-sectional surveys like this one being a good
10 example.

11 *Strengths of the Posttest-only No-treatment Control Group Design*

12 This design is capable of answering the question “Do clients who received
13 treatment X fare any differently compared with clients who did not
14 receive X?” but this ability is dependent on several assumptions. First,
15 the intervention is something that can be seen or at least determined to
16 be present or absent. Second, the outcome measure(s) must be both reli-
17 able and valid. Third, a sufficient number of participants exists in each
18 group (you can’t reliably use this design with just a few folks in each
19 condition). And, last, the two groups, treatment and no-treatment con-
20 trol group, are truly similar in all important respects, *except* that the
21 intervention group has received the treatment while the no-treatment
22 group has not. This last point can be a stickler.

23 *Limitations of the Posttest-only No-treatment Control Group Design*

24 You simply cannot assume that the two groups are similar simply because
25 they share similar problems. You *can* conduct a statistical analysis on the
26 members of the two groups, attempting to see if they do differ on any
27 variables. These variables can be demographic characteristics (e.g., age,
28 race, gender, etc.), or they can be clinical features (e.g., severity of symp-
29 toms). You do not need to analyze everything, only factors that can be
30 plausibly deemed as important. Client’s astrological signs, for example,
31 are likely unimportant. In the Carrillo and Thyer (1994) study, Ms. Carrillo
32 was able to show that, in terms of gender, marital status, race/ethnicity,
33 and income, the two groups were indeed very similar. This is good, in that
34 it argues that any differences observed during the observation period
35 are more likely due to the treated group’s exposure to treatment than to
36 any preexisting conditions in place before the treated group received inter-
37 vention. However, when Carrillo looked at age, she found a difference.

1 The advanced standing students' average age was 23.2 years, whereas the
2 2-year program students' was 29.2 years, a 6-year difference that was sta-
3 tistically significant. If she *had* found a difference posttraining, favoring
4 the 2-year students, this might have been a problem since the possibility
5 existed that the older students were somehow wiser, more seasoned in life,
6 more mature, or perhaps had more prior social work job experience than
7 the advanced standing students, and it was *these preexisting* factors that
8 accounted for their superior interviewing skills performance, not the first-
9 year foundation practice skills class that they had taken and the advanced
10 standing students had not. However, since Carillo did not find a difference
11 in interviewing skill ability, she was not faced with handling this possible
12 confound. But you can see how it might arise and present problems in
13 making a claim that intervention was responsible for such differences.

14 Another possibility to consider is that the null result of Carrillo's
15 study, the finding of no difference between the MSW and BSW students'
16 performance on interviewing skills, is that the outcome measure was too
17 blunt or insensitive an indicator to pick up on the actual differences
18 between the two groups. This is a common confound for any study that
19 fails to find anticipated differences. One interpretation is that the two
20 groups truly did not differ after training, but another is that the depen-
21 dent variable, the outcome measure, is too crude. One can guard against
22 this by choosing to use outcome measures of known reliability, validity,
23 and sensitivity to differences and changes.

24 You can go to the websites of many advocacy groups that find some
25 difference between two groups and claim that this difference exists because
26 of some factor that separates the two groups. These claims are then used
27 to buttress some policy position of the advocacy group. For example, you
28 may have heard a public service announcement on the radio advocating
29 for parents to eat more meals sitting down with their children because it
30 has been shown that nondelinquent, non-drug-using, non-sexually active
31 kids report having more sit-down meals with their parents than do kids
32 who engage in delinquency, drug use, and sex. Assume that this is true
33 (I am not claiming it is). The inferential logic is something like this:

- 34 1. A large group of kids are surveyed in terms of delinquency, sexual
35 activity, drug use, etc.
- 36 2. Some kids engage in a high amount of such problem behavior
37 and some do not.

1 3. When you examine the backgrounds of these two groups of
2 youth, those with high- and low-delinquency status, you find that
3 they differed in the amount of time they claimed to spend eating
4 family-style meals.

5 So far, we are on solid ground. But the next steps may be logically unten-
6 able, claiming, for example:

- 7 4. Eating meals family-style will protect kids from becoming
8 delinquent.
9 5. It should be a matter of policy to encourage parents to eat more
10 family-style meals with their children, and this will reduce the
11 risk of kids abusing drugs, committing crimes, or becoming
12 sexually active.

13 Can you detect the gap in logic between step 3 and steps 4 and 5? There
14 may be significant other preexisting differences or factors that lead to
15 fewer family-style meals *and* a higher risk of delinquency. For example,
16 one potential confounding factor might be having both biological par-
17 ents in the home versus having only one (e.g., a single-parent home). The
18 *real* cause might be this difference in overall parent availability for super-
19 vision, *not* the specific act of eating meals together.

20 While preparing this chapter, I went to the website of The Heritage
21 Foundation (www.heritage.org), a conservative think tank that regularly
22 produces policy papers based on contemporary social science research.
23 They highlighted a number of studies on their site, studies which, by
24 implication if not by explicit statement, supported generally conserva-
25 tive positions on various social welfare policy matters. Here are some
26 examples:

- 27 • “Compared with peers from intact families, adolescent and
28 young-adult women who experienced parental divorce were
29 significantly more likely to give birth out of wedlock” (based
30 on Martin, 2005).
31 • “Teens who were exposed to high levels of sexual content on
32 television were twice as likely to become pregnant during a 3-year
33 period than peers who had lower levels of exposure to sexual
34 content” (based on Chandra et al., 2008).

- 1 • “Compared to adolescents who were virgins, those who
 2 had initiated sexual activity were 58% more likely to
 3 engage in delinquent behavior in the year after they had
 4 become sexually active” (based on Armour & Haynie,
 5 2007).
 6 • “Compared with peers in intact families, children in
 7 blended or step families tended to have significantly lower
 8 GPAs and less positive engagement with school tasks
 9 and relationships” (based on Halpern-Meekin & Tach,
 10 2008).

11 Now, these individual studies, all using the posttest-only no-treatment
 12 control group design in the context of a survey study, may well be sound
 13 pieces of research. But logical and scientific problems arise if these
 14 summarized findings (all four quotes above were taken from the Heritage
 15 Foundation webpage on August 17, 2009) were used to justify social pol-
 16 icies advocating something like the following:

- 17 • Since youth from families whose parents were divorced
 18 were more likely to give birth out of wedlock, we can drive
 19 down the illegitimacy rate by enacting policies that
 20 discourage divorce.
 21 • Since exposure to sexual content on television promotes
 22 teenage pregnancy, let’s exact more stringent laws censoring
 23 such sexual content in order to reduce teenage
 24 pregnancy.
 25 • Since sexual activity in youth is associated with a much higher
 26 risk of subsequently engaging in crime, let’s promote the teenage
 27 norm of sexual abstinence as a means of reducing juvenile
 28 delinquency.
 29 • Since kids from divorced families have lower academic
 30 performance and social skills, let’s try and improve the well-being
 31 of youth by making it more difficult to get a divorce when
 32 children are involved.

33 Although the Heritage Foundation itself is too sophisticated to draw
 34 such explicit conclusions, causal linkages, and policy recommendations
 35 like those four just noted above, other groups are not so constrained.

1 The problem is again that of uncontrolled factors that are actually caus-
2 ally involved, rather than the particular one cited in a given study.
3 For example, although it may well be true that greater exposure to televi-
4 sion sexual content is associated with subsequent teenage pregnancy, the
5 true causal factor (this is only a hypothetical example) may be lack of
6 parental supervision. In the young, this may lead to unsupervised televi-
7 sion watching, with the kids naturally gravitating to the more salacious
8 shows (my own kids—aged 17, 15, 14, and 12—love *South Park*, which I
9 consider inappropriate for them but my wife thinks is fine). Then, with
10 puberty and even greater freedom, this lack of parental supervision pro-
11 vides hormonally charged adolescents with more opportunities not just
12 to *watch* sexy TV shows, but to do what comes naturally! Voila—more
13 babies born out of wedlock. But, the causal factor is not watching TV
14 shows; it's inadequate supervision by parents.

15 A more benign example of this type of fallacious reasoning can be
16 found in the high school band program two of my sons participate in.
17 The music programs distributed at their performances contain snippets
18 from studies that have shown that high school students who participate
19 in music education programs like bands are more likely to graduate from
20 high school, to have higher GPAs, and to enroll in college more often.
21 The implicit message is that high school music programs must be better
22 funded in order to promote the academic attainments of teenagers.
23 No mention is given to the possibility that the kids who gravitate to music
24 and band may be brighter, harder-working, or come from higher socio-
25 economic levels initially, and that these preexisting differences account
26 for disparate outcomes between band kids and non-band high-schoolers.
27 Once you become aware of such potential confounds and the errors in
28 reasoning they may lead to, you will be amazed at how frequently you
29 will encounter examples of precisely these types of errors in the media
30 and in everyday life. Would you be surprised to read a study that found
31 that BSW students hold more liberal and progressive values when they
32 graduate from college than do other graduating majors? What do you
33 think is more likely responsible for such differences—that BSW educa-
34 tion is inherently liberalizing, or that the more liberal and progressive
35 students gravitate to the social work major?

36 This inability to control for the equalization of groups on all mean-
37 ingful factors really complicates causal inference in designs such as this
38 and the others described later in this chapter.

1 The Posttest-only Comparison Group Design

2 Early on, Macdonald (1952) provided some suggestions relating to the
3 evaluation of social work services that are relevant here, noting:

4 I think we are ill-advised to cast about looking for no-treatment groups
5 as control groups. We cannot keep in touch with people who are not
6 being treated in order to learn about them and their problems. . . .
7 I think we are better advised to examine results within the group treated,
8 comparing subgroups who have different diagnoses or who have the
9 same diagnoses and are treated differently. We can evaluate new meth-
10 ods in comparison with the old. . . ." (Macdonald, 1952, p. 137)

11 Although I believe that Macdonald was overly pessimistic regarding the
12 feasibility of employing no-treatment control groups (for example, using
13 agency wait-lists of clients is one ethical option), she was spot-on in
14 terms of using quasi-experimental designs to evaluate clients who
15 received differing forms of intervention, or in comparing a novel therapy
16 with TAU.

17 The inferential logic of posttest-only comparison group design is
18 relatively simple. Compare the outcomes of a group of people who
19 received one type of intervention (X), program, or training, against the
20 outcomes of another group who received an alternative intervention (Y).
21 If the group X outcomes are better than the group Y outcomes, this cor-
22 roborates the hypothesis that treatment X is superior on some dimension
23 than treatment Y. If the outcomes between X and Y are no different, then
24 this weakens the hypothesis that X is a superior treatment. This design
25 can be diagrammed as follows:

26 $X - O_1$

27 $Y - O_1$

28 The design also permits more than one comparison group. For example,
29 an intervention program for male batterers might consist of two ele-
30 ments, individual counseling and group counseling. In the normal course
31 of service delivery, some referrals might get individual counseling alone,

1 others will get group counseling alone, and a third set of clients will
 2 receive both interventions. This could be diagrammed as follows:

3 $X_{\text{individual counseling alone}} - O_1$

4 $Y_{\text{group counseling alone}} - O_1$

5 $Z_{\text{combined counseling}} - O_1$

6 Over, say, a 12-month period, you would likely have unequal numbers
 7 being nonrandomly assigned to one of these three groups. That is accept-
 8 able, so long as no one group’s sample size falls below an acceptable limit.
 9 Clients would be enrolled throughout the year, and at some point be
 10 terminated at the conclusion of their course of treatment. The formal
 11 assessment could occur then; thus, these posttreatment assessments
 12 would occur at the same point in time procedurally, at the end of treat-
 13 ment for each client, but not at the same point of time chronologically, as
 14 these assessments, too, would occur throughout the year. Here are some
 15 actual published examples of social work researchers using this design.

16 *Are Regular Faculty Better Teachers Than Adjuncts or Doctoral Students?*

17 There is much hand wringing in contemporary academia because the pro-
 18 portion of teaching positions held by full-time, tenure-track faculty is
 19 declining in recent years, relative to instruction provided by community-
 20 based adjuncts hired to teach individual classes, or by doctoral students.
 21 Is this hand wringing justified? Are these lesser mortals somehow inherently
 22 less able teachers than full-time faculty? My colleagues and I decided to
 23 investigate this issue by taking advantage of the publically available course
 24 evaluations completed by students at my university. These evaluations
 25 included qualitative and quantitative sections, and we were able to retrieve
 26 these for several hundred BSW and MSW classes offered over several years.
 27 The courses were divided into those taught by regular full-time faculty, by
 28 adjuncts, and by Ph.D. students, after tossing out all those classes with fewer
 29 than ten respondents available. This design can be diagrammed as follows:

Classes taught by Regular Faculty (N = 181) $X - O_1$

Classes Taught by Adjuncts (N = 63) $Y - O_1$

Classes Taught by Ph.D. Students (N = 50) $Z - O_1$

1 We examined the quantitative scores on these evaluations and statisti-
2 cally compared the scores obtained by the three groups of instructors.
3 Basically, we found that the full-time faculty's course evaluations were
4 statistically the same as those earned by the adjuncts, and the evaluations
5 earned by the adjuncts were the same as those earned by the doctoral
6 students, but the full-time faculty had course evaluations that were sta-
7 tistically significantly better than those obtained by the Ph.D. students.
8 However, the size of this difference, although statistically reliable, was
9 quite small, and in fact practically meaningless. For all practical purposes,
10 the three groups of instructors earned similar course evaluations. This
11 null result suggests that, within this university and this social work pro-
12 gram, instructional quality is not suffering because of our use of adjuncts
13 and doctoral students. Parenthetically, this study is another example of
14 low-hanging fruit. The data were publically available on the university's
15 website. All we had to do was retrieve and analyze it. We did this, of
16 course, only after obtaining our university's institutional review board
17 approval for the study (see Thyer, Myers, & Nugent, 2011).

18 *Do Social Workers Make Better Child Welfare Workers?*

19 Social worker Robin Perry (2006) at Florida A & M University used this
20 posttest-only control-group design to test the hypothesis, so widely held
21 in the child welfare field, that having earned a professional social work
22 degree (BSW or MSW) is superior preparation for child welfare practice
23 than is receiving degrees in other disciplines. Working within Florida's
24 Department of Children and Families, Perry was able to obtain the semi-
25 annual evaluations of all child protective service (CPS) workers who
26 were employed within the state system as of March 2002—some 2,500
27 employees. He randomly selected about 25% of the employees and was
28 able to classify them as either having earned a BSW degree *or* another
29 degree (e.g., in a field such as psychology, criminology, sociology, busi-
30 ness, education, etc.). His outcome measure was the Career Service
31 Evaluation Form completed semi-annually by supervisors on each worker;
32 this form contains quantitative ratings across various important measures
33 of worker performance. Perry grouped all the evaluation form data for
34 the CPS workers with a BSW into one group (X), and compared their
35 data with those performance evaluation ratings received by workers with
36 the other educational degrees (Y). Basically, Perry found that BSWs did
37 not score higher than CPS workers with other professional backgrounds.

1 Thus, his hypothesis was disconfirmed, and the claim that social work
2 degrees make better preparation for responsible positions within child
3 welfare is weakened. This is a rather important finding. Although disap-
4 pointing for the profession of social work, it calls into question the ratio-
5 nale for the current practice of allocating large amounts of federal financial
6 support specifically designed to prepare BSWs and MSWs for careers in
7 child welfare. Perhaps such funding should be opened to students in all
8 disciplines, social work *and* non-social work, who are prepared to commit
9 to a career in public child welfare services? It would be premature to
10 advocate for opening up such funding now, particularly since Perry's pro-
11 vocative findings have not yet been tested, much less replicated, by others;
12 but, in science, the burden of proof rests on the person making an unusual
13 claim. In this instance, the unusual claim is that social work training is
14 superior to non-social work training in terms of preparing child welfare
15 workers. The evidence in favor of this hypothesis remains rather weak.

16 *Do School Social Worker Services Reduce Truancy?*

17 Newsome, Anderson-Butcher, Fink, Hall, and Huffer (2008) used a
18 posttest-only no-treatment control group design with 115 urban second-
19 ary school students in Ohio. About half received special school social
20 work services aimed at reducing absenteeism and about half did not.
21 The untreated half were matched as closely as possible to the treated stu-
22 dents, but they were not assigned to treatment versus nontreatment
23 using random assignment procedures, which makes this study a quasi-
24 experimental design, as opposed to a true randomized controlled trial.
25 The social worker provided an average of 14 direct or indirect interven-
26 tions for each referred student, including one-on-one counseling, group
27 counseling, phone contacts, and meetings with school personnel, parents,
28 or outside agency staff. Intervention lasted 9 weeks, and attendance data
29 were obtained from the school district's management information system
30 for each student, for the 9 weeks prior to the intervention and for the
31 9 weeks during which it was provided. This study was approved by the
32 human subject's IRB at both Ohio State University (to which some of
33 the authors were affiliated), as well as by the local school district's IRB.
34 Unfortunately, the students receiving school social work services did *not*
35 experience a statistically significant reduction in unexcused absences rela-
36 tive to the no-treatment group. Thus, this report could be considered one
37 of a treatment failure. Although personally, perhaps, and professionally

1 disappointing, finding out that certain forms of social work service are
2 *not* effective in producing desired outcomes is a good thing to know,
3 compared to not knowing.

4 *Does Early Childhood Intervention Impact Educational*
5 *Achievement and Arrest?*

6 The interdisciplinary research team of Reynolds, Temple, Robertson,
7 and Mann (2001), which included a social work author with a MSW,
8 conducted a 15-year follow-up of two groups of low-income children
9 attending public schools. A nonrandomized matched-group cohort
10 design was used, involving 1,539 children, most of whom were African
11 American. Nine hundred eighty-nine of the children received the Chicago
12 Child-Parent Program, which provided comprehensive education,
13 family, and health services, including half-day preschool at ages 3 and 4,
14 and full-day kindergarten at elementary schools for kids aged 6–9 years.
15 It was basically a very intensive family and educational support program.
16 The comparison group of children received a less intensive set of services,
17 basically kindergarten without the preschool and additional social ser-
18 vices. The outcome measures were assessed some 15 years later, a remark-
19 ably long follow-up period, when the children averaged 20 years of age.
20 Measures of educational attainment included high school/GED comple-
21 tion, grade retention (failures), and juvenile arrest records (including
22 numbers and type of arrests). The results? The intensive intervention
23 group had a significantly higher level of high school completion and
24 lower rates of dropout and grade retention, all relative to the comparison
25 group. The former also completed more years of schooling. The inten-
26 sive intervention group had significantly lower rates of arrest, lower rate
27 of multiple arrests, and fewer arrests for violent crimes. This was a
28 remarkably ambitious and large-scale study, with an extended follow-up
29 period, and it showed appreciable benefits for the intensive social ser-
30 vices intervention compared to more standard options. Perhaps it is not
31 surprising that this study was published in the *Journal of the American*
32 *Medical Association*, one of the most prestigious scientific journals in the
33 world, with the first author being a professor of social work! Can quasi-
34 experiments conducted by social workers be of high quality, provide
35 useful information, and appear in prestigious journals? Yes, indeed!
36 Other examples of social workers using this design can be found in Larsen
37 and Hepworth (1982), and in Sze, Keller, and Keller (1979).

1 These posttreatment-only comparison designs are exceedingly useful,
 2 but they can be markedly improved by adding formal pretreatment
 3 assessments to more rigorously ascertain whether the groups' function-
 4 ing changed over time, and if posttreatment differences are observed, to
 5 help rule out the possibility that they may be attributable to existing but
 6 unknown pretreatment differences. The next section explores some
 7 common variations on using this methodological refinement.

8 DESIGNS WITH CONTROL GROUPS AND PRETESTS AND POSTTESTS

9 The Pretest–Posttest No-treatment Control Group Design

10 This design is used to help determine if a given intervention produces
 11 any effects above and beyond those attributable to the passage of time,
 12 concurrent history, or the experience of being assessed. A diagram depict-
 13 ing this design shows:

14 $O_1 - X - O_2$

15 $O_1 \quad O_2$

16 We have two (or more) groups, not known to be equivalent on
 17 all possible factors, since the groups were not created using random
 18 assignment. Both groups are assessed at about the same point in time.
 19 The members of one group receive an intervention, whereas the mem-
 20 bers of the second group do not. Then, both are assessed at about the
 21 same point in time on a second occasion, after the first group receives
 22 intervention. If the treatment group changes and the no-treatment group
 23 does not, there is some modest logical justification to infer that it was the
 24 treatment, X , that produced these improvements.

25 In practice, with this and other designs incorporating control or com-
 26 parison groups, it may be logistically difficult to accumulate enough cli-
 27 ents at a single point in time so as to conduct the pretest assessment on
 28 everyone in both groups at roughly the same time. The pragmatic solution
 29 is often to use a *rolling enrollment protocol*, wherein new cases are added to
 30 each group over time, with perhaps months transpiring as a client is

1 enrolled in one group, assessed, treated, then reassessed, and the informa-
2 tion added to the dataset. Meantime, other clients who are not destined to
3 receive X are assessed, experience a similar length of time transpiring as do
4 clients receiving treatment, and are then reassessed. These clients' data too
5 are banked, until sufficient numbers of clients have accumulated in each
6 group, at which point the totality of the data are analyzed statistically. This
7 approach is scientifically less satisfactory than measuring all participants
8 at about the same time, since many different elements may come into play
9 during the time period covered by rolling enrollments and assessments,
10 elements that may change crucial features of an agency's operation, or that
11 may broadly affect the locale, state, or nation (and hence the features of
12 those enrolling in the project).

13 This pretest–posttest no-treatment control group design is a very
14 popular approach to evaluation research. Cook and Shadish (1994,
15 p. 566) go so far as to claim that “the most frequently employed quasi-
16 experiment still involves only two (nonequivalent) groups and two mea-
17 surement waves, one a pretest and the other a posttest measures on the
18 same instrument.” Here are some examples of evaluation research that
19 made use of this approach.

20 *Evaluating School Social Worker–Teacher Collaboration*

21 School children frequently experience academic, attendance, and behav-
22 ioral difficulties, and school social workers can be asked to help such chil-
23 dren. Viggiani, Reid, and Bailey-Dempsey (2002) evaluated the outcomes
24 of implementing a program wherein social work interns were placed in
25 elementary school classrooms for an entire day for 2 days a week to help
26 the teachers manage the class and resolve student difficulties before they
27 turned into a crisis. This was done in one kindergarten and one third-
28 grade classroom (total N = 36 children). Two comparable classes (total
29 N = 40) did not receive the added services of the social work interns.
30 Outcome measures related to student attendance, behavior, and grades
31 were obtained from student report cards pre- and post-intervention.
32 The researchers hypothesized that all three would improve in the two
33 “treated” classrooms, compared to the two classrooms that did not
34 receive the intervention. The treated and untreated classrooms' students
35 were statistically comparable in terms of gender and numbers. At the end
36 of the school year, children in the classes with the addition of the school
37 social worker for 2 days a week had significantly improved attendance

1 relative to the no-treatment classrooms, and achieved statistically signifi-
 2 cant improvements on 4 of 14 behavioral measures. Grades did not
 3 appear to be affected in any way. These mixed results were described by
 4 the authors as promising.

5 **The Pretest–Posttest Alternative-Treatment Control Group Design**

6 This design may be used to more rigorously test whether or not Treat-
 7 ment X produces outcomes different from those emerging from Treat-
 8 ment Y. The design can be diagrammed as follows:

9 $O_1 - X - O_2$

10 $O_1 - Y - O_2$

11 By now, you should understand what this outline means. Two groups of
 12 clients, in the natural course of events, differ in that some receive treat-
 13 ment X and some get treatment Y. Y may be an alternative legitimate
 14 treatment, such as TAU, or it may be a placebo treatment, something
 15 that provides exposure to the experience of being “helped,” but entails
 16 something that researchers believe will not simply be innocuous but
 17 actually ineffective. If some form of pretest and posttest assessments are
 18 available, it may be possible to make some inferences about the relative
 19 value of X versus Y in terms of how selected outcome measures are
 20 impacted.

21 *Evaluating a College-based Student Drinking Reduction Program*

22 Abusive drinking by students is a serious problem on some college cam-
 23 puses. Many students are underage, and any alcohol consumption by
 24 them is illegal. Being apprehended by law officers while consuming or in
 25 possession of alcohol, or using a fake ID to gain entry into bars, can result
 26 in a criminal record. Student alcohol use is associated with problems such
 27 as DUI arrests, vehicle accidents, alcohol poisoning up to and including
 28 death, fights, a higher risk for sexual activity and unprotected sex, and so
 29 forth. Many campus administrations are trying to reduce illegal and abu-
 30 sive use of alcohol through various approaches. One such approach being
 31 widely touted as effective is called a *social norms marketing approach*,

1 wherein a campus-wide marketing campaign is used to convey accurate
2 information about the modest amounts of alcohol actually consumed by
3 local college students, in the hope that by making moderate drinking or
4 alcohol abstinence be seen as the “norm,” students will moderate their
5 own alcohol intake. Faculty with the School of Social Work at San Diego
6 State University used a pretest–posttest comparison group design to eval-
7 uate an intensive social norms drinking reduction campaign undertaken
8 at their own university.

9 Students in one residence hall were exposed to the social norms
10 campaign, and students in another residence hall received a much less
11 intensive intervention: a 120-page booklet containing information about
12 alcohol-related laws and policies. The social norms approach included
13 posters, stickers, bookmarks, and notepads containing normative mes-
14 sages (e.g., “Seventy-five percent of students drink 0, 1, 2, 3, or 4 drinks
15 when they party”). The social norms intervention lasted 6 weeks. A total
16 of 476 students were exposed to the social norms approach, and 486 to
17 the informative booklet. Pre- and posttest surveys asked students about
18 their perception of normative drinking on their campus and about their
19 own drinking during the last 4 weeks. The surveys were completed anon-
20 ymously. Students in the two residence halls were equivalent in terms of
21 age and class standing. The outcomes? “[S]tudents exposed to the social
22 marketing campaign reduced their misperceptions of drinking norms but
23 drank *more frequently* at posttest than did their counterparts in the com-
24 parison group. The campaign had *no effect* on several other drinking indi-
25 cators . . . the frequency of drinking *actually increased* significantly over
26 time within the experimental group, while declining in the comparison
27 group” (Clapp, Lange, Russell, Shillington, & Voas, 2003, pp. 413–414,
28 emphases added). This was certainly an unexpected finding and, of
29 course, very disappointing.

30 *Evaluating Two Different Methods of MSW Instruction*

31 The opportunity to use this design occurred one semester when I was with
32 the University of Georgia School of Social Work. I was scheduled to teach
33 a foundation research MSW class, and a colleague was also slated to teach
34 a different section of the same course. I usually teach using a method of
35 instruction that places heavy reliance on using structured study questions.
36 I rarely prepared formal lectures, and I did not use midterm exams, final
37 exams, or term papers. My colleague taught using a different structured

1 pedagogical strategy, one called *learning with discussion*. We both had
2 similar syllabi in terms of course description, objectives, and the same
3 assigned textbook. We conducted a reliable and valid assessment of
4 the students' abilities to critique published social work research (one of
5 the common course objectives) at the beginning of the semester. We then
6 taught our courses using our respective and preferred method of instruc-
7 tion. At the beginning and again at the end of the term, we administered
8 a comparable, equivalent assessment of their critical research skills, assess-
9 ments that were blindly graded, with the two independent graders not
10 knowing which section the student assessment paper came from or if the
11 assessment paper was prepared at the beginning or end of the term.
12 We then broke the code for the students' grades and time of term (begin-
13 ning or end), and tabulated the pre and post-course grades for each sec-
14 tion. Both sections of students scored equivalently during the first
15 assessment. At the end of the semester, my students had made statistically
16 significant improvements, whereas my colleagues' students' grades had
17 not changed much. These results could be interpreted to mean that my
18 method of instruction is superior to the standardized approach that she
19 used, and of course that was the slant I placed upon them when this study
20 was published (Thyer, Jackson-White, Sutphen, & Carrillo, 1992)! (I hope
21 that the concept of allegiance effects occurred to you while reading the
22 previous sentence.) Of course, we also discussed some alternative inter-
23 pretations of the data and limitations of the study. By having me and the
24 other instructor not get involved in the scoring of the critical essays, by
25 having the essay graders blind as to which section an essay they were grad-
26 ing came from, or even if it was an essay completed at the beginning versus
27 the end of the semester, we tried to control for experimenter bias/alle-
28 giance effects. By using two independent graders, with neither grader
29 knowing how the other grader scored an essay they were reading, and
30 calculating interrater agreement (very high), we tried to control for test-
31 ing effects. Regression to the mean was not applicable in this study, stu-
32 dent attrition (mortality) was low and comparable between the two
33 sections, and concurrent history was at least partially controlled for by
34 conducting the pre- and postassessments at the beginning and end of the
35 same semester for both sections of the class. Selection bias (e.g., maybe
36 the smarter students gravitated to one instructor vs. the other) was par-
37 tially controlled for by showing that, during the pretest, the students in
38 both sections scored comparably poorly. So, all in all, we had a reasonably

1 tightly controlled study, even without random assignment and a true
 2 experimental design. I hope the concept of low-hanging fruit occurred to
 3 you while reading the previous paragraph.

4 **The Switching Replications Design**

5 This design is used to strengthen the ability to make a causal inference by
 6 trying to demonstrate, not just once but twice, that clients receiving
 7 X got better. Look over the diagram below and see if you can figure out
 8 the inferential logic behind this design:

9 $O_1 - X - O_2 \quad O_3$

10 $O_1 \quad O_2 - X - O_3$

11 Two groups of clients are processed a bit differently. The top group is
 12 assessed, immediately enrolled in treatment/therapy/intervention, and
 13 assessed again. During this same time period, the bottom group is simi-
 14 larly assessed but *does not* get treatment right away. Some time passes, and
 15 then the bottom group is assessed a second time, after which they begin
 16 the same treatment as the clients in the top group received. After treat-
 17 ment is completed for the second group, they are assessed a third time, as
 18 is the top group, as a form of follow-up measure, to see if treatment gains
 19 (if any) have been maintained. Ideally, following the implementation of
 20 this design, you would like to see a data pattern something like this: Both
 21 groups were essentially similar at O_1 . At O_2 , the treated top group would
 22 demonstrate significant improvements, and the bottom untreated group
 23 would not have improved. Then, at O_3 , the bottom group would have
 24 improved following receipt of treatment X to the degree seen in the top
 25 group at O_2 , while the top group would be shown to have maintained their
 26 treatment gains. In effect, this design strives for a *replicated* effect, with not
 27 just one but two demonstrations that clients improve following treatment X.
 28 The term *switching replications* refers to the fact that the top group is
 29 switched with the condition received by the bottom group initially (that of
 30 no-treatment), while the bottom group is switched to the initial condition
 31 received by the top group (namely, treatment X). This approach is also
 32 known as a *delayed treatment design* or as a *lagged-groups design*.

1 In the world of causal inference, science typically views results that
2 have been replicated as more credible than one-off effects, those docu-
3 mented with a single demonstration. The history of behavioral and social
4 science (and the natural and medical sciences too, for that matter), con-
5 tains many reports of marvelous discoveries reported in a single study,
6 findings that, although marvelous, could not be replicated by others. It is
7 this experience of being fooled, so to speak, that raises replicated findings
8 above the hoi-polloi of the effects observed but a single time.

9 This switching replications design was used in part by clinical social
10 worker Betsy Vonk in her evaluation of the outcomes of counseling pro-
11 vided to students at Emory University in Atlanta, services provided via
12 the university's Student Counseling Center where Betsy was employed.
13 In the normal ebb and flow of the Center's operation, not all new clients
14 could be enrolled in counseling right away. Some, due to a lack of coun-
15 selors with an open slot, had to be placed on a waiting list. In due course,
16 those on the waiting list were contacted and asked to make an appoint-
17 ment for a second assessment (a reliable and valid pencil-and-paper
18 measure of psychological symptoms), after which they could begin coun-
19 seling. All clients, when treatment was terminated, were asked to com-
20 plete the outcome measure again. Thus, Betsy had access to a naturally
21 occurring dataset that conformed to the parameters of the switching rep-
22 lications quasi-experiment. It is classified as a quasi-experiment because
23 the clients were not deliberately or randomly assigned to either the
24 immediate treatment condition or to the waiting list. It happened natu-
25 rally. Had assignment to the two conditions been truly random, say, on
26 the basis of a coin toss, then Betsy's study would rise to the level of a
27 genuine experiment. Nevertheless, at the time when she published her
28 report on this project (Vonk & Thyer, 1999), this quasi-experimental
29 design represented the most methodologically advanced outcome study
30 in the field of evaluating university student counseling programs avail-
31 able in the published literature.

32 **Dismantling Studies**

33 The purpose of a dismantling study is to try to determine the relative
34 contribution of one or more individual components of a social work
35 intervention that contains multiple elements. Typically, one group of
36 clients receives the "complete package" and another group receives the

1 complete package *minus one discrete element*. The logic is that if the two
 2 groups demonstrate equivalent results, the program element that was
 3 omitted is really not necessary to the entire package's success. Alternatively,
 4 if the group that received the entire program minus the one element dis-
 5 plays an impaired outcome, the extent of the impairment reflects the
 6 additive value of the missing element. One way of diagramming this type
 7 of study is as follows:

$$8 \quad O_1 - X - O_2$$

$$9 \quad O_1 - X_{-1} - O_2$$

10 wherein X represents the complete program and X_{-1} receipt of the pro-
 11 gram with one discrete element omitted.

12 An example of this type of quasi-experimental design is provided in
 13 Johnson and Stadel (2007), working in the field of hospital social work.
 14 The practice issue is getting patients (or their legal guardians) to provide
 15 what are called "health care proxies" prior to their receipt of elective
 16 orthopedic surgery. A health care proxy document designates someone
 17 who can make medical decisions on the patient's behalf in circumstances
 18 wherein the patient cannot make decisions about health care on their own
 19 (e.g., in a coma). At the hospital where this study was conducted, social
 20 workers introduced the option of executing a health care proxy by one of
 21 two methods. The complete package, so to speak, involved the social
 22 worker conducting a semi-structured face-to-face interview with patients
 23 who were scheduled for surgery and their families, informing them about
 24 the concept of a health care proxy *and* giving them written information
 25 and blank health care proxies to complete and return, if they so chose.
 26 Twenty-one patients received this complete package, labeled X. An addi-
 27 tional 36 patients were provided information about health care proxies
 28 solely by means of the written materials, without the face-to-face inter-
 29 views with the social worker. This latter group represented the compari-
 30 son condition, the X_{-1} group. The outcome measure was the percent of
 31 patients in each group who actually executed and turned in to the hospital
 32 a health care proxy prior to their surgery. The results? Forty-three percent
 33 (9 of 21) of the patients receiving the social work interview *and* the infor-
 34 mation completed proxies, versus only 6% (2 of 36) of those who received

1 the information alone. This quasi-experimental evidence certainly sup-
2 ports the hypothesis that the combined program is superior to the
3 information-alone intervention; to the extent that hospitals wish patients
4 to complete health care proxies, using social workers to conduct these
5 additional face-to-face educational interviews is an effective approach,
6 relative to providing information alone. This is good news for those who
7 advocate for the valued-added nature of providing social work services in
8 health care settings. But the good news is qualified in that these interviews
9 could likely be undertaken by non-social workers as well (e.g., by physi-
10 cians, nurses, or patient representatives) and by our awareness that the
11 nonrandom assignment of patients to the complete package versus the
12 partial one precludes an uncritical acceptance of the causal link between
13 complete treatment and outcome. Still, it is a good study, taking advantage
14 of naturally occurring differences in the interventions received by clients.

15 At this point, we can address the two final examples of important
16 research questions that quasi-experimental designs can address in the
17 evaluation of social work practice.

- 18 • Question 4. *What is the status of clients who have received a novel*
19 *treatment compared to those who received the usual treatment or*
20 *care?*
- 21 • Question 5. *What is the status of clients who have received a novel*
22 *treatment compared to those who received a credible placebo*
23 *treatment?*

24 These questions can be better answered using the designs presented in
25 the latter part of this chapter, than by the earlier ones presented here and
26 in Chapter 2. What is needed to answer Question 4 is a group of clients
27 who received the usual treatment or care, and to answer Question 5, indi-
28 viduals who received a placebo treatment. It may go against the grain of
29 social work to contemplate deliberately providing an intervention known
30 to be a placebo, but such a refinement is really necessary to come to grips
31 with the essential question “Is what we are doing better than nonspecific
32 influences?” Social worker Margaret Blenkner (1962) addressed this by
33 quoting Rosenthal and Frank (1956, p. 300), who observed:

34 [I]mprovement under a special form of psychotherapy cannot be taken
35 as evidence for: a) correctness of the theory on which it is based; or

1 b) efficacy of the specific technique used, unless improvement can be
2 shown to be greater or qualitatively different from that produced by the
3 patient's faith in the efficacy of the therapist and his technique. . . . To
4 show that a specific form of psychotherapy . . . produces results not
5 attributable to the nonspecific placebo effect it is not sufficient to com-
6 pare its results with changes in patients receiving no treatment. The only
7 adequate control would be another form of therapy in which the patient
8 had equal faith . . . but which would not be expected by the theory of
9 therapy being studied to produce the same effect.

10 Blenkner (1962, p. 58) went on to grapple bluntly with this core issue
11 for social workers:

12 Are we psychologically capable of entertaining the unpleasant idea
13 that workers can be placebos, and that our precious mystique—the
14 worker–client relationship—may be only the ubiquitous placebo effect?
15 Are we willing to give up our . . . prejudices long enough to find out
16 whether it is possible that regardless of theory, school, diagnosis, client
17 symptoms, or worker conceptualizations, if a worker has enthusiasm
18 and conviction about his way of helping, most clients will *feel* helped and
19 some will even *be* helped? If we are willing to do this we may finally get
20 to the really effective factors in technique and method and begin to
21 justify our claims to having a science-based art. (1962, p. 58, emphases in
22 original)

23 If a study is conducted that shows clients fare well after they received
24 a novel social work intervention, and a later study shows that clients who
25 received this intervention actually improved following treatment com-
26 pared to pretreatment levels of functioning, this is a good thing and
27 consistent with the idea that the intervention “works.” If further quasi-
28 experiments show that these gains are durable and that they compare
29 favorably with clients who received TAU, this too is a good thing to
30 know. But our enthusiasm for this remarkable new treatment should be
31 tempered by the knowledge that these impressive results are also consis-
32 tent with the hypothesis that both the novel treatment and the estab-
33 lished treatments produce improvements solely via placebo influences.
34 To truly show that the novel treatment works *better* than placebo, and

- 1 better than TAU, it must be compared not only against TAU, but also
 2 against a credible placebo. This requires a placebo comparison group
 3 design, perhaps formatted something like:

Novel Treatment Group	$O_1 - X - O_2$
Treatment as Usual Group	$O_1 - Y - O_2$
Placebo Treatment Group	$O_1 - Z - O_2$
No-treatment Group	$O_1 \quad O_2$

- 4 The no-treatment control group is needed to control for history effects,
 5 the passage of time, regression to the mean, and the like. Only if the novel
 6 treatment is statistically and clinically superior to no treatment, placebo
 7 treatment, *and* to treatment as usual can we have confidence that its
 8 effects are above and beyond those of existing care, placebo influences,
 9 and the passage of time alone. Some 50 years after Blenkner (1962) issued
 10 her challenge to the profession, such a well-designed placebo-
 11 controlled study has yet to be undertaken by the social work profession.
 12 Surely, it is time.

13 SUMMARY

- 14 This chapter has described in the abstract a series of progressively more
 15 elaborate and sophisticated quasi-experimental designs using variations
 16 on the theme of comparing novel treatment to no treatment, to com-
 17 parison treatments, and to placebo treatments. The design's abstract fea-
 18 tures were followed by presenting a series of actual, published studies
 19 illustrating their use. By adding pretreatment assessments and perhaps
 20 repeated pretests and posttests—including lengthy follow-up periods—
 21 these designs can be markedly improved upon to the point that they are
 22 capable of providing legitimate contributions to the empirical knowl-
 23 edge base of the human services. In some cases, these designs represent
 24 the only practical method of community-based research in environ-
 25 ments in which it is not possible to randomly assign clients to various

1 conditions of treatment or no treatment. In these instances, properly
2 controlled quasi-experiments represent the highest available form of evi-
3 dence that can be used to answer important questions, and these studies
4 can find publication outlets in some of the more prestigious journals in
5 the human services.

1

4

2

Interrupted Time Series Designs

3



Intervention research can be conducted using a variety of methodological approaches. Those studies involving large numbers of participants and analyzed (usually) with inferential statistics are known as nomothetic designs. These designs typically assess clients just a few times (e.g., pre- and posttreatment), maybe with an added follow-up period or two, under varying conditions. In other words, a large number of people are studied, but not very intensively. An alternative approach is variously called *single-system research designs*, *since-case research designs*, or *ideographic research*. This approach involves gathering data more intensively on a very small number of people, under varying conditions.

Both approaches are seen as viable methods to produce useful knowledge about the effects of social work interventions, but the scale is different. If a nomothetic study uses a sample of clients randomly selected from a larger population of interest, apart from the capacity to generate internally valid conclusions about the effects of treatment on your sample of clients, it may be legitimate to infer the effects of that same treatment if applied to other samples from the same populations, or even perhaps to the population itself. However, very few social work intervention research studies are able to obtain truly random samples of clients. Most rely on *samples of convenience*, which means that the generalizability of

1 most nomothetic findings is compromised. This is an important point—
2 no matter the sample size, 10, 100, or 1,000—generalizability is not legit-
3 imate unless the sample was obtained using true methods of random
4 selection. Thus, the major method by which generalizability is inferred is
5 via *replication*, conducting similar studies to see if initial findings can be
6 duplicated, then expanding the variety of clients, clinicians, and settings
7 in which the initially promising intervention is applied.

8 In single-system research, instead of, say, studying 30 people with one
9 or two assessments, one person is studied on 30 occasions, perhaps with 15
10 being prior to treatment (a baseline phase) and 15 after treatment (the
11 intervention phase). With proper design elements (e.g., replications with
12 several clients, systematically removing or introducing the same treatment)
13 it may well be possible to develop genuine causal inferences about the effects
14 of a given intervention for a particular client or small set of clients. However,
15 external validity is also very limited as one's "sample" of clients is typically
16 very small and not representative (in a statistical sense) of all clients with
17 a given condition or problem. Hence, one tries to generalize the results
18 from initially promising single-system designs via the same technique used
19 in nomothetic studies; namely, replications. Most social work research text-
20 books now contain chapters on the methodology of single-system designs
21 (Thyer, 2010a; Thyer & Myers, 2007; Rubin & Babbie, 2008; Yegidis,
22 Weinbach, & Myers, 2011), reflective of the acceptance of this approach
23 since its introduction to the field of social work in the mid-1960s.

24 The designs featured in the present chapter reflect a hybrid form of
25 intervention research strategy, one using the large numbers characteristic
26 of nomothetic studies while also incorporating the large numbers of
27 repeated assessments associated with idiographic research methods with
28 very small numbers of clients. These designs can be quite strong, although
29 most examples are classified as quasi-experiments since clients are not ran-
30 domly assigned to varying groups or conditions. Collectively, these are
31 labeled *interrupted time series designs*, and are very often used to evaluate
32 large-scale social welfare policies, as well as individual program outcomes.

33 THE INTERRUPTED TIME SERIES DESIGN

34 With the interrupted time series (ITS) design, the inferential logic is rela-
35 tively simple. An outcome measure assessing some variable of interest is

1 repeatedly measured a number of times prior to the initiation of an
 2 intervention, X. Then, after X is implemented or has occurred, assess-
 3 ments of the outcome measure are continued. This design can be dia-
 4 grammed as follows:

5
$$O_1 - O_2 - O_3 - O_k - X - O_{k+1} - O_{k+2} - O_{k+3} - O_{k+?}$$

6 with O_k referring to the total number of pretreatment assessments under-
 7 taken and $O_{k+?}$ being the final number of separate observations taken
 8 *after* the implementation of intervention X. The series of measurements
 9 taken prior to intervention is sometimes called a *baseline*, although this
 10 term is more appropriately used to refer to single-system research designs.
 11 These preintervention assessments provide us with a means of either
 12 intuitively or statistically projecting the level and slope of posttreatment
 13 measures, assuming that X did not happen. To the extent that posttreat-
 14 ment observations deviate from those “predicted,” so to speak, on the
 15 basis of the baseline measures, we may be able to infer an actual effect
 16 from X. In single-system research, this logical inference is usually made
 17 on the basis of visually inspecting the data. However, in the ITS design,
 18 inference is usually augmented through the use of specialized inferential
 19 statistics, with one of the more common tests being called *time series*
 20 *analysis*, or TSA. Do not confuse TSA with time series designs; TSA is a
 21 method of statistically testing for changes in data patterns pre- and
 22 postintervention, within the context of analyzing the outcomes of a time
 23 series design.

24 The most common method of TSA is known as the *univariate Box-*
 25 *Jenkins interrupted autoregressive integrated moving average* (ARIMA)
 26 analysis, which compares the values of some variable before the intro-
 27 duction of an intervention versus after that intervention. ARIMA takes
 28 into account the slope, not the just the average, of the values found in
 29 a series, and corrects for a mathematical problem known as *autocorrela-*
 30 *tion* (or *serial dependency*), wherein the value of a variable at one point
 31 in time can be used to predict the values of other variables in the series.
 32 The presence of significant autocorrelation within the time series
 33 data violates a major assumption of many parametric statistical tests;
 34 namely, that the data are independent. For example, the scores of a group
 35 of unrelated people who complete the Beck Depression Inventory (BDI)

1 are likely not dependent— knowing one person's scores does not help
 2 predict how someone else in that group will have scored. However,
 3 if someone completes the BDI on a weekly basis for many months,
 4 knowing their score on a given date *may* help in predicting their nearby
 5 scores in the series. This violates the important property of a lack of serial
 6 dependency and thus compromises the test to some extent. Similarly,
 7 weekly or monthly state-level data (say, numbers of clients receiving food
 8 stamps) violate the assumption of independence. Knowing the numbers
 9 for 1 month *does* help predict scores on subsequent months. ARIMA
 10 approaches compensate for this, something that simply using a t-test or
 11 analysis of variance (ANOVA) may not do, when used to examine average
 12 changes across phases (see McDowall, McCleary, Meidinger, & Hay, 1980).

13 A recently proposed alternative method of inferential analysis for
 14 interrupted TSA is known as *latent growth curve modeling* (Duncan &
 15 Duncan, 2004a, b), which has the advantage of being applicable with
 16 fewer than the minimum of 50 data points recommended for use
 17 with the Box-Jenkins approach to TSA. However, this newer method
 18 does not yet seem to have been applied by social work researchers using
 19 time series designs.

20 When using visual analysis to attempt to make inferences, one can
 21 look for:

- 22 • Changes following the introduction of the intervention, in terms
 23 of absolute magnitude (do the data obviously change up or down
 24 right away?)
- 25 • The slope of the graphed data (a steeper slope indicates a greater
 26 rate of change, which can be important)
- 27 • The direction of the data (do the data reverse the direction
 28 taken in the baseline? That is, change from rising to falling,
 29 or vice versa)
- 30 • Does the amount of variability in the data increase or decrease
 31 following the implementation of the intervention?

32 Visually obvious changes in one or any combination of these poten-
 33 tial patterns are a pretty good indication that *something* is different.
 34 This is a conservative test; thus, if you cannot simply see it, the real effect
 35 of the intervention, if any, is probably weak and not practically impor-
 36 tant. Using inferential statistics enables one to more reliability detect

1 small changes that are not visually obvious, so researchers end up learn-
2 ing more about weak interventions.

3 Cook and Shadish (1994, p. 562) describe these types of designs as
4 follows:

5 In interrupted time series, the same outcome variable is examined over
6 many time points. If the cause–effect link is quick acting or has a known
7 causal delay, then an effective treatment should lead to change in the
8 level, slope or variance of the time series at the point where treatment
9 occurred. The test, then, is whether the obtained data show the change
10 in the series at the prespecified point. . . . Internal validity is the major
11 problem, especially because of history (e.g., some other outcome-causing
12 event occurring at the same time as the treatment) and instrumentation
13 (e.g., a change of record keeping occurring with the treatment).

14 Time series designs like this are not infrequently used in the investi-
15 gation of psychosocial interventions provided to various groups of
16 clients, but they are more often employed in the analysis of possible
17 changes induced by social policy and in the field of community-wide
18 interventions (Biglan, Ary, & Waagenaar, 2000). Their use within social
19 work has been reviewed by Tripodi and Harrington (1979), DiNitto
20 (1983), DiNitto, McDaniel, Ruefli, and Thomas (1986), and Bowen and
21 Farkas (1991). These designs can be very powerful in terms of internal
22 validity. “Although considered quasi-experimental, the ITS design has
23 been noted as representing one of the strongest alternatives to the ran-
24 domized experiment” (Duncan & Duncan, 2004a, p. 271). The following
25 paragraphs describe some specific examples of using ITS designs as a
26 quasi-experimental approach to evaluating social interventions.

27 A hypothetical example of a simple ITS design is presented in Fig-
28 ure 4.1. The 15 or so preintervention data points are relatively stable,
29 and after the intervention is introduced, there is an abrupt discontinuity
30 in the data, which is maintained over the next 15 or so data points.
31 The more resistant to change the data could be projected to be, based
32 on prior research or good theory, the stronger the inferences that can
33 be made that change really did occur and that this change was due to
34 the intervention. For example, HIV-related deaths would be a more
35 difficult outcome to reduce than say, teenage attitudes toward Brittany
36 Spears.

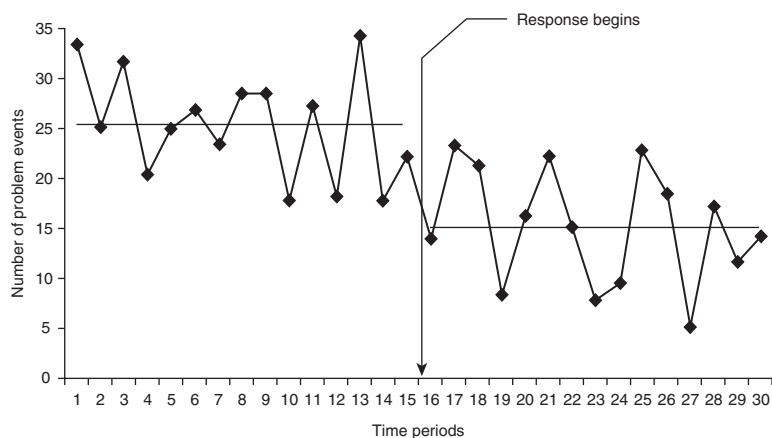


Figure 4.1. Hypothetical Example of a Simple Interrupted Time Series Design.

1 Server Intervention to Reduce Alcohol-related Traffic Accidence

2 In 1986, Oregon introduced legislation requiring specialized training
 3 for all alcohol servers (bartenders, waiters, waitresses, etc.), teaching
 4 them to recognize intoxication, how to politely stop serving intoxicated
 5 patrons, and how to encourage a customer to call a taxi, use designated
 6 drivers, etc. The intent was to reduce alcohol-related traffic accidents.
 7 By the end of 1989, when most servers had been so trained, the state
 8 found statistically significant reductions in single-vehicle nighttime traf-
 9 fic crashes (those with the highest percentage of alcohol involvement),
 10 which was interpreted as support for the effectiveness of this piece of
 11 social legislation (see Holder & Wagenaar, 1994).

12 Controversial Studies Using Time Series Analysis

13 Alcohol not to your taste? How about evaluating changes in laws pertain-
 14 ing to reporting and sentencing sexual assault? If that is of interest to
 15 you, read Schissel's (1996) TSA of such laws. Want to read something
 16 really controversial? Check out Berk, Sorenson, Wiebe, and Upchurch's
 17 study (2003); they used ITS designs to evaluate the presumptive effects of
 18 more liberalized abortion laws on homicides among young people aged
 19 15–24 years old. Their conclusion? "We conclude that the 1990s decline

1 in the homicide of young men is statistically associated with the legaliza-
2 tion of abortion” (Berk et al., 2003, p. 45). These authors may be right,
3 inferring that the decline in the numbers of people born is related to
4 increased availability of abortion, which caused a reduction of the size of
5 a more criminally prone age cohort. But they really can’t be sure, due to
6 the possibility of other explanations accounting for the decrease in crime.
7 Still, the use of an ITS design is a great way to provide for initial tests of
8 hypotheses, no matter how provocative.

9 **Does Raising the Legal Drinking Age Reduce Teenage Traffic Accidents?**

10 What happens when states raise or lower the minimum drinking age,
11 in terms of alcohol-related traffic accidents among young drivers? This
12 question lends itself very nicely to the ITS design since states already
13 gather such accident data, and it is relatively easy to determine when
14 the independent variable (enforcement of the new law) goes into effect.
15 Not surprisingly, raising the minimum drinking age is clearly followed
16 by a reduction in alcohol-related traffic accidents among young drivers
17 (Wagenaar & Toomey, 2002). Although such a conclusion would seem
18 both logical and self-evident, many well-intended social policies have
19 unintended consequences, some of which are harmful, and it is always
20 valuable to obtain actual *post-policy data* to confirm or disconfirm the
21 anticipated effects of new policies and laws.

22 **Does Banning Smoking in Public Buildings Reduce** 23 **Psychiatric Emergency Room Visits?**

24 Kurdyak, Cairney, Sarnocinska-Hart, Callahan, and Strike (2008) wished
25 to evaluate the possible effects of a smoking cessation policy on visits to
26 psychiatric emergency rooms, initially at one center in Toronto, and then
27 when the smoking ban was implemented province-wide. Psychiatric
28 emergency room visits at a specific hospital in Toronto were recorded
29 from March 1, 2002 to December 31, 2006. On September 21, 2005, the
30 specific hospital imposed a smoking ban (e.g., no smoking was allowed
31 in or near the hospital buildings), and on May 31, 2006, the ban was
32 imposed across the entire province for all public buildings. In a nutshell,
33 the hospital-specific ban on smoking had no effect on psychiatric ER
34 visits at that particular hospital, but when the smoking ban was extended

1 province-wide to all public buildings, ER visits by individuals with a
2 psychotic disorder dropped by over 15%. Rates of smoking are very high
3 among individuals with a psychotic disorder, and the authors expressed
4 the concern that this reduction might increase adverse psychiatric events
5 (attempts at self-harm, suicide, etc.). In the authors' words: "Our find-
6 ings suggest that if a smoking cessation policy is implemented in a psy-
7 chiatric emergency department setting consideration must be given as
8 to whether this will disadvantage some patient groups or populations.
9 The smoking cessation policy may act as a barrier to crisis services in
10 people with psychotic disorders" (Kurdyak et al., 2008, p. 782).

11 **Do Sex Offender Registration and Community Notification** 12 **Laws Reduce Sex Crimes?**

13 Few offenses evoke as much outrage as do sex crimes. As a result, many
14 jurisdictions have enacted mandatory registration and community noti-
15 fication policies for sex offenders. Sex offenders released from prison
16 must notify law enforcement officials as to their place of residence,
17 and law enforcement must in turn notify the community as to where
18 registered sex offenders reside. These laws are intended to help prevent
19 the reoccurrence of sex crimes by past offenders. Do they? Sandler,
20 Freeman, and Socia (2008) examined this hypothesis by obtaining the
21 criminal history files of every offender arrested for a registrable sex
22 offense in New York state between 1986 and 2006. This involved over
23 160,000 different individuals and over 170,000 sex offense-related arrests.
24 In January 1996, the state of New York enacted a Sex Offender Registra-
25 tion Act (SORA) requiring registration and community notification
26 related to sex offenders, with this law serving as the study's independent
27 variable or intervention. There were about 10 years of data available
28 before the law was passed, and about 11 years after the enactment of the
29 New York state SORA. A number of crime statistics were recorded,
30 including registrable sex offenses (e.g., rape, incest, sodomy, child mole-
31 station, etc.) assessed individually by type of crime and the totals. Arrests
32 were separated according to whether they involved a registered sex
33 offender or an individual not previously convicted of a sex crime. The
34 authors found no positive effect for the SORA—there were no significant
35 effects on total sexual offending, rape, or child molestation. "This finding
36 casts doubts upon the ability of sex offender registration and notification

1 laws, as well as residency and occupational restriction laws, to actually
 2 reduce sexual offending” (Sandler et al., 2008, p. 297). Although disap-
 3 pointing, it is still good to know whether certain laws are working as
 4 intended. Considerable resources go into implementing SORA laws,
 5 resources that could perhaps be diverted to more effective law enforce-
 6 ment policies related to public safety and sex offenses.

7 The above designs are not usually sufficient to completely rule
 8 out rival hypotheses, perhaps those arising from the threats of concur-
 9 rent history, the passage of time, general changes in the population of
 10 interest, and the like. Quite simply, something apart from the introduc-
 11 tion of X may have occurred at about the same time as X was imple-
 12 mented, and it could be this other “something” that produced any
 13 observed changes. One way to attempt to control for such threats is
 14 found in the next design.

15 **THE INTERRUPTED TIME SERIES DESIGN WITH A REMOVAL PHASE**

16 This design uses inferential logic similar to that found in the switching
 17 replications design discussed earlier. Some outcome measure is repeat-
 18 edly assessed, then an intervention, X, is introduced. Assessments are
 19 conducted in the same manner after X is in effect, for a given period
 20 of time, then X is *removed* and assessments continue as before. If X is
 21 truly effective, when it is introduced, there should be a change in the
 22 outcome measures. If X is genuinely responsible for any improvements
 23 (or perhaps deterioration), and X is subsequently removed, then the
 24 post-(not X) condition’s data should reflect a change in the data so that
 25 the data revert to the level, slope, and variability seen during the first
 26 baseline phase, prior to the introduction of X. This design can be dia-
 27 grammed as follows:

28
$$O_1 - O_2 - O_k - X - O_{k+1} - O_{k+2} - O_{k+n}$$

$$- (\text{not X}) - O_{k+n+1} - O_{k+n+2} - O_{k+n+p}$$

29 Here, k refers to the final measurement taken during the first baselines,
 30 $k + n$ the final measurement taken when X is in effect, and $k + n + p$ the
 31 final measurement taken after X was withdrawn. If a sufficient number of

1 data points are gathered for each phase of this design, the outcome mea-
2 sures are reliable and valid, and the effects of X are prospectively
3 predicted to be temporary, then this is a very powerful design indeed,
4 especially when the unit of analysis is something large, like a county,
5 state, or nation. It is much more difficult to change large systems (think
6 of the Titanic) than very small ones, so if effects such as those described
7 above are observed, the internal validity of such a study is likely high,
8 since rival explanations are pretty implausible (e.g., history, regression,
9 maturation). Even placebo influences are unlikely when assessing large-
10 scale systems, since placebo factors generally occur at the level of the
11 individual (although phenomena such as mass hysteria do occur and
12 cannot be discounted—witness the so-called “Obama effect” following
13 President Obama’s election on the generally increased level of optimism
14 observed in the United States, even before he had undertaken any mean-
15 ingful initiatives. One unpublished study found that a white-black
16 achievement gap on a GRE-like test all but disappeared following Barak
17 Obama’s election (see Dillon, 2009). However, such powerful placebo-
18 like influences are rarely impactful on large-scale systems, and by using
19 this ITS design with a removal phase, placebo-like factors can usually be
20 dismissed.

21 THE INTERRUPTED TIME SERIES DESIGN WITH AN 22 EXTENDED INTERVENTION PHASE

23 Most of the discussion on the ITD design has assumed that the interven-
24 tion is a one-time event, with potentially lingering influences. Sometimes
25 the intervention is best construed as being applied repeatedly over a
26 period of time, and, prospectively, perhaps being predicted to lose any
27 influence once the intervention is discontinued. Such was the situation
28 found in Chu, Frongillow, Jones, and Kaye (2009), in the area of improv-
29 ing the dietary selection of students eating at university food service
30 operations. The study was conducted at Ohio State University’s dining
31 center for students. For 14 days pre-intervention, the researchers posted
32 simple descriptions of each day’s 12 hot entrées on the menu board and
33 monitored which meals were chosen by the students, unobtrusively
34 taking into account the nutritional values of each entrée (e.g., total
35 energy, serving size, fat, protein, and carbohydrates). Then, for 14 days,

1 the same descriptions of each entrée were posted, *along with* each entrée's
 2 nutritional information. The students' entrée selections continued to be
 3 monitored unobtrusively. In the last phase, lasting 13 days, the nutri-
 4 tional information was removed from the description of the entrées,
 5 leaving the students with the same information as during the first phase
 6 of the study. This study could be diagrammed as:

7 $O_1 \dots O_{14} - X_1 \dots X_{14} - O_{15} \dots O_{37}$

8 with 14 days of data pre-intervention (the initial baseline phase), fol-
 9 lowed by 14 days of the intervention, followed by 13 days of the baseline
 10 condition. The results?

11 We observed an immediate drop in the energy content of patrons' entrée
 12 selections from the first day of posting nutrition labels for entrées at the
 13 dining center; this drop was maintained throughout the treatment
 14 period. When nutritional labels were removed, patrons reverted to select-
 15 ing entrées with higher energy content relatively soon. These changes
 16 occurred without negative impact on overall sales and revenue for the
 17 dining center. (Chu et al., 2009, p. 2002)

18 The authors pointed to the many strengths of their study. Data were
 19 collected electronically via point-of-sale machines, so instrumentation bias
 20 was controlled for. The reversion to initial baseline patterns when the
 21 intervention was removed argues for the beneficial and causal effect of
 22 posting the nutritional information, the intervention was inexpensive to
 23 implement, and it did not reduce food service revenue. They linked their
 24 study to America's obesity epidemic and the possible value of laws mandat-
 25 ing the posting of nutritional information on franchise restaurant menus.
 26 There appeared to be a clear functional relationship between posting nutri-
 27 tional information and students' choosing lower caloric value entrées. This
 28 has obvious public health implications. Whether the nutritional informa-
 29 tion would continue to exert a positive influence over a longer period of
 30 time remains open to question. It is possible that patrons would habituate
 31 to them, and the signs would have less effect on their behavior.

32 Researchers in North Carolina implemented a program intended to
 33 encourage low-income, older (40+ years) minority women to obtain

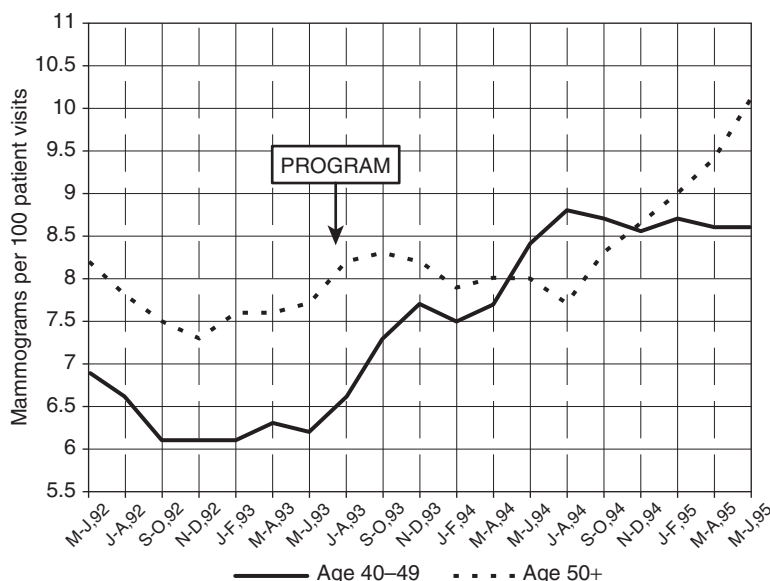


Figure 4.2. Mammograms Per 100 Patient Visit; 12 Month Moving Averages.

Source: Figure reproduced from Michielutte, Shelton, Parskett, Tatum, and Velez (2000, p. 619), with the permission of the publisher.

1 mammograms. A pretreatment phase looked at the numbers of mam-
 2 mograms obtained per 100 patient visits over about a 1-year period at a
 3 given facility, followed by about 2 years' worth of data following the
 4 introduction of a deliberate program designed to encourage mammo-
 5 grams. The data are depicted in Figure 4.2 and clearly show, visually, that
 6 mammograms increased after the program was implemented (see
 7 Michielutte, Shelton, Paskett, Tatum & Velez, 2000, for a full report of
 8 this project).

9 THE REPLICATED INTERRUPTED TIMES SERIES DESIGN

10 Our confidence in the conclusions drawn from any individual study are
 11 enhanced if the effects, positive or negative, can be replicated. With ITS
 12 designs, it is sometimes possible to examine the effects of an intervention

1 applied concurrently to two or more systems (e.g., counties, states), as in
 2 the instance when a statewide (with counties) or national (with states)
 3 policy initiative is introduced. Finding out that, following some policy
 4 becoming law, positive effects were observed in one state, such as Florida,
 5 is a good thing to know. However, Florida may be an idiosyncratic state,
 6 and the effects observed in Florida may not apply to other states. When
 7 measuring across two or more systems, when an intervention is applied
 8 to all these systems at the same time, if you observe congruent effects in
 9 both systems following the introduction of X, the confidence you have
 10 that X induced these changes is enhanced. One replication is good, two is
 11 even better, and so forth. This design may be diagrammed as:

12 State A $O_1 - O_2 - O_3 - O_k - X - O_{k+1} - O_{k+2} - O_{k+3} - O_{k+?}$

13 State B $O_1 - O_2 - O_3 - O_k - X - O_{k+1} - O_{k+2} - O_{k+3} - O_{k+?}$

14 Here, states A and B are measured similarly over the same time period,
 15 and intervention X (e.g., a national policy?) is introduced in both states
 16 at the same time. If both states change in similar ways following X, then
 17 we have more confidence that X caused these changes than if we tested
 18 for the effects in one state only.

19 This type of ITS design was used by Palmgreen, Lorch, Stephenson,
 20 Hoyle, and Donohew (2007) to evaluate the effects of the National Youth
 21 Antidrug Media Campaign (based on radio and television public service
 22 announcements), a 5-year, \$1 billion initiative of the federal government
 23 to prevent and reduce drug abuse among youth. The authors obtained
 24 interview-based data on anonymously self-reported recent marijuana
 25 use among adolescents in Lexington, Kentucky, and in Knoxville,
 26 Tennessee, monthly for 42 months before the start of this new federal
 27 campaign and for 6 months after it began. In both cities, reported mari-
 28 juana use significantly declined, leading the authors to conclude that the
 29 media campaign was effective and causally responsible for these decreases.
 30 In the authors' own words "We used data from a 48-month, independent
 31 sample interrupted time series project (one which tests trends before and
 32 after an intervention). . . . The interrupted time series design is one of the
 33 strongest quasi-experimental designs for inferring causal effects of an
 34 intervention" (Palmgreen et al., 2007, p. 1645).

1 **THE REPLICATED TIME SERIES DESIGN WITH A LAGGED**
 2 **INTERVENTION GROUP**

3 Another way to enhance the internal validity of an ITS is by assessing
 4 some outcome measure in two or more large-scale systems, introducing X
 5 into one system (e.g., State A) only, and seeing if hypothesized changes
 6 occur. However, X is *not* introduced into the other system(s) (e.g., State B),
 7 although these other systems continue to be monitored. If A changes and
 8 B does not, our confidence that X induced the observed changes in A is
 9 enhanced. If the same intervention X is subsequently introduced into
 10 State B, which was heretofore unchanged, and *then* State B changes in a
 11 manner similar to that observed previously in State A, then our confi-
 12 dence is greatly enhanced that X is causally responsible for the observed
 13 changes. This design can be diagrammed as follows:

14 State A $O_1 - O_2 - O_3 - O_k - X - O_{k+1} - O_{k+2} - O_{k+3} - O_{k+?}$

15 State B $O_1 - O_2 - O_3 - O_4 - O_5 - O_6 - O_k$
 $- X - O_{k+1} - O_{k+2} - O_{k+3} - O_{k+?}$

16 It is good for State B’s baseline to be a good bit longer than that of
 17 State A, but there are no hard and fast rules involved as to how much
 18 longer. Sometimes data may be collected daily, weekly, monthly, or even
 19 annually. The only principle is that the additional length of time should
 20 be of sufficient duration to provide a fair appraisal as to whether or not
 21 State B’s data remained stable, even after X was introduced into State A.
 22 It is not necessary that the postintervention data collection period
 23 between States A and B be identical in length, only that the available
 24 number of data points permits legitimate inferences.

25 **THE INTERRUPTED TIME SERIES DESIGN WITH A NO-TREATMENT**
 26 **CONTROL GROUP**

27 This design is an improvement over the first ITS study design described ear-
 28 lier in this chapter, and it improves upon this prior approach through incor-
 29 porating a no-treatment control series. A sequence of observations are made

1 of some outcome measure. Then an intervention, X, is introduced, and the
 2 series of observations are continued. As an added refinement however, a
 3 comparison group is similarly assessed during the same time period, but this
 4 second group is *not* exposed to X. This design can be diagrammed as below:

5
$$O_1 - O_2 - O_3 - O_k - X - O_{k+1} - O_{k+2} - O_{k+3} - O_{k+n}$$

6
$$O_1 - O_2 - O_3 - O_k \quad O_{k+1} - O_{k+2} - O_{k+3} - O_{k+n}$$

7 The exact number of pretest observations is noted as k , simply to indicate
 8 that this number can vary from study to study, and the final number of
 9 posttest observations, $k + n$, is also open. Again, note that k and $k + n$ do
 10 not have to be the same number of observations pre- and posttreatment,
 11 but it is helpful if they are the same, in terms of inferential symmetry and
 12 treatment by statistical tests. Helpful, but not essential.

13 This design can be used in, say, the investigation of state policies that
 14 are implemented in one state but not in another. Some outcome measure
 15 can be tracked over similar time periods in both states, then a law or
 16 policy change is introduced into the group depicted at the top of the
 17 diagram. However, the law or policy is *not* put into effect in the bottom
 18 state. The logic is that if changes are observed post-X in the top state, but
 19 not in the bottom, the inference is strengthened that it is the law or policy
 20 that was responsible for this change.

21 **Community Response to a Racist Murder**

22 On June 6, 1998, James Byrd, Jr., a 49-year-old African American,
 23 was walking home after attending a family event. Three white men
 24 offered him a ride but the cloaked gesture became apparent as they
 25 assaulted, savagely beat, and chained him to their pick-up truck, eventu-
 26 ally dragging him to his death. The trauma besetting the community
 27 grew intense in the days after the murder, as the severity of the crime
 28 quickly ignited a political, social, and media storm that gripped Jasper.
 29 (Wicke & Silver, 2009, p. 234)

30 This murder in Jasper, Texas, rocked the community, with virtual una-
 31 nimity among the local population regarding its heinousness. Wicke and

1 Silver (2009) examined the community-level response to this social
 2 trauma using an ITS design with a control condition—a similar com-
 3 munity in which the murder did not occur. Through theory and prior
 4 literature, the writers hypothesized that there would be several reactions
 5 to the murder in the areas of economic, criminal, and social indicators;
 6 specifically enhanced levels of cohesiveness and altruism; and a crum-
 7 bling of racial and class barriers. Archival data were taken from public
 8 records and other reliable sources of information during the years 1995–
 9 2003, with 42 observation points prior to the murder of James Byrd, Jr.
 10 and for 66 months after. The control community, Center, Texas, was
 11 similar in terms of size, racial and ethnic composition, geography,
 12 and economics. Graphs were prepared for various outcome measures
 13 within the two communities, before and after the murder, to see if
 14 Jasper’s data changed in ways following the murder not evidenced in
 15 Center. One such measure was the rate of violent crime in the two
 16 communities, and these data are presented in Figure 4.3. Apart from
 17 visual inferences, the authors used a statistical test called an *ordinary*
 18 *least squares regression* to see if the level and slope of the data lines
 19 changed, pre- versus posttest, and between the two communities. Violent
 20 crime did not significantly increase in the two communities during
 21 the 12 months prior to the murder in Jasper, but it did significantly

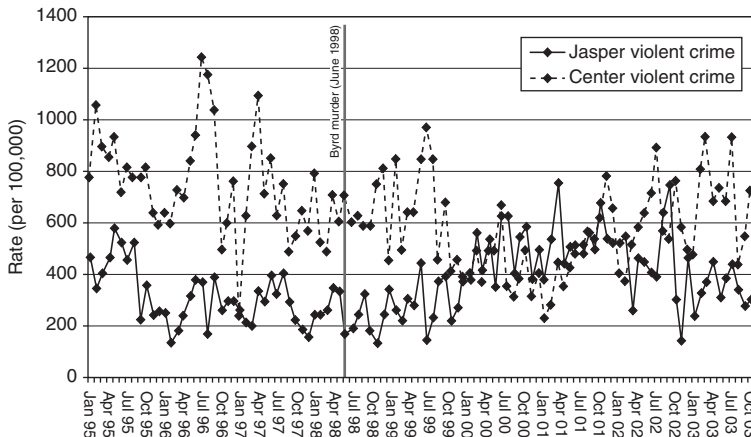


Figure 4.3. Monthly Violent Crime Incidents in Jasper and Center 1995–2003.

Source: Figure reproduced from Wicke and Silver (2009, p. 238), with permission.

1 increase after the first 12 months post-murder. The authors attempted
2 to explain this rise by linking their data to prior research (e.g., Raphael,
3 1986) on how communities respond to disasters: initially with an altruis-
4 tic response, followed by a “second disaster” as the first cohesive reac-
5 tions dissipate (the community rallies together) and the original existing
6 divisive forces among segments within the town reinstate themselves.
7 Here, the delayed reaction (an increase in violent crime) was predicted
8 in advance of its observed effect, which is actually more satisfactory
9 than changes observed immediately postintervention (in this instance,
10 the murder).

11 The authors’ overall conclusions:

12 [O]ur analysis identified several “negative” changes in Jasper in the
13 months and years following the murder. The divorce rate increased
14 and the housing market (as measured by the numbers of houses sold)
15 softened; both are negative indicators of community well-being. Jasper
16 also experienced an increase in violent crime and in its jail population.
17 While the results seem to indicate that Jasper changed for the worse after
18 the Byrd murder, the larger picture presented by the data suggests a
19 remarkable degree of resilience. . . . (Wicke & Silver, 2009, p. 244)

20 Do Longer Bar Hours Cause More Traffic Accidents?

21 In 1996, the provincial government on Ontario passed a law extending
22 the hours during which bars could remain open from 1 A.M. until 2 A.M.
23 Vingilis et al. (2006) examined alcohol-related motor vehicle casualties
24 and fatalities for the 4-year period prior to the passage of this Ontario
25 law and for 3 years afterward. They compared data trends in Windsor,
26 Ontario, with those obtained from Detroit, Michigan, located just across
27 the river from Windsor, with easy bridge and tunnel connections between
28 the two cities. In Windsor (with a legal drinking age of 19), a significant
29 increase in motor vehicle-related casualties was found following the
30 extended drinking hours. Detroit (with a legal drinking age of 21) found
31 a decrease. No similar trends were found for other areas of Ontario or
32 Michigan. Also, accidents in Windsor involving automobiles with
33 Michigan license plates increased, but there were no differences in
34 Ontario-licensed automobile accidents occurring in Detroit. The overall
35 pattern of the data indicated that Detroit drinkers were driving to

1 Windsor to take advantage of the extended drinking hours in bars,
2 as well as the younger drinking age, and that this was increasing accidents
3 in Windsor and actually decreasing them in Detroit.

4 **Do Mandatory Drivers' License Suspensions Reduce Drunk Driving?**

5 This method of analysis was employed by Wagenaar and Maldonado-
6 Molina (2007) to investigate the effects of suspending the drivers' licenses
7 of drivers involved in alcohol-related crashes. Such social policies
8 have been put into effect in 48 states, and they are intended as a deter-
9 rence policy: Don't drink and drive; if you do, you *will* lose your license.
10 These suspension policies have been put into place in various states
11 at various times, and this staggered pattern of implementation lent
12 itself very nicely to using the ITS design with a no-treatment control
13 group. These authors found that these policies do work; they are fol-
14 lowed by statistically significant and important reductions in alcohol-
15 related crash involvement, estimated to save about 800 lives per year in
16 the United States. Moreover, the authors were able to show statistically
17 that the rapidity of the punishment (time from arrest to trial and sen-
18 tencing) was more of a deterrent than the severity of the post-conviction
19 sentences.

20 The same approach was used by Wagenaar, Maldonado-Molina,
21 Erickson, Ma, Tobler, and Komro (2007) to examine the effects of DUI
22 fines and jail penalties on single-vehicle nighttime crash involvements,
23 and by Seekins et al. (1988) to see if legislation requiring the use of child
24 safety seats improved child safety when riding in automobiles.

25 **STRENGTHENING THE TIME SERIES DESIGN**

26 The usual suspects are involved in efforts to strengthen the internal valid-
27 ity of a time series design. Choose outcome measures of well-established
28 reliability and validity. Try to ensure that the intervention was really
29 delivered as intended and that the target group truly came into contact
30 with it. Phases containing larger numbers of data points are more
31 credible than those containing fewer. *Prospectively planned* ITS designs
32 have more scientific legitimacy than do *retrospectively conducted* ones.
33 With the former, one can develop predictive hypotheses in advance of

1 knowing the outcomes. Such predictions are inherently “riskier” than
2 those developed after the fact, perhaps after an examination of the data.
3 In the latter case, such studies may be little more than a “fishing expedi-
4 tion.” Cook and Shadish (1994, p. 562) emphasize the following:

5 Plausible threats are best ruled out by using additional time series.
6 Especially important are (a) control group time series not expected to
7 show the hypothesized discontinuity in level, slope, or variability of an
8 outcome; and (b) additional treatment series to which the same treat-
9 ment is applied at different times so we expect the obtained data to recre-
10 ate the known differences in when the treatment was made available.

11 Most of the ITS designs described in this chapter include variations
12 of the above elements to strengthen our confidence that any observed
13 changes can be legitimately attributable to the intervention. This is
14 admittedly an inexact science in that we strive to eliminate *plausible* rival
15 explanations, not possible but wildly implausible ones.

16 SUMMARY

17 Interrupted time series designs are widely used in the quasi-experimental
18 appraisal of the effects of social welfare, health, and public policy inter-
19 ventions. They particularly lend themselves to the analysis of archival
20 data typically maintained by city, county, and state governments, as well
21 as data gathered at the federal level. The simpler of the ITS designs
22 are capable of providing answers to very simple questions, whereas the
23 more complex and controlled designs may possess sufficiently high inter-
24 nal validity to justify (cautious) causal inferences about the effects of
25 interventions.



5

Evaluating and Reporting Quasi-Experimental Studies



This concluding chapter will touch on a number of additional elements to be considered in evaluating and undertaking quasi-experimental research designs used in the evaluation of social work practice. Included will be a discussion of ways in which the data resulting from such studies may be presented and analyzed, including descriptive and inferential statistics; the interpretation of negative outcomes; contemporary editorial standards that are increasingly being used by journals to help authors structure manuscripts reporting quasi-experimental studies; and an overview of ethical principles that must be followed in the conduct of these designs.

14 EVALUATING THE QUALITY OF QUASI-EXPERIMENTAL STUDIES

A wide variety of published tools are available to assess the quality and susceptibility to bias in quasi-experimental studies, with one recent review locating 53 checklists and 33 scales (Sanderson, Tatt, & Higgins, 2007). Thyer (1991) presented one checklist specific to social work for use by authors in evaluating and preparing research studies, and this is reproduced in Table 5.1. Most of the guidelines contained in Table 5.1

Table 5.1 Guidelines for Assessing the Adequacy of Reports on Research

Introduction

1. Does the report appropriately cite earlier, relevant studies drawn from the social work and other disciplinary literature?
2. Does the introduction conclude with one or more explicitly stated testable hypotheses?

*Methods**Clients*

1. Is a clear, potentially replicable description provided of the sampling procedure used to recruit clients for the study?
2. Are salient characteristics (demographic, clinical, diagnostic, etc.) of the sample of clients described in detail to permit comparisons of this sample with those used in prior (and future) studies?
3. Is a description provided of the nature of the informed consent process used to obtain client agreement to participate in the study?

Outcome Measures

1. Did the outcome measures employed in the study possess acceptable levels of reliability and validity?
2. Were the outcome measures *clearly* pertinent to the target problem?
3. Did the outcome measures possess treatment validity?

Intervention

1. Is the intervention program (treatment) described in sufficient detail to permit replication? If not, does the author provide a source to obtain a treatment manual or more explicit description of the intervention?
2. Were measures taken to assess practitioner compliance with intended interventions? If so, were the interventions carried out as intended?
3. If blind conditions were imposed on clients or practitioners (or both), were measures taken to assess the integrity of the blind nature of the study participants?

Research Design

1. Do the authors provide a clear description of the research design employed?
2. If the clients were assigned to various conditions, is the nature of this assignment process described in sufficient detail to permit replication?

(Continued)

Table 5.1 (Continued)

3. Are pretreatment measures taken of the clients' problems, strengths, or situation? If so, were the groups of clients assigned to differing experimental conditions roughly equivalent to each other pretreatment?

Results

1. Are the results obtained from the various outcome measures consistent with one another? Is the pattern of improvement (or deterioration) clear across all outcome measures?

2. Are the results presented in the form of graphs or tables? If so, are the data comprehensible without recourse to the narrative text?

3. If the results are presented in the form of descriptive statistics, is each mean accompanied by a standard deviation?

4. If inferential statistics are employed, are the data shown to meet the assumptions the tests are based upon (e.g., normal distribution, similar standard deviations, no significant autocorrelation, etc., in the case of parametric tests)?

5. If correlational measures are employed, are the *N*, correlation coefficient, and alpha level reported for each such analysis?

6. If a *t*-test or analysis of variance is used, does the report of each such test contain the degrees of freedom, actual *t* or *F* value, and alpha level?

7. If a statistically significant difference is found, is it accompanied by an appropriate effect size?

8. If multiple inferential statistical tests are performed, are the alpha levels appropriately adjusted to account for the numbers of such tests?

9. Apart from statistically significant changes and effect sizes, is the *clinical* significance of any improvements discussed?

Discussion

1. Does the author clearly address alternative explanations (e.g., threats to internal validity) for the results, apart from the hypotheses that were tested?

2. Does the author report only conclusions supported by the data? Are speculations clearly described as such, rather than as facts?

3. Are suggestions to improve future research in this area described?

4. Are clear *applications* to practice derived from this study described, with special reference made to the unique aspects of *social work* practice?

Note: Adapted Thyer, B. A. (1991). Guidelines for evaluating outcome studies on social work practice. *Research on Social Work Practice*, 1, 88–89. Copyright 1991 by Sage Publications, Inc.

1 will be familiar to the reader, as they have been addressed earlier in the
 2 present volume. The rationale for each suggestion should seem evident,
 3 and this listing remains pertinent some two decades later.

4 Holosko (2006, pp., 452–453) also provides a checklist for use
 5 by authors who are preparing an article manuscript for submission to
 6 the journal *Research on Social Work Practice*, a portion of which is
 7 reproduced here:

8 **Method**

9 A. Sample

- 10 • Are the techniques used in the sample selection process specified?
- 11 • Is the time frame for sampling specified?
- 12 • If there are other unique features of the sample, are they
- 13 mentioned?

14 B. Design

- 15 • Is the type of study design mentioned?
- 16 • Is the time frame to complete the study mentioned?

17 C. Data Collection

- 18 • Do you specify where and how data were collected?
- 19 • Do you mention the ethical considerations of data collection,
- 20 including whether institutional review board (IRB) approval
- 21 was obtained or why it was not necessary?

22 D. Outcome Measures

- 23 • Are all outcome measures used in the study specified and
- 24 referenced, as appropriate?
- 25 • Do you comment on the reliability and/or validity of the
- 26 outcome measures?

27 E. Intervention

- 28 • Is the intervention described in sufficient detail to permit
- 29 replication, or are citations provided to primary sources
- 30 fully describing the intervention?

31 **Results**

32 A. Statistics

- 33 • Are all statistics used to analyze the data mentioned?
- 34 • Do all inferential tests include levels of significance and/
- 35 or are effects sizes or proportions of variance accounted for
- 36 (if appropriate)

1 These questions are intended to be answered in a yes or no manner,
2 with authors asked to go back and address any issue that is not responded
3 to in the affirmative. Making sure that all of this information is included in
4 any research write-up of the results of a quasi-experimental study will go
5 far toward ensuring some consistency in reporting the essential informa-
6 tion needed to properly understand and appraise a given investigation.

7 Additional guidelines regarding the preparation of social work arti-
8 cles can be found in Thyer (2002, 2008), Schilling et al. (2005) and
9 in Holden et al. (2008). However, for practical purposes, contempo-
10 rary social work researchers contemplating designing, writing up, and
11 attempting to publish a quasi-experiment can focus on only two major
12 sets of criteria, the STROBE statement and the Journal Article Reporting
13 Standards produced by the American Psychological Association (APA).

14 **THE STROBE STATEMENT**

15 STROBE stands for *Strengthening the Reporting of Observational Studies*
16 *in Epidemiology*. “The STROBE statement was developed to assist authors
17 when writing up analytical observational studies, to support editors and
18 reviewers when considering such (quasi-experimental) articles for publi-
19 cation, and to help readers when critically appraising published articles”
20 (von Elm et al., 2007, p. 801). In the general field of health care, the term
21 “observational study” is often used in lieu of quasi-experimental study,
22 but the logic remains the same when it comes to evaluating interven-
23 tions. Observational studies are construed as having several methods,
24 one of which is known as the *cohort study*. Here is how one text describes
25 these designs:

26 A cohort study follows a group of people from one point in time to
27 another and observes changes that occur during that period. The study
28 can be retrospective . . . or prospective. . . Cohort studies may use
29 routine data or data specially collected for the purpose of the study,
30 or both. . . A cohort study is one in which subjects who . . . receive a
31 particular treatment are followed over time. They may be compared
32 with another group . . . without treatment. . . Cohort studies can be
33 very powerful. . . This type of research is useful for studying: the out-
34 come of treatment where a randomized controlled trial is impossible . . .

1 different approaches to service delivery and management when these
2 cannot be tested by a randomized controlled trial. . . . In addition, cohort
3 studies are also a useful means of studying “natural experiments” . . .
4 where different patterns of care exist in similar settings as a result of his-
5 tory or tradition. (Moore & McQuay, 2006, p. 163)

6 So, the term “cohort study” is used to describe what in health care is
7 commonly called an *observational study*. This latter term is used to indi-
8 cate that there was no deliberate manipulation of treatment assignments,
9 which is characteristic of randomized controlled trials (RCTs). Wherever
10 the STROBE statement refers to an observational study, just think of it
11 meaning a quasi-experimental design. All cohort studies are observa-
12 tional (or quasi-experimental, if you prefer) designs, but not all observa-
13 tional (or quasi-experiments) are cohort studies.

14 Another variety of observational study is called a *case-control study*.
15 In this approach, clients who have one outcome (e.g., recovery from
16 severe alcoholism) are compared with clients who did not recover from
17 alcoholism. Careful case histories are taken to try to ascertain if any dis-
18 tinct factors can be used to figure out what may have been responsible
19 for these disparate outcomes. If, for example, it was found in a large
20 sample of *former* heavy drinkers that most had become active partici-
21 pants in Alcoholics Anonymous (AA), and that in a demographically
22 similar group of heavy drinkers who did not stop drinking very few had
23 joined AA, the tentative hypothesis might be drawn that AA involvement
24 leads to sobriety. You can see why this example would be considered
25 very tentative evidence in terms of making any casual inferences because
26 of the possibility of controlled confounding factors.

27 A real-life example occurred to test the hypothesis, widely prevalent
28 at one point, that certain infant vaccinations triggered the development
29 of autistic disorder. Researchers in Denmark (where medical records are
30 very well maintained) examined the incidence of autism in 440,655 youth
31 who had received the vaccinations in question and compared them
32 to 96,648 nonvaccinated youth. In the former, the percent with autism
33 was 0.06, and in the latter 0.055, a negligible and non-statistically sig-
34 nificant difference. In this case, the case-control study was powerful evi-
35 dence indeed to disconfirm the hypothesis that infant vaccination causes
36 autistic disorder (see Moore & McQuay, 2006, p. 170). It is my impres-
37 sion that case-control observational (e.g., quasi-experiments) designs

1 are rarely used by social work researchers, relative to cohort-type obser-
 2 vational studies.

3 Social workers undertaking a quasi-experiment would be well advised
 4 to be familiar with the STROBE statement to be sure that no important
 5 features are omitted from consideration in either the study’s design or in
 6 the write-up of the final report or journal manuscript. A large number
 7 of journals (primarily biomedical in nature) have adopted the STROBE
 8 statement as a portion of their editorial policy, and the applications of
 9 this checklist to behavioral and social science research, including social
 10 work studies, are obvious and compelling.

11 The STROBE statement is supported by a website ([www.strobe-](http://www.strobe-statement.org)
 12 [statement.org](http://www.strobe-statement.org)) that goes into the details of these recommendations
 13 and provides a rationale for each. To date, no social work journal has
 14 appeared to incorporate the STROBE statement into its editorial policy.
 15 This does not diminish the usefulness of the guidelines, however, which
 16 are reported in their entirety in Table 5.2.

Table 5.2 STROBE Statement—Checklist of Items that should be Included in Reports of Observational Studies

	<i>Item</i>	<i>Recommendation</i>
	<i>No</i>	
Title and abstract	1	(a) Indicate the study’s design with a commonly used term in the title or the abstract (b) Provide in the abstract an informative and balanced summary of what was done and what was found
Introduction		
Background/ rationale	2	Explain the scientific background and rationale for the investigation being reported
Objectives	3	State specific objectives, including any prespecified hypotheses
Methods		
Study design	4	Present key elements of study design early in the paper
Setting	5	Describe the setting, locations, and relevant dates, including periods of recruitment, exposure, follow-up, and data collection

(Continued)

Table 5.2 (Continued)

Participants	6	<p>(a) <i>Cohort study</i>—Give the eligibility criteria, and the sources and methods of selection of participants. Describe methods of follow-up</p> <p><i>Case-control study</i>—Give the eligibility criteria, and the sources and methods of case ascertainment and control selection. Give the rationale for the choice of cases and controls</p> <p><i>Cross-sectional study</i>—Give the eligibility criteria, and the sources and methods of selection of participants</p> <p>(b) <i>Cohort study</i>—For matched studies, give matching criteria and number of exposed and unexposed</p> <p><i>Case-control study</i>—For matched studies, give matching criteria and the number of controls per case</p>
Variables	7	Clearly define all outcomes, exposures, predictors, potential confounders, and effect modifiers. Give diagnostic criteria, if applicable
Data sources/ measurement	8*	For each variable of interest, give sources of data and details of methods of assessment (measurement). Describe comparability of assessment methods if there is more than one group
Bias	9	Describe any efforts to address potential sources of bias
Study size	10	Explain how the study size was arrived at
Quantitative variables	11	Explain how quantitative variables were handled in the analyses. If applicable, describe which groupings were chosen and why
Statistical methods	12	<p>(a) Describe all statistical methods, including those used to control for confounding</p> <p>(b) Describe any <i>methods</i> used to examine subgroups and interactions</p> <p>(c) Explain how missing data <i>were</i> addressed</p> <p>(d) <i>Cohort study</i>—If <i>applicable</i>, explain how loss to follow-up was addressed</p> <p><i>Case-control study</i>—If applicable, explain how matching of cases and controls was addressed</p> <p><i>Cross-sectional study</i>—If applicable, describe analytical methods, taking account of sampling strategy</p> <p>(e) Describe any sensitivity analyses</p>

(Continued)

Table 5.2 (Continued)

Results

Participants	13*	(a) Report numbers of individuals at each stage of study—e.g., numbers potentially eligible, examined for eligibility, confirmed eligible, included in the study, completing follow-up, and analyzed (b) Give reasons for nonparticipation at each stage (c) Consider use of a flow diagram
Descriptive data	14*	(a) Give characteristics of study participants (e.g., demographic, clinical, social) and information on exposures and potential confounders (b) Indicate number of participants with missing data for each variable of interest (c) <i>Cohort study</i> —Summarize follow-up time (e.g., average and total amount)
Outcome data	15*	<i>Cohort study</i> —Report numbers of outcome events or summary measures over time <i>Case-control study</i> —Report numbers in each exposure category, or summary measures of exposure <i>Cross-sectional study</i> —Report numbers of outcome events or summary measures
Main results	16	(a) Give unadjusted estimates and, if applicable, confounder-adjusted estimates and their precision (e.g., 95% confidence interval). Make clear which confounders were adjusted for and why they were included (b) Report category boundaries when continuous variables were categorized (c) If relevant, consider translating estimates of relative risk into absolute risk for a meaningful time period
Other analyses	17	Report other analyses done—e.g., analyses of subgroups and interactions, and sensitivity analyses
Discussion		
Key results	18	Summarize key results with reference to study objectives

(Continued)

Table 5.2 (Continued)

Limitations	19	Discuss limitations of the study, taking into account sources of potential bias or imprecision. Discuss both direction and magnitude of any potential bias
Interpretation	20	Give a cautious overall interpretation of results considering objectives, limitations, multiplicity of analyses, results from similar studies, and other relevant evidence
Generalizability	21	Discuss the generalizability (external validity) of the study results
Other information		
Funding	22	Give the source of funding and the role of the funders for the present study and, if applicable, for the original study on which the present article is based

*Give information separately for cases and controls in case-control studies and, if applicable, for exposed and unexposed groups in cohort and cross-sectional studies.

Note: An Explanation and Elaboration article discusses each checklist item and gives methodological background and published examples of transparent reporting. The STROBE checklist is best used in conjunction with this article (freely available on the websites of *PLoS Medicine* at <http://www.plosmedicine.org/>, *Annals of Internal Medicine* at <http://www.annals.org/>, and *Epidemiology* at <http://www.epidem.com/>). Information on the STROBE Initiative is available at www.strobe-statement.org.

1 As with the guidelines found in Table 5.1, it is hoped that the sugges-
 2 tions provided by STROBE's checklist of items are self-evident. Just as
 3 even the most experienced pilots complete a checklist prior to taking
 4 off in an airplane, social work authors and research consumers can ben-
 5 efit from a careful review of the STROBE standards prior to submitting
 6 an article for publication.

7 JOURNAL ARTICLE REPORTING STANDARDS

8 Social work research will be more directly impacted by the new guidelines
 9 appearing in the sixth edition of the *Publication Manual of the American*
 10 *Psychological Association* (American Psychological Association [APA],
 11 2009, pp. 247–253), since most social work journal *do* follow the APA's
 12 publication guidelines. Several sets of documents are included in the

1 APA's new *Journal Article Reporting Standards* (JARS). One set of report-
2 ing standards includes information recommended for inclusion in all new
3 data-based manuscripts (pp. 247–248), and a second set includes report-
4 ing standards for studies using nonrandom assignment of participants to
5 experimental groups (p. 250); for example, in quasi-experiments.

6 The general recommendations are broken out by the section of the
7 manuscript: for example, Title and title page, Abstract, Introduction,
8 Method (Participant Characteristics, Sampling Procedures, Sample Size,
9 Power and Precision, Measures and Covariates, Research Design), Results
10 (Participant flow, Recruitment, Statistics and data analysis, Ancillary
11 analyses), and Discussion. Additional reporting standards for studies
12 using nonrandom assignment of participants to experimental groups
13 (e.g., quasi-experiments) are presented in the Methods section and
14 involve providing details as to the assignment method (e.g., unit of
15 assignment, as in individuals, groups, communities), Procedures used to
16 help minimize potential bias, Masking (e.g., whether or not those assess-
17 ing the outcomes were aware of the condition assignments), and Statistical
18 Methods (e.g., statistical methods used to compare study groups on pri-
19 mary outcomes).

20 It is an advance for the field to have such recommended standards
21 clearly established, as this will aid in the reporting and understanding of
22 the design and conduct of quasi-experimental studies. To the extent that
23 social work journal editors rigorously adhere to the latest publication
24 standards found in the APA's *Publication Manual* (and they should, for
25 consistency's sake), social work researchers will be guided by these new
26 standards and write up their quasi-experiments in a manner consistent
27 with these guidelines. The APA actually provides a very useful, free
28 multimedia tutorial on preparing papers in APA style (see <http://www.apastyle.org/learn/tutorials/basics-tutorial.aspx>), although it does not
29 cover the specifics of the JARS. Basically, every social work researcher
30 should acquire a personal copy of the APA's *Publication Manual* and
31 learn its guidelines well. This is an essential skill for any researcher.
32

33 THE COALITION FOR EVIDENCE-BASED POLICY GUIDELINES

34 The Coalition for Evidence-based Policy has issued guidelines titled
35 *Which Comparison-Group ("Quasi-experimental") Study Designs Are*

1 *Most Likely to Produce Valid Estimates of a Program's Impact* (see [http://coalition4evidence.org/wordpress/wp-content/uploads/2009/11/](http://coalition4evidence.org/wordpress/wp-content/uploads/2009/11/Validity-of-comparison-group-designs-updated-Nov09.pdf)
 2 [Validity-of-comparison-group-designs-updated-Nov09.pdf](http://coalition4evidence.org/wordpress/wp-content/uploads/2009/11/Validity-of-comparison-group-designs-updated-Nov09.pdf)). These offer
 3 the following standards for appraising the credibility of a given quasi-
 4 experiment in terms of its ability to yield credible answers:
 5

-
- *A number of careful investigations have been carried out to address this question.* In other words, more than a single quasi-experimental has been conducted to answer this question, and similar conclusions have been arrived at.
 - *The comparison-group designs most likely to produce valid results contain all of the following elements:*
 - (i) The program and comparison groups are highly similar in observable preprogram characteristics, including:
 - *Demographics* (e.g., age, sex, ethnicity, education, employment, earnings)
 - *Preprogram measures of the outcome the program seeks to improve* (at the beginning of the study, the groups are roughly equivalent in terms of their scores on the outcome measures)
 - *Geographic location* (both studies obtain participants from the same general area)
 - (ii) Outcome data are collected in the same way for both groups
(e.g., the same survey administered at the same point to both groups)
 - (iii) Program group and comparison group members are likely to be similar in motivation
 - (iv) Statistical methods are used to adjust for pretreatment differences between the two groups
 - *Preferably, the study chooses the program and comparison groups "prospectively" (i.e., before the program is administered)*
 - *The study follows the same practices that a well-implemented RCT follows in order to produce valid results (other than the actual random assignment):* For example, the study should have adequate sample size, use valid outcome measures, prevent cross-overs to or contamination of the comparison group, have low sample attrition, use an intent-to-treat analysis, and so on.
-

1 The similarities among the JARS, the STROBE statement, and the
2 guidelines provided by Thyer, Holosko, and the Coalition for Evidence-
3 based Policy, are all highly congruent, with few contradictory recom-
4 mendations. Each can be useful in appraising the adequacy of a written
5 report of the design and conduct of a quasi-experimental study. Given
6 the wide currency of the *APA Style Manual*, however, I believe the JARS
7 standards to be the single most useful resource in this regard.

8 REPORTING RESULTS

9 Basically, quantitative data may be presented in two ways: visually via
10 tables, graphs, and figures; and numerically via descriptive and inferential
11 statistics. Both approaches have their strengths and weaknesses, and can be
12 seen as complementary approaches, not competing ones. Various sources
13 describe how descriptive statistics should generally be reported (not spe-
14 cific to quasi-experimental designs). For example, Thyer (1991) suggested:

15 Whenever descriptive statistics are employed, they should be reported in
16 their entirety. It is common for articles to include information about
17 means (averages) but to omit the accompanying standard deviation.
18 A mean without a standard deviation does not allow the reader a clear
19 understanding of the variation present in the data, and it precludes other
20 scholars from conducting replication or secondary analyses of the data.
21 It is common for some analyses to not include all of the clients present
22 in a particular group. For example, there may have been dropouts from
23 the study or the client failed to complete all outcome measures. Thus,
24 a report of the exact numbers of clients should also accompany each
25 mean value and standard deviation, as well as any data presented in the
26 form of percentages. (p. 84)

27 It is important to pay particular attention to presenting relevant
28 demographic information on your sample(s) of clients. Their mean age,
29 races/ethnicities, diagnoses, gender, marital status, and other germane
30 characteristics should be reported, as this information is crucial to any
31 scientific investigation. It also aids other researchers in determining if
32 different studies made use of client samples that were relatively similar or
33 widely disparate from each other.

1 The above suggestions are a good start, but the APA *Publication*
2 *Manual* (2009) is perhaps the best resource you have in terms of learning
3 the conventions of presenting results. Here are some of the more perti-
4 nent guidelines contained in this manual:

5 In the results section, summarize the collected data and the analysis per-
6 formed on those data relevant to the discourse that is to follow. Report
7 the data in sufficient detail to justify your conclusions. Mention all rele-
8 vant results, including those that run counter to expectation; be sure
9 to include small effect sizes (or statistically nonsignificant findings) when
10 theory predicts large (or statistically significant ones). Do not hide
11 uncomfortable results by omission. Do not include individual scores
12 or raw data, with the exception, for example, of single-case designs
13 or illustrative examples. . . . When reporting the results of inferential
14 statistical tests or when providing estimates of parameters or effect sizes,
15 include sufficient information to help the reader fully understand the
16 analyses conducted and possible alternative explanations for the out-
17 comes of those analyses. Because each analytic technique depends on
18 different aspects of the data and assumptions, it is impossible to specify
19 what constitutes a “sufficient set of statistics” for every analysis. However,
20 such a set usually includes at least the following: the per-cell sample sizes;
21 the observed cell means (or frequencies of cases in each category for
22 a categorical variable); the cell standard deviations, or the pooled within-
23 cell variance. . . . For inferential statistical tests (e.g., t , F and X^2 tests),
24 include the obtained magnitude or value of the test statistic, the degrees
25 of freedom, and the probability of obtaining a value as extreme or more
26 extreme than the one obtained (the exact p value), and the size and
27 direction of the effect. When point estimates (e.g., sample means or
28 regression coefficients) are provided, always include an associated mea-
29 sure of variability (precision), with an indication of the specific measure
30 used (e.g., the standard error). (APA, 2009, pp. 32–33)

31 Few things are more frustrating for a reader than looking for some
32 important feature (say, the racial makeup of the clients in a given study)
33 and not being able to find it. Or, for a researcher wishing to undertake
34 a meta-analysis of a given report, to find that the standard deviations or
35 sample sizes are not included along with the mean values for a given
36 variable. Appropriately including relevant demographic and outcomes

1 information versus obsessively reporting irrelevant minutia in tedious
2 detail can be a difficult balancing act. Further guidance can be found in
3 reviewing particularly well done studies in areas relevant to the one you
4 are undertaking or reviewing, and learning from the models presented by
5 these published articles.

6 INTENT-TO-TREAT VERSUS EFFICACY SUBSET ANALYSIS

7 In social work research, the conventional practice has been to analyze
8 outcomes based on who actually receives a treatment. For example,
9 if in a pretest–posttest study, 100 clients are initially assigned to receive
10 a treatment, and 80 of the original 100 participants remain for the post-
11 test assessment (e.g., 20 dropped out), the common technique is to look
12 at the average score pretest (with $n = 100$), compared to posttest (when
13 $n = 80$). One may hope, ideally, that the group mean at posttest is statisti-
14 cally improved versus the pretest mean, and one may judge this as reflect-
15 ing the possible positive effects of intervention. A problem with this
16 method is that the drop-outs may alter the posttest mean scores in ways
17 *unrelated* to the possible effects of treatment. For example, if the most
18 impaired clients tended to be more likely to drop out, the group means
19 at posttest could look much improved, when in reality the changes in the
20 average score are really due to the most impaired clients being omitted
21 from the posttest analyses. This approach is known as *efficacy subset anal-*
22 *ysis* (although the term is not widely used) and, again, is the most
23 common way in which social work interventions are evaluated. A prob-
24 lem with this approach is that it introduces bias into the statistical analy-
25 sis and inflates the likelihood of making a type I error (concluding that
26 treatment exerted an effect, when in reality it did not).

27 One recommended way to compensate for this problem is to include
28 *all* participants at each assessment, regardless of whether they received
29 the intervention or not. This is called an *intention-to-treat* analysis, and is
30 attempted by obtaining the same outcome measures from those who
31 dropped out from treatment as from clients who completed treatment
32 (Lachin, 2000). In other words, in an intention-to-treat analysis, even if
33 20 clients had dropped out from treatment at posttest, the social work
34 researchers would try very hard to get them to complete the posttreat-
35 ment assessments, and to then include these data along with those

1 obtained from those who successfully engaged in treatment and com-
2 pleted intervention. Thus, the pre- and posttest comparisons involved
3 100 and 100 persons, not 100 and 80, respectively.

4 The more common approach to the problem of dropouts is for the
5 researchers to compare salient demographic and outcome measures
6 between those who completed treatment and those who dropped out.
7 If there are no statistically significant differences between the two groups,
8 the problem of attrition is considered to have been addressed, and it is
9 assumed that the completers did not differ from the dropouts in any
10 meaningful way. Another way is to impute missing data, using any of
11 a variety of methods of statistical legerdemain. The more Calvinistic
12 research methodologists consider these approaches (retrospectively
13 comparing completers vs. noncompleters, or imputing missing data)
14 to be a less than satisfactory resolution of the problem, however, and
15 urge adoption of intent-to-treat analysis as the most rigorous solution.
16 If this approach is followed, then one can analyze the data both ways,
17 using the intent-to-treat approach (including all original clients) or via
18 an efficacy subset analysis (limited to those who actually completed
19 treatment). An intent-to-treat analysis is a much more conservative
20 approach to analyzing the effects of an intervention. For example, receipt
21 of treatment may have many effects in addition to those directly assessed
22 by the outcome measures. Treatment could cause terrible side effects,
23 be psychologically very grueling, or cause a great deal of family or spousal
24 discord. This could be missed if the outcome measures solely looked
25 at presenting symptoms or problems, such as group mean scores on the
26 Beck Depression Inventory (BDI), or illicit drugs consumed. By includ-
27 ing treatment dropouts in the outcome analyses, the researcher will
28 likely obtain a more complete picture of how people react to a given
29 intervention, not simply how those who were able to complete it reacted
30 to it symptomatically. To my knowledge, no social work outcome
31 study has yet to include an intent-to-treat analysis as a part of its evalua-
32 tion design. It is likely that such intent-to-treat analyses will be intro-
33 duced and eventually become common within the social work outcomes
34 literature. Recall that most social work journals follow the APA's
35 *Publication Manual*. This guide now clearly states:

36 In studies reporting the results of experimental manipulations or
37 interventions, clarify whether the analysis was by intent-to-treat. That is,

- 1 were all participants assigned to conditions included in the data analysis
- 2 regardless of whether they actually received the intervention, or were
- 3 only participants who completed the intervention satisfactorily included?
- 4 Give a rationale for the choice. (APA, 2009, p. 35)

5 INTERPRETING NEGATIVE OUTCOMES

6 Research that fails to reject the null hypothesis, or which fails to find
7 a predicted difference (e.g., a given treatment was not followed by mean-
8 ingful improvements, or the treatment group did not improve its func-
9 tioning more than the nontreatment control or comparison group) is
10 often seen as a disappointment by the researchers. This may occur because
11 of one's personal investment in a project or a desire to see the research
12 hypotheses corroborated. We have a natural tendency to want to see posi-
13 tive results, results that truly do make a difference in people's lives, and
14 when this does not happen, we are disappointed. The consumers of
15 research—social work practitioners, administrators, policy-makers, and
16 clients—all wish to find results with a difference, something they can take
17 away and perhaps apply to their agency-based practice or policy, or use as
18 guidance in seeking out genuinely effective therapies. This is, of course,
19 understandable. From the viewpoint of the behavioral or social scientist
20 however, keep in mind that the purpose of research is to reveal the truth,
21 to discover lawful relationships in nature and among human beings. Social
22 work research is not intended so much to *prove* a particular point as it is
23 to simply discover what the point is, what exists. Ideally, research is not
24 driven by a quest to prove that a given hypothesis is true, but rather to
25 determine what the results *are*, positive, negative, or neutral. It is more
26 objective to state that the research question is “How do clients who receive
27 rectification therapy (RT) fare posttreatment, compared to clients who
28 received treatment as usual (TAU)?” as opposed to “Clients who receive
29 RT will display statistically and clinically greater improvements, compared
30 to clients who received TAU.” (As you'll recall from Chapter 1, rectifica-
31 tion theory is an imaginary treatment for juvenile delinquency.)

32 When faced with a negative result, for example a finding that a
33 given therapy did not prove to be effective, one can examine the study
34 to help determine the possible reasons for this outcome. There are
35 several possibilities. One is that, truly, the therapy really does not work.

1 This is the default position, since most therapies do not really prove help-
2 ful, relative to no treatment, to credible placebo treatments, or in the
3 long run. In other words, the null hypothesis usually *is* the true state of
4 affairs. However, any study whose conclusions support the null hypoth-
5 esis should be carefully examined to make sure that it *really is* a good
6 study, a well designed and fair appraisal. This is because any poorly
7 designed study can find no differences, and one's task, when interpreting
8 a finding of no difference, is to critically appraise the quality of the study
9 to be sure it has a reasonable chance of actually finding any differences,
10 if they existed. This brings up the other possibility alluded to above;
11 namely, that the study was so poorly designed and conducted that its
12 finding of no difference cannot be trusted as legitimate.

13 TYPE I AND TYPE II ERRORS

14 When drawing conclusions from an outcome study, one may make
15 a true conclusion or a false conclusion. If treatment really works, and you
16 conclude from a study that it does work, this is a true conclusion. If treat-
17 ment does not really work, and you conclude that it does not really work,
18 this is a true conclusion as well. However, if a treatment *really does not*
19 work, but on the basis of your study you conclude that it does work, this
20 is an incorrect conclusion, and this type of mistake has been labeled
21 a type I error. If treatment *really does* work, and you conclude on the
22 basis of a study that it does not work, this incorrect conclusion is called
23 a type II error. When one commits a type I error—claiming that some-
24 thing worked but it really did not—one unfairly promotes ineffective
25 treatments. When one commits a type II error—saying something did
26 not work when it really did—then genuinely effective treatments can be
27 overlooked or prematurely discarded. Both mistakes are problematic.
28 The diagram helps explain these concepts further.

29 Through using a conventional alpha level of $p < .05$, roughly speaking
30 about 1 in 20 scientific conclusions will consist of a type I error (claiming
31 a difference exists when it really does not). Science tries to reduce the
32 perpetuation of type I errors via replication. If you conclude, on the basis
33 of a study with one statistically significant difference on an outcome mea-
34 sure between a treatment and control group, that treatment is more effec-
35 tive than the control condition, you have a 1 in 20 chance of being wrong.

	Reality	
	<i>Treatment Does Not Work</i>	<i>Treatment Works</i>
<i>You conclude that treatment works.</i>	Type I Error	True Conclusion
<i>You conclude that treatment does not work.</i>	True Conclusion	Type II Error

1 If someone replicates this study, the chance of obtaining two similarly
 2 wrong conclusions is $.05 \times .05$, or $.0025$, a dramatically smaller risk, and
 3 this clearly shows why single studies should be replicated before a finding
 4 of an effect (e.g., treatment X “works”) should be accepted as legitimate.
 5 The chance of three type I errors like this is $.05 \times .05 \times .05$, or $.000125$, a
 6 very small risk indeed. The importance of replication has been asserted by
 7 Ziman (1978, p. 56): “The results of repetitions of the same experiments
 8 are fundamental to the creation of any body of knowledge.”
 9 And by Thomas (1975, p. 278):

10 The results of replication may be essentially positive, in which case
 11 confidence in the reliability of the procedures used is greatly increased.
 12 Indeed, each successful, positive replication increases plausibility multi-
 13 plicatively, because the chance occurrence of such results becomes much
 14 more improbable with each additional replication. The same may be
 15 said, of course, for replicated failures.

16 Going back further, we can turn to Chapin (1949, p. 135), who
 17 generously claimed:

18 For only by replication in numerous similar studies may we escape from
 19 the dilemma of whether the obtained differences were due to the non-
 20 randomness of the samples, or to the fact that they were drawn from
 21 different universes. . . . Should the same results be found on many trials,
 22 then generalization from even non-random samples to a universe might
 23 be valid and justified. (p. 135)

24 And

25 [T]hrough the replication of experimental design studies, which attempt
 26 to measure the effectiveness of specific means–ends schemes planned to

1 attain specific goals, it may be possible to develop a systematic mosaic of
2 nonrandom samples that will possess a degree of representativeness to
3 compensate for lack of randomization, and thus to supply a basic represen-
4 tiveness upon which reliable scientific generalizations may rest. (p. 139)

5 The mathematics linking type I and type II errors dictates that when
6 the likelihood of one type of error declines, the other goes up, and vice
7 versa. If you use a more stringent alpha level to determine statistical sig-
8 nificance (reducing the possibility of committing type I error), you make
9 it more difficult to “find” differences, and you may overlook real effects.
10 The risk of type I errors is that “findings” are discovered that are really
11 false. In intervention research, this means that some therapies are said
12 to be effective when they are really not. The danger of type II errors is
13 that effective treatments may be ignored. However, type II errors are in
14 some ways *less problematic*, because any true effects “missed” are liable
15 to be real, reliable in a statistical sense, but pragmatically small, exerting
16 little clinical impact. Given two groups of sufficient size, a treatment
17 group and a no-treatment group, a difference of 1% favoring the treat-
18 ment will be determined to be statistically significant, and indeed, it will
19 be real in a probabilistic sense. However, an intervention that reduces
20 clients’ average scores on something like the BDI by 1% is unlikely to be
21 a clinically useful treatment, even if the effect is genuine. Thus, a true
22 effect could be claimed from the study (treatment X statistically signifi-
23 cantly reduces BDI scores), but it would be a type I error, claiming a true
24 effect when one does not actually exist (that is clinically useful). Indeed,
25 our journals are filled with individual studies (maybe 1 in 20?) that
26 reached statistical significance, but in reality represent type I errors.

27 When a type II error occurs, when we miss a true effect, it is usually
28 because the effect is small and clinically unimportant. Thus, in the world
29 of intervention research, type II errors are less problematic in developing
30 a scientific knowledge base of meaningfully effective treatments than are
31 the plethora of type I errors, wherein teeny effects are artificially elevated
32 to importance when they are merely statistically reliable.

33 One common design problem that can yield an incorrect conclusion
34 that a given treatment is ineffective is using too few clients to adequately
35 support the statistical analysis (e.g., an underpowered study). Dattalo
36 (2007) provides an excellent overview on this topic. Other obstructions
37 to finding differences may involve the use of inappropriate or insensitive

1 outcome measures. For example, the Michigan Alcoholism Screening
2 Test-Revised (MAST) is a 22-item client self-report instrument designed
3 to detect alcohol abuse. Some of the items consist of questions whose
4 positive answers will be insensitive to change, such as:

5 6. Have you ever attended a meeting of Alcoholics Anonymous?

6 ___ Yes

7 ___ No

8 7. Have you ever gotten into physical fights when drinking?

9 ___ Yes

10 ___ No

11 17. Have you ever gone to anyone for help about your drinking?

12 ___ Yes

13 ___ No

14 A MAST score is based on the total of yes answers. You can see how
15 questions using the words “Have you ever . . .” would not change pre-
16 and posttest results. If a researcher used the MAST as a pretest and post-
17 test measure in a quasi-experimental study, even if treated clients became
18 completely sober, many of their responses to the MAST items would
19 not change. If participants’ scores on the MAST did not change between
20 the pre- and posttest (after treatment), one could conclude, perhaps
21 erroneously, that the treatment was ineffective due to this lack of change.
22 In reality, however, the measure was simply incapable of picking up on
23 any real changes in drinking habits.

24 Other possible flaws that can result in failing to find an effect of treat-
25 ment would be issues such as the therapist being incompetent to perform
26 the intervention he is charged with providing; the clients failing to attend
27 enough sessions or to otherwise adequately engage in treatment; or perhaps
28 a blurring of treatment conditions, wherein some clients are assigned to
29 receive treatment X only and others treatment Y only, but in the actual
30 conduct of treatment the therapists inadvertently or deliberately (perhaps
31 for what he considered to be sound clinical reasons) provided elements of Y
32 to clients assigned to condition X, or vice versa. This contamination would
33 also result in the finding of no differential effect between the two groups and
34 the erroneous conclusion that X and Y do not vary in their effectiveness.

35 Thus, when faced with a null finding, an important task is to criti-
36 cally review the article itself and see if the study was adequately designed.

1 If it is not, then it is useless because its negative findings cannot be relied
2 upon, since you do not know if they are a legitimate conclusion or a type II
3 error. But if the study is methodologically rigorous, then the nega-
4 tive findings are more likely to reflect the true state of affairs; namely,
5 that the treatment does *not* really work. And this is a good thing to
6 know.

7 How can it be a good thing to know that certain interventions or
8 practices *do not* work? Look over the examples below and ask if these
9 studies with negative results actually advanced the human services in
10 positive directions:

- 11 • Comprehensive summaries of the outcomes literature on social
12 casework were published by Fischer (1973, 1976), Segal (1972),
13 and Grey and Dermody (1972) finding that the available research
14 indicated that, for the most part, professional social work services
15 provided no positive effects, and some cases provided negative
16 ones. Over the next couple of decades, this led to a surge in
17 better-designed outcome studies, many of which documented
18 more positive results.
- 19 • In the face of widespread and inaccurate reports that early
20 childhood vaccinations caused autistic disorder, leading to
21 measurable declines in vaccinations and increases in vaccine-
22 preventable childhood diseases, comprehensive reviews of this
23 quasi-experimental evidence were published clearly showing
24 no credible link between autistic disorder and childhood
25 vaccinations (Smeeth et al., 2004; Demicheli et al., 2008).
- 26 • A comprehensive review of antidepressant drug trials,
27 co-authored by a licensed clinical social worker, found that these
28 powerful medications were not useful in the treatment of mild to
29 moderate depression, but were useful only in the most severe
30 cases, and that they are widely overprescribed (Turner, Matthews,
31 Linardatos, Tell, & Rosenthal, 2008). This important study
32 appeared in the prestigious *New England Journal of Medicine*.
- 33 • Although medications are widely prescribed for the treatment
34 of persons suffering from anorexia nervosa (AN), a large-scale
35 review of the available studies found that “Pharmacotherapy
36 provides little benefit in the treatment of AN at present”
37 (Crow, Mitchell, Roerig, & Steffen, 2009, p. 1).

- 1 • Facilitated communication (FC) is a widely used treatment to
2 try to help persons with severe autism and other developmental
3 disorders communicate via typing on a keyboard. A “facilitator”
4 holds the client’s hand over the keyboard as the client supposedly
5 pecks out words and sentences. Many thousands of persons were
6 trained in and provided this therapy. Careful investigations
7 revealed that the facilitators were unconsciously guiding the
8 typing, and that it was not being done by the client. It was a
9 manifestation of the *Ouija board effect*, not an effective therapy
10 (Herbert, Sharp, & Guidano, 2002). Leading professional
11 organizations have called for a ban on the clinical use of FC.
- 12 • One supposedly crucial component of treatment in the popular
13 psychotherapy called *eye movement desensitization and reprocessing*
14 (EMDR) includes having the client use her eyes to track the
15 therapist’s finger as it is waved back and forth in front of the
16 client. Considerable time and training go into doing these eye
17 movements “correctly.” Dismantling studies of EMDR provided
18 with and without these supposed crucial eye movements have
19 shown that they have no effect on outcomes (e.g., Carrigan &
20 Levis, 1999), thus undercutting the neurophysiological basis said
21 to be responsible for EMDR’s effects.
- 22 • So-called *reparative therapy* attempts to convert gay men or
23 lesbians to a heterosexual orientation. The National Association
24 of Social Workers (NASW) has determined that, thus far,
25 the available evidence indicates that reparative therapy does
26 not work. Because of this, as well as because of issues related to
27 respect for diversity and wanting to avoid pathologizing gay and
28 lesbian orientations, the NASW has issued a position statement
29 claiming that the practice of reparative therapy by social workers
30 is unethical and should not be provided to clients. Similar
31 statements have been issued by other human service organizations.
32 Is it important to be aware of the outcome studies demonstrating
33 that reparative therapy does not work? Obviously.
- 34 • Considerable health care resources, public and private, go
35 into providing clients with the treatment known as acupuncture,
36 which consists of inserting thin needles into precise positions
37 on the body called *meridians*. Much research shows that people
38 who receive legitimate acupuncture feel better. However,

1 a considerable number of experiments have been conducted
 2 comparing real acupuncture, involving the accurate placement
 3 of the needles into the correct meridians, versus fake
 4 acupuncture, in which the needles are placed into randomly
 5 chosen spots. Typically, both groups of patients improve
 6 equally, suggesting that acupuncture is essentially a powerful
 7 placebo treatment (Novella, 2011). Is it important to know
 8 if acupuncture exerts specific effects, other than placebo
 9 influences? Obviously.

10 The bottom line here is that research with negative outcomes or with
 11 findings of no difference can be quite important and valuable. This per-
 12 spective is emphasized in Holden et al. (2008, p. 68), who stated that
 13 “Neither reviewers nor editors should consider studies reporting nega-
 14 tive results as inherently inferior to studies reporting positive results. It is
 15 the conceptualization and conduct of the study, combined with the
 16 interpretation and write up of the results that are important—not the
 17 direction of the results.”

18 Quasi-experiments with negative results can serve to head off an ulti-
 19 mately unproductive line of research. This is why federal requests for
 20 proposals for large research grants frequently ask that investigators
 21 include any pilot data in the results of their grant application. Suppose
 22 one wished to do a large-scale randomized controlled clinical trial of RT
 23 as a treatment for juvenile delinquency, and it was proposed to test the
 24 effectiveness of RT by comparing its results against TAU for adjudicated
 25 kids within the community. The outcome measure is recidivism, which
 26 would be examined at 3, 12, and 24 months following the termination of
 27 treatment. The proposed study could be diagrammed as:

28 R N = 100 $X_{RT} - O_{1(3months)} - O_{2(12months)} - O_{3(24months)}$

29 R N = 100 $X_{TAU} - O_{1(3months)} - O_{2(12months)} - O_{3(24months)}$

30 Note that with the relatively large sample size there is no real need for any
 31 pretreatment measures. And with the lengthy follow-up periods, this
 32 would be seen as a very strong design to evaluate the relative effectiveness
 33 of RT versus TAU—but also a *very* expensive one. If there were no prior

1 evidence that RT produces any positive results, much less that it is any
 2 better than TAU, federal (and other) funders might be reluctant to spend
 3 the hundreds of thousands of dollars it could take to carry out this com-
 4 plex randomized controlled experiment. Remember, the default hypoth-
 5 esis most likely to be true *is* the null hypothesis—there will be no
 6 differences across time within groups or between groups during the
 7 various follow-up periods. Given this, the grant is unlikely to be funded.
 8 It is simply not a good bet for the funders.

9 Now imagine that the above grant proposal was accompanied
 10 by pilot data with positive results from a pre-experimental or quasi-
 11 experimental study, maybe one looking as simple as this:

$$12 \quad N = 100 \quad X_{RT} - O_{1(3\text{months})} - O_{2(12\text{months})} - O_{3(24\text{months})}$$

13 Here, we have a single imaginary group of 100 youth who received RT.
 14 Further imagine that the recidivism rate at 3 months was 4%; 6% at
 15 12 months; and 7% total at 24 months. Whoa! These are remarkably low
 16 recidivism rates for any intervention applied in the field of juvenile jus-
 17 tice, and for them to remain low for 2 full years posttreatment would be
 18 unheralded in the annals of delinquency research. With very strong
 19 results like this, when dealing with an intractable problem, the need for
 20 a no-treatment control group or for a comparison treatment condition is
 21 less stringent. Pilot data with strong results like this, even in the context
 22 of a simple posttreatment-only group design, really augments the legiti-
 23 macy of one's request for sizable funding to conduct an evaluation of far
 24 greater methodological rigor (and cost). This illustrates one of the
 25 strengths of quasi-experimental studies.

26 Take the converse. You are a researcher interested in examining the
 27 effects of RT on recidivism rates among juvenile delinquents, and you do
 28 the simple pre-experimental posttest-only study noted above and find
 29 that recidivism rates are 70% at 3 months, 80% at 12 months, and 85% at
 30 24 months. Even without a control group, these rates would be seen as
 31 quite high, and certainly not supportive the hypothesis that RT is an effec-
 32 tive intervention to reduce recidivism among delinquents. Faced with this
 33 disappointing information, you may well be inclined to drop further
 34 investigations into RT for juvenile offenders, and to forego any effort to
 35 design and seek funding for a large-scale RCTs. This is another strength of

1 quasi-experiments—they can serve as a filter or screen, useful in weeding
2 out the obviously ineffective and useless. By doing so, one can stop pursu-
3 ing lines of inquiry that will ultimately prove to be a dead end (thus saving
4 lots of time, energy, and professional disappointment). “A study not
5 worth doing is not worth doing well” (Holland, 1997, p. 2585).

6 Tremendously powerful interventions need neither control groups
7 nor statistical analysis to judge their effects. Also, issues that are well rec-
8 ognized as being relatively intractable, those that respond little or not at
9 all to placebo influences, that do not change with the passage of time,
10 that are not prone to maturation effects, and that tend to be relatively
11 steady-state, do not require control or comparison conditions to be use-
12 fully researched. Some examples might include individuals diagnosed
13 with obsessive-compulsive disorder (OCD), autistic disorder, severe
14 Down syndrome, or schizophrenia of the paranoid type. Now, this is not
15 to say that the natural history of these conditions suggests that the pic-
16 ture is completely hopeless, but it must be admitted that the vast major-
17 ity of persons with these conditions, reliably diagnosed, are unlikely
18 to dramatically improve absent something akin to a miracle (which
19 *can* happen, occasionally). So, for example, if one conducted a pretest–
20 posttest study with lengthy follow-up periods for 100 persons with one of
21 these disorders (say, OCD); found well-established, severe, and unabated
22 psychopathology lasting for years; then provided these individuals with
23 RT and found that immediately posttreatment and at 1 and 2 years later,
24 not a single person met the diagnostic criteria for OCD and for all intents
25 and purposes seemed “cured,” this could be a Nobel-prize-winning
26 study. No statistics and no control groups would be needed. Holland
27 expressed it this way:

28 One example of a trial that should not be a randomized controlled
29 trial is when the initial results are so striking and the database of prior
30 experience so uniform that the conclusion is inescapable. . . . Where the
31 observation represents a sea change, based on unmistakable objectivity . . .
32 the wisdom and experience of the observer reach the goal sooner and
33 with a shorter causality list. (Holland, 1997, p. 2585)

34 In a satirical article titled *Parachute Use to Prevent Death and Major*
35 *Trauma Related to Gravitational Challenge: A Systematic Review of*
36 *Randomized Controlled Trials*, Smith and Pell (2003) facetiously pointed

1 out that no RCTs existed demonstrating that using parachutes saves
2 the lives of people falling out of airplanes. Their point is that with
3 obviously powerful interventions, there is no need for RCTs or quasi-
4 experimental studies to demonstrate the value of the approach. Regret-
5 tably, for research purposes, such clear-cut situations are relatively rare
6 in the human services. Many of the psychosocial problems we address via
7 intervention research wax and wax in severity, and clients have a distress-
8 ing habit of sometimes getting better all on their own, without any pro-
9 fessional intervention (e.g., stopping smoking and abusive drinking,
10 losing weight, overcoming phobias or unemployment, leaving abusive
11 relationships and developing much more productive lives). This is good
12 for the persons concerned but makes the task of drawing legitimate causal
13 inferences about the effects of treatment much more difficult for the
14 poor researcher.

15 One tool that researchers have to help detect small but reliable effects
16 of interventions is known as *inferential statistics*, the use of various statis-
17 tical tests to determine if the results obtained from a study significantly
18 deviate from those expected by chance or random variation in the data
19 alone. With quasi-experimental studies, the usual purpose of inferential
20 statistics is to derive conclusions about the clients seen in *our* particular
21 research project, and not to try to generalize any results to larger popula-
22 tions of interest. For example, if you conduct a pretest–posttest study
23 on an Individual Development Account program for 50 poor families
24 that you were able to recruit from within your local community due
25 to their convenience, it is most likely that your 50 families are not some-
26 how “representative” of all poor people in your area. Therefore, you
27 cannot legitimately (in a scientific sense) extrapolate any conclusions
28 from your study to all local poor people. But you can use inferential
29 statistics to tell you if the mean amount of savings for your 50 families
30 significantly increased following the program. If you find out that it has,
31 and the savings are meaningful—not merely statistically significant—
32 as a practical matter, you could be tempted to apply this same program
33 to other local poor people and see if the positive findings can be repli-
34 cated. If this happens, each successful replication enhances your confi-
35 dence that, yes indeed, you do have an intervention that helps the poor
36 save, but this enhanced generalizability is based on successful replica-
37 tions with different groups of poor people, not from conducting the
38 original study on a representative sample of the local poor.

1 REPORTING RESULTS

2 Let's begin with simple descriptive statistics, since these are the most
3 appropriate way of describing the results of very simple studies. In the
4 case of the posttest-only design ($X - O_1$), the group of clients (individu-
5 als, families, couples, organizations, communities, etc.) is exposed to an
6 intervention and systematic data are obtained on client functioning
7 following receipt of the intervention. However, no pretreatment mea-
8 sures are formally taken. An example of this might be a group of middle
9 schoolers who received the Drug Abuse Resistance Education (DARE)
10 program, and then, some years later, their drug use is assessed. Here,
11 the data can be presented in a very simple descriptive manner. If, 5 years
12 later, 100% scored negative on a drug screening, this would be evidence
13 consistent with the hypothesis that the DARE program did help protect
14 kids from using drugs. If 95% scored positive for drug use, we could be
15 pretty sure DARE was not very useful in this regard. With less extreme
16 results, interpretation is more difficult because this design offers no
17 group to compare against the youth who received DARE. If 35% of the
18 kids who received DARE turned up positive for drug use, is this a rela-
19 tively good or bad result? Lacking information about the extent of drug
20 use from comparable kids who did not get DARE, it is difficult to tell if
21 DARE is protective or not. To some extent, such a determination is a
22 judgment call.

23 We may be able to evaluate the results of a given program offered in
24 the context of a posttest-only design if very strong public claims have been
25 made for the expected results of this program. For example, one authority
26 has claimed that EMDR is effective in reducing and even eliminating mil-
27 itary combat-related posttraumatic stress for 85% or more of patients
28 after only a few sessions. This is a very strong, even remarkable, claim, one
29 that may sound too good to be true. If a posttest-only study was done with
30 military combat veterans, and posttreatment posttraumatic stress disor-
31 der (PTSD) rates were found to be much higher than 15%, the strong
32 claims made by EMDR's proponents would be weakened.

33 An example of the purely descriptive analysis of a posttest-only study
34 is provided in Thyer (1988). I had adopted a new method of instruction
35 that I called *teaching without testing*. Basically, it involves having my stu-
36 dents bring their written answers to detailed study questions to each
37 weekly class they had with me, with the study questions being based on

1 that week's particular assignment. During class, I would call upon indi-
 2 vidual students for their answers to particular questions, often digressing
 3 to elaborate on some point or another and facilitating class discussion
 4 of each question. I graded each week's assigned set of questions and
 5 opted to not use a mid-term or final examination or term paper assign-
 6 ment, since the students were apparently working very hard every week.
 7 My impression was that this was a better method of promoting learning
 8 than using tests or term papers. Student had to read the assigned material
 9 each week, write out the content, and then engage in discussion about
 10 it during class. This made it very difficult to escape coming into close
 11 contact with the course materials. I used this method in several bache-
 12 lor's degree, master's degree, and doctoral-level classes and, at the end
 13 of each term, asked the students to anonymously answer some ques-
 14 tions about my method of instruction. The general results are depicted
 15 in Table 5.3:

Table 5.3 Percentage of Students (N = 40) Who "Strongly Agreed" or "Agreed" with the Anonymous Survey's Questions.*

79% "I found answering the study questions an excellent way to learn the course content."

88% "Answering the study questions helped me to keep up to date in my readings."

85% "I found answering the study questions a better learning tool than having to prepare for mid-term and final examinations."

94% "I found answering the study questions a better learning tool than having to write a term paper."

65% "I devoted more time studying my class readings in this course than in my other courses."

37% "I attended this class more regularly than my other classes."

*Reproduced from Thyer (1988, p. 51).

16 This looks good, but you can readily see the limitations of this type of
 17 descriptive analysis. The students' favorable endorsements may have been
 18 influenced by their desire to please me, even though I tried to minimize
 19 this by keeping their appraisals anonymous. However, if I had gotten really
 20 bad appraisals, that too would have been really informative, basically

1 telling me to change my method of teaching. As it was, this preliminary
 2 positive endorsement led me to conduct future quasi-experimental stud-
 3 ies on my own teaching using this method of instruction.

4 A more recent example of taking a purely descriptive approach
 5 to analyzing the results of a posttest-only study is reported by DeWalt
 6 et al. (2009), who evaluated a goal-setting intervention in the area of
 7 patient self-management of diabetes. The brief 15-minute structured
 8 intervention was intended to help patients with diabetes establish small,
 9 realistic, but meaningful goals in helping them manage their diabetes,
 10 in areas such as diet, exercise, blood glucose monitoring, medication
 11 adherence, and insulin use. The patient chose an area on which to focus
 12 and was helped to create an action plan intended to lead to healthy
 13 behavior change. Follow-ups occurred some 3–4 months following
 14 the intervention, and patients were asked if they remembered the action
 15 plan and whether they had achieved the behavioral goal. Of an initial
 16 250 patients who received the intervention, 20 did not complete the
 17 study, which is fairly low attrition. One set of purely descriptive results
 18 are presented in Table 5.4 (from Dewalt et al., 2009, p. 221):

Table 5.4 Number of Subjects Who Achieved and Sustained a Given Number of Goals

<i>Number of Times Goals Achieved/ Behavior Sustained</i>	<i>Frequency</i>	<i>Percent</i>
0	17	7
1	44	19
2	92	40
3	76	33
Total	229	

19 Other information was provided in this study, but the core analy-
 20 sis, the extent to which patients reported positive behavioral changes
 21 following the goal-setting intervention, was presented purely descrip-
 22 tively. The researchers were pleased with this result, inasmuch as the
 23 intervention was low-cost, brief, and seemingly resulted in patient
 24 changes. Although this level of analysis may seem rather low-grade ore

1 for scientific investigations, keep in mind that such studies are best seen
2 as preliminary or pilot work that serves as a precursor to more sophisti-
3 cated investigations.

4 We can also compare the outcomes of a posttest-only study with
5 those expected on the basis of chance alone. This was the approach used
6 by Albright and Thyer (2010) in their simple test of the validity of the
7 national examination used in most states to license clinical social work-
8 ers (the LCSW exam). These authors obtained a copy of the LCSW
9 sample or practice test, a test said to be similar in difficulty and content
10 to the real test. This sample LCSW test was completed by 59 first-year
11 MSW students. However, the actual questions were blanked out, leaving
12 only the four possible answers to each question visible, and the students
13 were told to pick the correct response. They knew this was a study on
14 the guessability of the LCSW exam. Now, with four possible options per
15 item, but only one correct one, it could be predicted that the average
16 score would be about 25% correct. In reality our students answered on
17 average 52% of the questions correctly, more than double the score
18 expected by chance alone.

19 Now Albright and Thyer could have simply reported these data
20 descriptively, as above, but this could have left lingering the possibility
21 that perhaps these students scored so much higher on the basis of chance
22 alone; maybe these students were particularly lucky. A simple inferential
23 test can be used in situations like this, a method called the *Z test*. The *Z*
24 test is used when you have a sample of clients who score some given value
25 on a particular measure, and this value can be known or predicted from
26 the larger population of interest. The *Z* test statistic can tell you if the
27 obtained score is significantly different from the predicted score, and it is
28 often used in standardized tests of this nature. In this case, the obtained
29 score is the average of how our 59 students scored, or 52%. The expected
30 population score is 25%: how well “everyone” should score when guess-
31 ing randomly. The *Z* test is calculated using the following formula:

$$32 \qquad Z = (X - m) / SE$$

33 Where *X* is the mean of the sample to be standardized, *m* (μ) is the
34 population mean, and *SE* is the standard error of the mean. $SE =$
35 s / \sqrt{n} , where *s* is the population standard deviation and *n* is the

1 sample size (see <http://changingminds.org/explanations/research/analysis/z-test.htm>). The value of Z tells you how much the sample score differs
2 from the population mean, in terms of standard deviation units. With
3 the Z score in hand, one looks up a Z table in a statistics book or online
4 and determines if the result is statistically significant. If it is, then you
5 know that the results are unlikely (usually at the .05 level) to be due
6 to chance. In other words, a real effect is present, and in the case above,
7 yes, the 59 MSW students scored significantly better than chance.
8 Therefore, the sample test (a proxy for the real exam) is very guessable,
9 and hence may not be a legitimate evaluation of one's ability to practice
10 social work safely. This type of study, using blanked-out questions from
11 sample tests, has been used in a variety of areas to examine the validity
12 of standardized tests. Other examples within social work include the
13 validity of the GRE as an admissions requirement for MSW programs
14 (Donohue & Thyer, 1992), the School Social Work Examination
15 (Johnson, Thyer, Daniels, Anderson, & Bordnick, 1996), the Academy
16 of Certified Social Workers examination (Thyer & Vodde, 1994), and the
17 advanced practice examination also used to license social workers at
18 a lower level than the LCSW (Randall & Thyer, 1994). All these studies
19 used the posttest-only design and the Z test as an inferential statistic.

20
21 Another method of analysis was used by Thyer, Sowers-Hoag, and
22 Love (1986) in their analysis of how BSW and MSW student perceptions
23 of field instruction quality varied according to gender mix. At the end
24 of an internship, all students completed a standardized measure evaluat-
25 ing their field experience. Over the course of a number of semesters,
26 students received their primary supervision from a supervisor of a given
27 gender, and we were interested in seeing if the students' perception of the
28 quality of their field experience varied by their own gender and that of
29 their supervisors. When we did this study, we had post-internship super-
30 vision satisfaction scores for 413 students. This posttest-only design,
31 including the numbers of students per group and their mean satisfaction
32 score (with the standard deviation for each mean), could be diagrammed
33 as follows:

34 Because we had data that was scored using an interval scale, the
35 appropriate inferential test is one called a *one* (time period) *by four*
36 (groups) *analysis of variance* (ANOVA). Basically, we found no meaning-
37 ful differences across the four groups of students and supervisors. Thus,
38 suggestions that had appeared in prior literature indicating that students

# of Students	Type of Supervision They Received	M Score (SD)
N = 217	W Female Student/Female Supervisor – O ₁	68.4 (8.6)
N = 122	X Female Student/Male Supervisor – O ₁	63.9 (11.6)
N = 30	Y Male Student/Female Supervisor – O ₁	63.0 (10.0)
N = 44	Z Male Student/Male Supervisor – O ₁	64.6 (9.4)

1 be placed on the basis of gender with supervisors of the same gender in
 2 order to achieve an optimal internship experience were not supported by
 3 our data.

4 The situation is a bit more complicated in the case of the pretest–
 5 posttest design [O₁–X–O₂], with the choice of test also being dependent
 6 upon the nature of the data being collected. In both descriptive and
 7 inferential statistics, data can be roughly grouped into the following
 8 methods of measurement:

- 9 • **Categorical** (also known as nominal) data, classifies one’s data
 10 into groups according to some identifiable feature. Descriptively,
 11 one might think of variables such as gender (male, female), race
 12 (white, black, Hispanic, Asian, etc.), food stamp status (yes or
 13 no), or religion (Protestant, Catholic, Jewish, Muslim, Hindu,
 14 etc.). Examples drawn from outcome research might include
 15 a binary categorization of *Cured* versus *Not Cured*, *Positive*
 16 versus *Negative* (think of drug screening results), *Pass* versus
 17 *Failed* (think of school performance), etc. For the purposes of data
 18 analysis, numbers may be assigned to categorical data (e.g., for
 19 entry into a database), with say 1 = male and 2 = female,
 20 or 1 = Protestant, 2 = Catholic, 3 = Jewish, 4 = Muslim, and 5 =
 21 Hindu. However, these numbers have no mathematical meaning;
 22 for example, a female (scored as a 2) is somehow not twice the
 23 value of a male (scored as a 1). Nor would it make sense to calculate
 24 the average gender or religion of your clients using the above
 25 coding schemes. Categorical data are usually reported in terms of
 26 numbers and percentages, not means and standard deviations.
- 27 • **Ordinal** data occurs when values are ordered in ranks that
 28 represent some sort of meaningful hierarchy, as in first place,

- 1 second place, and third place; college status (freshman,
2 sophomore, junior, senior); or levels of impairment (highly
3 impaired, moderately impaired, mildly impaired). Knowing
4 the *order* of something provides more information than does
5 simply knowing a category because it also conveys a sense of
6 hierarchy. But the order does not provide information as to
7 the magnitude of differences. A first-place horse can lead the
8 second-place horse by a nose, a length, or by several lengths.
9 Knowing that one was in first place and the other in second
10 tells you about their order, but not the extent of their differences.
11 Ordinal data are usually reported in terms of numbers and
12 percentages, and their preferred measure of central tendency
13 is the *mode* (most common value in a series) or *median*
14 (the midpoint in a range of values).
- 15 • **Interval** data can be used to categorize values as well as place
16 them in a hierarchy, but they convey still more information
17 since the values assigned have a mathematical meaning in
18 relation to each other, with differences representing meaningful
19 and consistent distinctions. Examples include clients' weights,
20 heights, or cholesterol levels, or a student's SAT score or her
21 score on some scale, test, or measure whose range of values does
22 not include a meaningful zero value. One cannot weigh zero
23 pounds, have zero height, or even earn a zero score on the SAT
24 (if you take it you have some sort of non-zero score). The values
25 of something measured on an interval scale have arithmetic
26 meaning in relation to each other. Someone weighing 200 pounds
27 weighs twice as much as someone who weighs 100 pounds.
28 One student's SAT score of 1,000 is 100 points below the score
29 of a student who scored 1,100, and 100 points above the score
30 of a student with a 900 score. The intervals are basically
31 equivalent to each other. The preferred measures of central
32 tendency for interval data are the mode, mean, or *arithmetic*
33 *mean*, with the latter being the most commonly used one in
34 inferential statistics.
 - 35 • **Ratio** data possess all of the attributes of interval data, except that
36 the scaling system possesses a meaningful zero value, representing
37 the complete absence of the attribute. The number of children
38 a client has, the number of crimes committed or hospitalizations

1 experienced, or the amount of money in the bank reflect
2 examples of data that can be analyzed using ratio scales.
3 The value of each could be zero. The central tendency of a
4 variable measured on a ratio scale may include the mode,
5 median, or arithmetic mean (e.g., average). In physics, the Kelvin
6 scale, which includes a zero value reflecting the complete absence
7 of warmth, is an example of a ratio level of measurement.
8 However the Celsius temperature scale is an example of an
9 interval measure, since, although it does include a value of zero,
10 this number was arbitrarily set because it represents the value
11 at which water freezes, not the complete absence of temperature.
12 Cold as it is, water at 0° Celsius retains warmth/temperature and
13 can grow colder still.

14 Now, in returning to the example of the one-group pretest–posttest
15 design, in order to decide what statistic to use, we must determine on
16 what level of measurement to scale the outcome value. Let take as an
17 example psychiatric patients who are admitted for inpatient treatment
18 and, upon admission, are asked to complete a measure of psychiatric
19 symptomatology that yields scores with the interval level of measure-
20 ment. Then, when they are about to be discharged, they complete the
21 same measure. The social work researcher wishes to test the hypothesis
22 “Clients who are treated on our unit will display statistically significantly
23 lower levels of psychiatric symptomatology on discharge, relative to their
24 scores on admission.” This is a good hypothesis. It can be falsified (clients
25 might on average grow worse), and it is directional, calling for changes in
26 one direction only (they will get better, not worse). A directional hypoth-
27 esis is a riskier hypothesis in that it is easier to falsify, compared to a
28 nondirectional one (“I predict that clients will *change* on average over the
29 course of their treatment on this unit.”) A directional hypothesis that
30 also asserts a certain level of change (“Clients will improve, on average,
31 by at least 5 points”) is even riskier than a purely directional one that
32 lacks an additional prediction as to the extent of change.

33 In the case of the one-group pretest–posttest design ($O_1 - X - O_2$),
34 the appropriate test to see if scores have significantly changed between
35 the two assessments is called the *paired sample t-test* (also called the pair-
36 wise *t-test*), if you have an outcome measure that is scaled as an interval
37 or ratio measure. It is called *paired* because the two sets of scores are

1 from the same group of people. This test basically examines the mean
2 (average) score at pretreatment, compares it to the average score at post-
3 treatment, and lets you know, with a certain level of probability, if the
4 observed difference is likely due to chance or to some other nonrandom
5 factor. It does not tell you what *caused* any differences. It is highly unlikely
6 that the pretest and posttest scores will be exactly the same—people do
7 change, and random and systematic errors occur in measurement—so
8 the mean scores will very likely differ. The question for the researcher is,
9 “Is this difference due to chance or not?” If the *t*-test does fall below a
10 certain threshold of probability (the convention is the .05 level, meaning
11 less than 1 chance in 20 that the difference was due to random factors,
12 the data’s natural variability, or chance), the null hypothesis is rejected,
13 and you can assume that *something else* is responsible for the observed
14 changes. This something else may be the treatment, but, as we have seen,
15 it may be due to a number of other factors, such as threats to internal
16 validity, those confounding variables described earlier that may really
17 have caused these changes, not the treatment. So, a statistically signifi-
18 cant *t*-test does not mean that *treatment* caused any improvements
19 (or deterioration for that matter), only that true changes did occur and
20 treatment *may* have been the reason.

21 Back in the day before safety belt use was required by law in most
22 states, I used this design (a pretest–posttest design) to evaluate changes in
23 my students’ seat belt use by offering them extra credit in class if they
24 would sign a safety belt use pledge, agreeing to wear their safety belts each
25 time they drove in a car, in return for some extra credit. However, accord-
26 ing to the contract, if I ever saw them, during the course of the semester,
27 riding in a car and appearing to be unbuckled, they agreed to accept
28 a final grade of F in the class! Thirty-five of forty students signed the
29 pledge at the beginning of the term and provided anonymous estimates
30 of the percent of time they wore their safety belts when driving (about
31 83% [SD = 29%]). At the end of the term, anonymous belt use was
32 reported to be 94% (SD = 17%). The paired sample *t*-test result was
33 [$t(34) = -2.02$; $p < .05$]. (The astute reader will have noted that I should
34 have reported the *exact p* value here.) Because the *t*-statistical was
35 significant at the $<.05$ level, I could conclude that my students’ safety
36 belt use *did increase* over the term, which is consistent with the hypoth-
37 esis that the safety belt pledge had an effect. However, I could not be
38 certain that the increase was the result of the agreement they signed, since

1 other factors could have been responsible. For example, during the term,
2 there might have been a horrible and well-publicized local accident involv-
3 ing college students who died because they were not using their safety
4 belts and it was this publicity that actually caused my students to increase
5 their safety belt use. Or, maybe the state passed a mandatory safety belt
6 use law, and this policy change was really responsible for their greater
7 reported seat belt use. Basically, one has no way of knowing if it was inter-
8 vention (the pledge) that caused the change, but we do know safety belt
9 use increased, and that is a good thing. One student told me at the end of
10 the term that after she signed her pledge she coaxed her boyfriend into
11 wearing his safety belt also. During the term, he was in a serious car acci-
12 dent in which his car was demolished, and both he and the highway patrol
13 officer attributed his survival to his wearing his safety belt. This is a satisfy-
14 ing bonus to my undertaking this small project (Thyer, 1987).

15 A more recent example of this design and the use of the paired-
16 sample t -test is the analysis undertaken by Jones, Chancy, Lowe, and
17 Risler (2010), who looked at the possible effectiveness of residential treat-
18 ment on sexually abusive youth. On intake, all youths (ages 9–18) com-
19 pleted a reliable and valid measure of psychosocial functioning called
20 the Child and Adolescent Assessment of Functioning Scale (CAFAS),
21 and they completed this measure again at discharge (average length of
22 stay was 30 months). A total of 58 youth had pretest and posttest scores
23 available for analysis. Jones et al. (2010, p. 177) posed the following
24 research question: “Do youths’ functional impairment scores and sexual
25 interest scores *change* from intake to discharge from a residential treat-
26 ment program?” (It would have been stronger, scientifically, to not just
27 ask this question, but also to pose the more risky direction prediction
28 included in the hypothesis: “Do CAFAS scores statistically significantly
29 *improve* over the course of treatment?”)

30 Mean CAFAS scores at pretest for the 58 youth were about 145 points
31 ($SD = 39$) and at discharge 74 points ($SD = 46$). With the CAFAS, lower
32 scores imply higher functioning; thus it was found that this change of
33 some 70 points in a positive direction was statistically significant, when
34 examined by the paired-sample t -test. Jones et al. cannot conclude that
35 the treatment program caused these changes. They may be due to matu-
36 ration (teenagers were tested over a 30-month period) or to the passage
37 of time alone. Again, inferential tests can detect reliable differences, not
38 the source of those differences.

1 This same design and statistic was used by Parrish and Rubin (2011)
2 in their analysis of the effectiveness of continuing education (CE) pro-
3 grams they offered on the topic of evidence-based practice. These writers
4 provided a series of CE workshops, and assessments involved having
5 participants complete a reliable and valid measure called the Evidence-
6 based Practice Process Assessment Scale (EBPAS) at the beginning of the
7 workshop and again some 3 months after it was concluded. Delaying
8 assessment for 3 months provided a more valid evaluation of the work-
9 shops' effectiveness since you could look at long-term retention, not
10 knowledge only retained immediately after the conclusion of the training.
11 For all participants, combined across four different workshops, the mean
12 EBPAS score pretraining was about 27 points ($SD = 7$), and 3 months after
13 the 7-hour training program it was about 32 points ($SD = 7$), with higher
14 scores indicating greater knowledge. The t -test [$t(57) = -3.4$; $p < .001$]
15 was significant, demonstrating that these improvements were not due to
16 chance. These authors properly did not make any unwarranted causal
17 inferences; that is, they did not state that they could be sure these
18 improvements were *caused* by the workshop they provided, but they did
19 provide a good discussion of rival explanations and why they believed
20 a case could be made for ascribing the changes to the workshops attended.
21 For example, the pretest–posttest design does not usually control for the
22 threat to internal validity called maturation. However, given that their
23 workshops were only 7 hours in length, the follow-up period was only
24 3 months, and all participants were adults, it is pretty unlikely that matu-
25 ration is a viable rival explanation. Allen Rubin is one of social work's
26 most distinguished researchers, and it speaks to the value of the quasi-
27 experimental pretest–posttest design and paired-sample t -test that this
28 approach was used by him and his colleague to evaluate CE training.
29 Simple designs *do* have value in answering simple questions, and some-
30 times it is very important to answer simple questions first.

31 If an outcome measure is categorical or ordinal in nature, then a
32 different test statistic should be used to analyze the results of a pretest–
33 posttest design. Take the case of a study that compared the outcome of
34 57 patients with panic disorder who received either panic control training
35 (a cognitive behavior therapy), alprazolam (an antianxiety medication),
36 waiting list control (no treatment), or a placebo medication (Klosko,
37 Barlow, Tassinari, & Czerny, 1990). There were several categorical out-
38 come measures assessed posttreatment including end-state functioning

1 (e.g., cured vs. not cured) and the complete absence of experiencing fur-
 2 ther panic attacks. These results are summarized in Table 5.5, along with
 3 the associated chi-square (X^2) analysis. Overall, the study’s results gener-
 4 ally favored the behavioral therapy on a number of outcome measures.

Table 5.5 Selected Categorical Outcome Measures Reported by
 Klosko et al. (1990)

	<i>Cured</i>	<i>Experienced Zero Panics</i>
	N(%)*	N(%)**
Alprazolam (N = 16)	8(50%)	8(50%)
Placebo (N = 11)	5(45.5%)	4(36.4%)
Panic Control Training (N = 15)	11(73.3%)	13(86.7%)
Waiting List (N = 15)	2(20%)	5(33.3%)

* X^2 (3, N = 57) = 8.62, $p < .05$.
 ** X^2 (3, N = 57) = 10.42, $p < .02$.

5 Here, the use of the inferential test helpfully augments simply eyeballing
 6 the data.

7 Jainchill, Hawke, and Messina (2005) used chi-square analyses
 8 to examine possible differences in outcomes among male and female
 9 adjudicated adolescents who received treatment in a therapeutic com-
 10 munity (TC). At 5 years post-TC treatment, the follow-up sample
 11 included 70 males and 51 females who were assessed on an array of psy-
 12 chosocial, criminal, and drug use variables. Very simply put, this study
 13 could be diagrammed as follows:

14
$$N = 70 \text{ males } X - O_1$$

15
$$N = 51 \text{ females } Y - O_1$$

16 Among the statistically significant differences that appeared 5 years after
 17 TC treatment were included (among many variables assessed) whether
 18 the clients was arrested for drug possession, drug sales, or property
 19 crimes. These outcomes are broken down by gender and presented in
 20 Table 5.6.

Table 5.6 Selected 5-year Categorical Outcomes Following Therapeutic Community Treatment for Drug Abuse

	<i>Males</i>	<i>Females</i>
	(N = 70)	(N = 51)
Involved in Drug Possession	63%	34%*
Involved in Drug Sales	56%	18%**
Involved in Property Crimes	29%	10%***

* $\chi^2 = 6.97, p < .01$ ** $\chi^2 = 14.45, p < .001$ *** $\chi^2 = 4.23, p < .04$

From Jainchill, Hawke, & Messina, 2005, p. 984.

1 For each variable, it appears that males are more likely have engaged
 2 in selected illegal activities 5 years following TC treatment, compared to
 3 female clients.

4 Suggested inferential tests suitable for each type of quasi-experimental
 5 design are noted in Table 5.7. There are other appropriate ones that can
 6 be used, and those mentioned are presented as one suggested course of
 7 analysis, not the sole appropriate or definitive approach to statistical
 8 inference with these designs. They are, however, those most commonly
 9 employed.

10 EFFECT SIZES

11 It is now widely recognized that in studies with a sufficiently large sample
 12 size, very small differences can be shown to be statistically significant,
 13 which only means that the difference is reliable or not likely (within
 14 a certain probability) due to random variation in the data. It does not
 15 refer to the clinical importance of the changes or differences observed.
 16 An effect size (ES) should accompany any report of a statistically signifi-
 17 cant difference as it provides more of an estimate of the meaningfulness
 18 of any difference or change. One measure of ES that may be familiar
 19 to the reader is associated with reporting Pearson correlations of paired
 20 quantitative data, correlations that can range from -1 to $+1$. The correla-
 21 tion, r , when squared, yields a measure called the *coefficient of deter-*
 22 *mination*, and it estimates the proportion of variance shared by the
 23 two measures. If two measures are correlated $+0.40$, the coefficient of

Table 5.7 Chart of Designs/Diagrams Outlining Design Categories, Objectives, Statistics Used

Pre-Experimental Designs

1. The Posttest-Only Single Group Design

$$X - O_1$$

This design controls for virtually no threats to internal validity. Its data are usually presented descriptively. Inferential statistics (e.g., Z-test) may be applied if there are known values for the outcome measure available on a larger population of interest.

2. The Pretest-Posttest Single Group Design

$$O_1 - X - O_2$$

This design controls for very few threats to internal validity. Its data can be discussed descriptively and analyzed using inferential tests, such as the paired-sample *t*-test if the outcome measure is scaled as an internal or ratio variable, or the X^2 test if the data are categorical or ordinal.

3. The Pretest-Posttest Single-Group Design with Repeated Pretests

$$O_1 - O_2 - X - O_3$$

This design may partially control for regression to the mean. Interval/ratio outcomes can be evaluated using the analysis of variance for repeated measures (ANOVA), and categorical/interval scaled data can be analyzed using a one (groups) by three (time periods) X^2 test.

4. The Pretest-Posttest Single-Group Design with Repeated Posttests

$$O_1 - X - O_2 - O_3$$

This design may partially control for relapse or improvements that are temporary.

Interval/ratio outcomes can be evaluated using the ANOVA, and categorical/interval scaled data analyzed using the X^2 test.

Quasi-Experimental Designs, with a Control or Comparison Condition

1. The Posttest-Only No-Treatment Control Group Design

$$X - O_1$$

$$O_1$$

(Continued)

Table 5.7 (Continued)

This design may partially control for the passage of time, concurrent history, and maturation. Interval/ratio scaled outcome measures may be evaluated using the t -test for independent samples, comparing posttreatment group means, and the X^2 test applied to the frequency data of categorical or ordinal variables. This design can also involve more than one control or comparison group.

2. The Pretest–Posttest No-Treatment Control Group Design

$$O_1 - X - O_2$$

$$O_1 \quad O_2$$

This design may partially control for the passage of time, regression to the mean, concurrent history, the existence of pretreatment differences between groups, and maturation. Interval/ratio scaled measures may be analyzed using the two (groups) by two (times) ANOVA, and categorical/ordinal scale data using a two (groups) by two (times) X^2 test.

3. The Pretest–Posttest Alternative Treatment Comparison Group Design

$$O_1 - X - O_2$$

$$O_1 - Y - O_2$$

Where X indicates a group that received an experimental intervention, and Y indicates a group that received some alternative treatment, treatment as usual, or a placebo intervention. This design may partially control for placebo effects, social desirability factors (wanting to please the therapist), concurrent history, and existence of pretreatment differences between the two groups. Use a two (groups) by two (times) repeated measures ANOVA for interval/ratio data, or a two by two X^2 test for categorical/ordinal data.

4. The Pretest–Posttest Alternative Treatment/No Treatment Control Comparison Design

$$O_1 - X - O_2$$

$$O_1 - Y - O_2$$

$$O_1 \quad O_2$$

This design may partially control for the passage of time, concurrent history, social desirability factors, regression to the mean, and pretreatment differences among the groups. Interval/ratio level data may be analyzed using a two (times) by three (groups or condition) ANOVA, with a similar X^2 test applied to categorical/ordinal data.

(Continued)

Table 5.7 (Continued)

Note: Each of the above quasi-experimental designs may be modified by using more than one pretest assessment period, more than one posttest assessment, or both, typically strengthening the basic design’s internal validity.

Time Series Designs

1. The Posttreatment-Only Time Series Design

$$X - O_1 - O_2 - O_3 - O_k$$

This design takes a very large number of posttreatment assessments after an intervention has been introduced. It controls for very few threats to internal validity, and its results are usually graphed and interpreted visually.

2. The Simple Interrupted Time Series Design

$$O_1 - O_2 - O_k - X - O_{k+1} - O_{k+2} - O_{k+n}$$

These designs typically have a large number of pretests and posttests. O_k indicates the final pretest assessment, O_{k+1} indicates the first posttreatment assessment, and O_{k+n} the last posttreatment assessment. Time series designs with very large numbers of data points may be analyzed with a test statistic known as time series analysis, for which a minimum of 50 data points per phase is recommended, pre- and postintervention. This design may control for regression to the mean, maturation, and repeated testing.

3. The No-Treatment Control Group Interrupted Time Series Design

$$O_1 - O_2 - O_k - X - O_{k+1} - O_{k+2} - O_{k+n}$$

$$O_1 - O_2 - O_k \quad O_{k+1} - O_{k+2} - O_{k+n}$$

Two similar groups (states, counties, organizations) are repeatedly assessed on some variable of interest. The top group receives an intervention (X) after a number of assessments, and the bottom group does not. This design controls for regression, maturation, repeated testing, concurrent history, and the passage of time. Use time series analysis to investigate changes within this design and for other forms of interrupted time series data.

- 1 determination is $.40 \times .40$, or $.16$. This means that up to 16% of the vari-
- 2 ance in one measure can be predicted from the other. If you have two
- 3 variables with a Pearson correlation of $.80$, then from $.80 \times .80$, we know
- 4 that at most $.64\%$ of the variance in one measure can be predicted from
- 5 the other. In the social sciences, ES around $.10$ – $.20$ are called small, those
- 6 $>.20$ – $.35$ as medium, and those $>.35$ as large.

1 Generally speaking, measures of ES when looking at differences (not
2 correlations) provide an estimate of the extent to which the average
3 member of the experimental treatment group is better (or worse) off
4 compared to the average member of the comparison (TAU, no-treatment,
5 placebo control condition) group, as expressed in standard deviation
6 units. An ES of .30 favoring RT over a no-treatment control group would
7 mean that the average RT client was .30 standard deviation (SD) units
8 better off than persons who did not receive treatment; an ES size of
9 .80 would mean that the treated clients were, on average, .80 SDs better
10 off than those not treated, etc. Typically ES in social work intervention
11 research is rather small, reflecting that our interventions are not excep-
12 tionally potent; however, given the intractability of many of our clients'
13 problems, reliable small gains may be important to demonstrate. Simply
14 knowing an ES by itself not a sufficient measure of importance. It must
15 be interpreted within the context of the overall study design. For exam-
16 ple, an ES of .30 favoring RT clients versus untreated clients is not as
17 impressive as is effect size of .30 favoring treated RT clients versus place-
18 bo-treated or clients who received TAU. For many problems, most thera-
19 pies are capable of yielding some small benefits, comparing to getting
20 nothing. The more robust comparison is to compare experimental treat-
21 ment versus TAU or placebo. Also, outcome measures that are labile, or
22 of weak reliability and validity, may lend themselves to yielding stronger
23 ES than do those obtained from studies using more rigorous measures.

24 Effect sizes may be calculated for all inferential statistical tests, with
25 one known as Cohen's d being the most commonly used for t -tests. It is
26 calculated by taking the difference between the two means (pre- versus
27 post, or between two groups posttreatment), divided by the pooled stan-
28 dard deviation of the data. A similar effect size for use with t -tests is called
29 Hedge's g . For ANOVAs and multiple regression results, another test may
30 also be used, called Cohen's f^2 , whereas for X^2 tests, a measure called
31 Cramer's phi is appropriate. There is a large literature on the importance
32 of calculating ES and of including this information in statistical report-
33 ing. The *Publication Guidelines* (APA, 2009) of the American Psychologi-
34 cal Association now requires it, as do an increasing number of journals
35 (e.g., *Research on Social Work Practice*). A number of articles can be found
36 in the social work (Hudson, Thyer, & Stocks, 1985; LeCroy & Krysik,
37 2007) and related literatures (Cohen, 1994) that address the topic, and
38 common statistical software programs include options for reporting ES.

1 Effect size information is also very important when attempting to
2 systematically review a number of studies evaluating the effects of a given
3 treatment. If sufficient primary statistical information is included in an
4 original study, other researchers will be more able to aggregate the results
5 of relevant studies using a technique called *meta-analysis* to arrive at con-
6 clusions made possible through combining a larger number of small
7 studies. Littell, Corcoran, and Pillai (2008) provide a good review of
8 designing and conducting systematic reviews and meta-analyses, with
9 the latter being based on ES calculations.

10 ETHICAL CONSIDERATIONS IN THE DESIGN AND CONDUCT 11 OF QUASI-EXPERIMENTS

12 As members of a profession, social workers, including researchers,
13 are guided by various codes of ethics. The code of ethics promoted by
14 the NASW (2008) is among the more widely recognized, but other
15 social work organizations (e.g., the Clinical Social Work Association),
16 other interdisciplinary groups that individual social workers may
17 choose to affiliate with, in addition to or in lieu of the NASW (e.g., the
18 Association for Behavior Analysis – International, or the American
19 Evaluation Association), and various state licensing boards may also
20 promulgate specific codes of ethics. So, although the NASW code of
21 ethics is not the sole appropriate standard that may cover a social work-
22 er's activities, it will be referred to here due to its widespread acceptance.
23 The NASW COE standards relating to evaluation and research appear in
24 Table 5.8.

25 It is clear that conducting research, especially evaluation research
26 related to practice and policy, is an expected role of professional social
27 workers. Whenever possible, clients should provide informed consent,
28 without penalty or deprivation, prior to their being enrolled in a research
29 project, and underage or otherwise impaired individuals should
30 have their consent provided by an appropriate proxy (e.g., parent, guard-
31 ian *ad litem*, etc.). Clients must be able to withdraw their participation at
32 any time during a project, and must be protected from undue risk of
33 harm. Information gathered from clients must be protected and treated
34 respectfully and confidentially, not promiscuously disclosed to others
35 unconnected with the research project. Data must be reported honestly.

Table 5.8 National Association of Social Workers Code of Ethics Standards Pertaining to Evaluation and Research Activities

- (a) “Social workers should monitor and evaluate policies, the implementation of programs, and practice interventions.
- (b) Social workers should promote and facilitate evaluation and research to contribute to the development of knowledge.
- (c) Social workers should critically examine and keep current with emerging knowledge relevant to social work and fully use evaluation and research evidence in their professional practice.
- (d) Social workers engaged in evaluation or research should carefully consider possible consequences and should follow guidelines developed for the protection of evaluation and research participants. Appropriate institutional review boards should be consulted.
- (e) Social workers engaged in evaluation or research should obtain voluntary and written informed consent from participants, when appropriate, without any implied or actual deprivation or penalty for refusal to participate; without undue inducement to participate; and with due regard for participants’ well-being, privacy, and dignity. Informed consent should include information about the nature, extent, and duration of the participation requested and disclosure of the risks and benefits of participation in the research.
- (f) When evaluation or research participants are incapable of giving informed consent, social workers should provide an appropriate explanation to the participants, obtain the participants’ assent to the extent they are able, and obtain written consent from an appropriate proxy.
- (g) Social workers should never design or conduct evaluation or research that does not use consent procedures, such as certain forms of naturalistic observation and archival research, unless rigorous and responsible review of the research has found it to be justified because of its prospective scientific, educational, or applied value and unless equally effective alternative procedures that do not involve waiver of consent are not feasible.
- (h) Social workers should inform participants of their right to withdraw from evaluation and research at any time without penalty.
- (i) Social workers should take appropriate steps to ensure that participants in evaluation and research have access to appropriate supportive services.
- (j) Social workers engaged in evaluation or research should protect participants from unwarranted physical or mental distress, harm, danger, or deprivation.
- (k) Social workers engaged in the evaluation of services should discuss collected information only for professional purposes and only with people professionally concerned with this information.
- (l) Social workers engaged in evaluation or research should ensure the anonymity or confidentiality of participants and of the data obtained from them. Social workers should inform participants of any limits of confidentiality, the measures that will be taken to ensure confidentiality, and when any records containing research data will be destroyed.

(Continued)

Table 5.8 (Continued)

- (m) Social workers who report evaluation and research results should protect participants' confidentiality by omitting identifying information unless proper consent has been obtained authorizing disclosure.
- (n) Social workers should report evaluation and research findings accurately. They should not fabricate or falsify results and should take steps to correct any errors later found in published data using standard publication methods.
- (o) Social workers engaged in evaluation or research should be alert to and avoid conflicts of interest and dual relationships with participants, should inform participants when a real or potential conflict of interest arises, and should take steps to resolve the issue in a manner that makes participants' interests primary.
- (p) Social workers should educate themselves, their students, and their colleagues about responsible research practices."

Reprinted from the National Association of Social Workers (1999, Section 5.02).

1 These are all sensible standards, and there is seemingly little to quibble
2 about. However, the devil is in the details.

3 When data are gathered retrospectively and anonymously, perhaps
4 obtained from state or federal agencies, the principle of informed con-
5 sent for participation in research is largely a moot point. In many
6 instances, the data collected and analyzed in time series designs are not
7 even derivable down to the level of individuals. It is not people who are
8 directly being measured but more conceptual phenomena such as acci-
9 dents, visits, numbers of births, high school dropouts, etc. It may even be
10 possible that formal institutional review board (IRB) approvals are not
11 necessary for such studies, if you are not interacting with human beings
12 or gathering personally identifiable information. This may also be the
13 case if you are making use of data that are publicly available, as such
14 studies too may be exempt from IRB oversight. However, if a researcher
15 is employed at an institution that receives federal funding, it always a
16 wise policy to check with the chair of the local IRB regarding the possibly
17 exempt status of your project. In the case of IRB oversight, the best policy
18 is to get permission first, not seek forgiveness afterward. IRBs can be very
19 touchy on this issue (see Holosko, Thyer, & Danner, 2009).

20 When it might be technically possible to obtain informed consent
21 from research participants but is impractical, the IRB may grant you
22 a dispensation from tracking down individuals whose data comprises

1 the variables you are analyzing. This can be especially useful when the
2 data may be years old and locating individuals would be very difficult.
3 This dispensation is more likely to be granted if the data are innocuous,
4 not sensitive; if your sample of clients is not a protected group (e.g., pris-
5 oners, pregnant women, minorities of color, children), if the risks are
6 otherwise low, and the identity of respondents is either not known or will
7 be kept confidential.

8 Generally speaking, according to federal policy, conducting a quasi-
9 experimental evaluation of practice or policies will be considered research,
10 if the project meets *both* of the following standards:

- 11 1. The project involves a systematic investigation, *and*
- 12 2. The design, goal, purpose, or intent of the project is to contribute
13 to generalizable knowledge.

14 Generalizable knowledge is interpreted by the federal government
15 to mean that the researcher plans to publish his or her results in a journal
16 or present it at a professional or academic meeting or conference. This is
17 an important caveat: If your purpose in conducting an evaluation is
18 solely to develop an internal report for an agency or organization, per-
19 haps with the intent of improving the agency's services, with no plans to
20 distribute the report publicly, then the project does not rise to the thresh-
21 old of the federal definition of research, and no external oversight or
22 approval from a Human Subjects Protection board is required. You can
23 do all the quasi-experimental evaluations you wish, so long as you have
24 no intent (and do not eventually) to publish or present them publicly.
25 Technically, such projects are not research!

26 A *human subject* is a living person from whom an investigator/
27 researcher obtains data (1) through intervention or other interaction
28 with them or (2) through identifiable personal information (e.g., name,
29 address, social security number, etc.). *Intervention* refers to the physical
30 methods by which data are collected, and any type of manipulation of the
31 client or his or her environment for the purposes of the research.
32 *Interactions* includes communication or interpersonal contact between
33 the researchers and the clients. What these features include is fairly clear,
34 but what they exclude is often overlooked. For example, data related to
35 deceased individuals may not be technically construed as research, but
36 could be the source of information comprising a quasi-experimental study.

1 Data gathered by others, say agency staff, and provided to the researcher
2 with the data de-identified (no personal information provided) could
3 conceivably not be called doing research. Analyzing publicly available
4 data, from state databases for example, lacking personal information
5 (say, mortality statistics before and after some new policy is enacted)
6 would not be engaging in research, according to federal guidelines. There
7 is no personal information, and you are not interacting with human
8 beings. Even if personal information is gathered, it is conceivable that
9 the activity may not be construed as research in certain circumstances.
10 For example, my local paper publishes color mug shots of people arrested
11 in our community on a weekly basis. I could use these photos (which
12 are *very* personal) of real live people in some sort of research project,
13 for example to examine if white or black felons received disparate sen-
14 tences for the same type of crime, or to see if males with beards or facial
15 tattoos committed different types of crimes than do clean-shaven men.
16 Earlier, I described how I accessed individual teacher's course evalua-
17 tions from my university's publicly available website and looked for
18 possible differences in teaching effectiveness. In that case, although
19 I technically did not need to gain approval from my university's IRB,
20 I chose to do so, just in case any irate person questioned my using this
21 perhaps sensitive information.

22 Even if one does not belong to a professional association that pro-
23 motes a particular code of ethics or is a licensed social worker, a part
24 of being a professional social worker consists of adhering to ethical
25 practices. This includes persons who conduct research such as quasi-
26 experiments. The principle of "First, do not harm" is of paramount
27 importance, as are the general values of respecting clients, protecting
28 privacy, obtaining informed consent when appropriate (with provision
29 for withdrawing such consent without any penalty), and beneficence
30 (hopefully, clients or the field will benefit in some manner, from your
31 study). Quasi-experimental outcome studies, because they involve inter-
32 ventions provided to real live clients, must be based on an ethical bar set
33 higher than more benign and less intrusive forms of research (e.g., sur-
34 veys, correlational investigations). One of the most infamous studies ever
35 conducted in the United States was the Tuskegee Study of the natural
36 history of untreated syphilis involving low-income African American
37 men in the rural south. This could be construed as a quasi-experimental
38 study, a posttest-only time series analysis lasting decades. It was a social

1 worker, Peter Buxtun, who served as the whistle-blower on this project
 2 and forced its eventual termination and the provision of treatment (and
 3 financial compensation) to the participants and their families. (see [http://](http://en.wikipedia.org/wiki/Peter_Buxtun)
 4 en.wikipedia.org/wiki/Peter_Buxtun). Even relatively unsophisticated
 5 quasi-experimental designs can be the context for unethical studies. It is
 6 crucial that social workers undertake quasi-experiments in a manner
 7 consistent with the highest ethical standards. This requires familiarity
 8 with appropriate codes of ethics and due consideration of these ethical
 9 standards from the inception of any such study.

10 THE FUTURE OF SOCIAL WORK AND QUASI-EXPERIMENTAL RESEARCH

11 Across the social sciences, we see that a far greater proportion of quasi-
 12 experiments are published relative to true RCTs, and that this distribution
 13 holds true within the social work disciplinary literature devoted to empir-
 14 ically evaluating the outcomes of practice. Our field has made use of such
 15 designs since the early part of the 20th century, and they form an impor-
 16 tant core body of research investigating what has worked, and not worked,
 17 in serving clients effectively. Several predictions may be ventured:

- 18 1. Social work will continue to make comparatively extensive use
 19 of quasi-experimental designs, although they, like empirical
 20 research as a whole, will remain a minor form of disciplinary
 21 scholarship.
- 22 2. Recognizing that the weaker of these designs pose significant
 23 limitations in terms of permitting true causal inferences, those
 24 that are published will include more stringent cautionary
 25 language so that readers will avoid exaggerating the results or
 26 extend causal claims beyond those legitimately permitted by
 27 the data. Rubin and Parrish (2007) documented the degree to
 28 which problematic phrases involving causal inferences are found
 29 in published experimental and quasi-experimental social work
 30 outcome studies. Rubin himself honestly noted his own mistakes
 31 in this area, citing a specific example, and vowed to be more
 32 conservative in his future writing. With a noble example like that,
 33 we can hope that journal editors as well as authors will be more
 34 vigilant in excising unjustifiable claims (“Treatment X *caused*

1 clients to get better”) from articles before they appear in
2 print. More conservative language might say something like
3 “The results of this study are *consistent* with the hypothesis
4 that treatment X caused the clients to get better.”

5 3. Large-scale funded quasi-experiments may come under increasing
6 pressure to be registered. For many years, RCTs have been
7 encouraged, and in some cases required by funders, to be
8 registered with a system sponsored by the World Health
9 Organization (see www.clinicaltrials.gov). Listing on this clinical
10 trials registry encourages transparency in design and reporting,
11 aids in recruiting research participants, and serves as a large-scale
12 database of experimental intervention research. Over 80,000
13 studies from over 150 countries are on the ClinicalTrials.gov
14 registry, most of which are drug trials for various medical
15 conditions. However, thousands of trials are listed in the general
16 area of mental health, and over 400 appear when “psychotherapy”
17 is used as a search keyword. This is a great resource for
18 psychosocial intervention researchers, as well as pharmaceutical
19 investigators. It has been recently suggested that a similar separate
20 registry be specifically developed for quasi-experimental studies
21 (Staff, 2010). However, it should be noted that almost 14,000
22 quasi-experimental studies are already listed on www.clinicaltrials.gov.
23 Whether a new registry for quasi-experiments is developed or
24 not, it seems likely that the public registration of quasi-
25 experimental outcome studies of psychosocial interventions on
26 such trial registries will become increasingly common in the years
27 to come, especially if governmental funding sources require this as
28 a precondition of receiving research dollars.

29 SUMMARY

30 Pre-experimental and quasi-experimental research designs are major
31 investigatory tools in the evaluation of the outcomes of social work prac-
32 tice and in beginning enquiries into the causal effects of specific psycho-
33 social interventions. These designs provide excellent ways to get credible
34 answers to some simple questions facing each practitioner, administra-
35 tor, and policy analyst. The designs reviewed in this chapter range from

1 the simple and parsimonious to the complex and elegant. Each possesses
2 strengths and limitations and should not be accepted for use without a
3 full consideration of a given design's potential to provide the answers to
4 the questions being posed by the social worker.

5 Some of these designs were used in the very earliest published evalu-
6 ations of the effects of social work, and they remain in widespread use
7 today. In this volume, I have tried to portray the essential features of
8 the major varieties of quasi-experiments, the inferential logic behind
9 them, and to present an array of examples of their use. I have highlighted
10 many such studies authored by social workers that appeared in some of
11 the highest-quality scientific journals, to illustrate the utility and accep-
12 tance of these designs within the scientific toolbox we all have access to.
13 All social workers, at a minimum, should be able to recognize these
14 designs, understand their logical foundations, and provide an informed
15 critique of contemporary research that makes use of them. Some social
16 workers will find it within the scope of their professional activities to
17 actually undertake using some of these designs in the evaluation of their
18 own practice. This book was written to encourage such efforts, as such
19 studies will promote the empirical foundations of our discipline, as envi-
20 sioned over a century ago by many of our founders.

Glossary



Alpha level Alpha is the probability of making a type I error (rejecting the null hypothesis when the null hypothesis is true) when using an inferential statistical test. Most inferential tests set alpha at or less than .05.

ANOVA A parametric inferential statistic that examines differences between the means of three or more groups in a study, groups exposed to different independent variables (e.g., treatment vs. no treatment), or longitudinally at least three times for a single group (e.g., pretest, posttest, and at follow-up).

Assessment/Treatment interaction Changes in a study's outcome measures induced by an interaction between the assessment procedures used and interventions received. This may be a threat to the internal validity of a study.

Attrition/mortality Clients sometimes drop out of a research study before it is completed, and this drop out is known as attrition or mortality. This may be a threat to the internal validity of a study since the participants remaining at the study's conclusions are only a subset of the individuals who began the study.

Beneficence A primary ethical concern of social research. It refers to both doing no harm to people you are studying and, at the same time, promoting a common good for individuals in the research community because of your study. Its origin in present-day social research in America can be traced back to the Belmont Report.*

* *Terms with an asterisk were reproduced with permission from M. J. Holosko and B. A. Thyer (2011). Pocket glossary for commonly used research terms. Thousand Oaks, CA: Sage.*

- Categorical data** Data (variables) that differ only in kind, not in amount or degree. Nominal data are categorical: for example, female versus male, true versus false.*
- Causal inference** Drawing conclusions about the effects of an independent variable by ruling out rival explanations apart from the intervention under investigation.
- Chi-square test** A nonparametric test of statistical significance, appropriate when the data are in the form of frequency counts. It compares frequencies actually observed with expected frequencies to see if they are statistically different.*
- Cohen's d** A widely used measure of effect size.
- Cohort study** An observational (e.g., quasi-experimental) study in which a defined group of people (the cohort) is followed over time. The outcomes of people in subsets of this cohort are compared, to examine those who were exposed or not exposed (or exposed at different levels) to a particular intervention or other factors of interest.*
- Comparison group** A group of clients in a study who receive treatment as usual, partial treatment, or a placebo treatment. Changes observed in a comparison group can be “subtracted” from changes observed in the group of clients who received “real” treatment, to help determine the effects of “real” treatment absent changes induced by receiving the usual treatment, placebo treatment, or partial treatment.
- Control group** A group of clients in a study who do not receive any formal intervention. This is used to control for various threats to internal validity such as the passage of time, concurrent history, and regression to the mean. Changes observed in the control group can be “subtracted” from changes observed in the treatment group, to help determine the “real” effects of the treatment, absent changes induced by non-treatment-related factors.
- Demographic features** Background information relating to statistical characteristics of a study's groups (e.g., age, gender, race, income, etc.).
- Dependent variables** What is measured in a study, and what is affected during the study. The dependent variable (e.g., outcome measure) responds to the independent variable (e.g., treatment or intervention). It is called dependent because it depends on the independent variable.*
- Descriptive statistics** Numbers used to describe the basic features of sample data in a study. They provide simple summaries about the sample and its measures; for example, mean, median, mode, variance, or standard deviation. Descriptive statistics are given at the beginning of most quantitative studies' data analysis processes.* Sample descriptive statistics should include information such as a group's gender distribution, age, race, ethnicity, socioeconomic status, and diagnosis (if relevant).

- Differential attrition** Occurs when clients drop out from participation in a study to varying degrees across groups (e.g., treatment versus no-treatment). Thus, at the end of the study, the proportions of clients remaining in each group may differ, making it difficult to draw any inferences about the effects of treatment. This may be a threat to the internal validity of a study.
- Diffusion/contamination of treatments** This occurs when a comparison group learns about a research program from other program participants, thus preventing the control and experimental groups from remaining distinct. This may be a threat to the internal validity of a study.*
- Differing treatment credibility** This may occur when clients receiving differing treatments (e.g., real treatment vs. placebo therapy) perceive that the treatment they are receiving may be more or less believable (e.g., effective). More credible treatments may exert more powerful placebo effects than less believable interventions. This may be a threat to the internal validity of a study.
- Dismantling study** An outcome study in which one group of clients receives a “complete” treatment package, and their results are compared to clients who receive a subset of the complete treatment. Differential outcomes may be ascribed to the absent components of complete treatment.
- Effect sizes** An index used to indicate the magnitude of an obtained change, result, or relationship between time 1 and time 2 observations. Cohen’s *d* is the most commonly used statistic to compare mean score differences. Approximately speaking, effect sizes (ES) can be small, $<.03$; medium, $+.05$; or large, >0.75 .* Each time a statistically significant difference is reported, its associated ES should be included.
- Effectiveness study** A study evaluating a treatment, conducted under clinically representative or real-life conditions. It is usually used in the later stages of evaluating a new intervention.*
- Efficacy study** A study conducted under conditions of maximum experimental control (e.g., carefully screened clients, highly trained therapists using detailed treatment manuals).* Such studies maximize potential internal validity at the expense of external validity (e.g., generalizability to real-life practice). Interventions found useful using efficacy studies should be replicated in effectiveness studies in real world settings.
- Efficacy subset analysis** Examining various subgroups of participants who received a given intervention to ascertain any possible differences in outcomes (e.g., Do females respond more or less than males?).
- Evaluation study** A systematic inquiry to describe or assess the intervention impact of a specific program or intervention on individuals by determining its activities and outcomes. These can be evaluations of practice or programs.*
- Experimental design** A research study in which one or more independent variables are systematically varied by the researcher to determine their effects on

- dependent variables.* Randomized experiments randomly assign participants to various treatment, control, or comparison groups.
- F test** A statistical test of the equality of the variances of two or more populations. The test compares the differences between groups and within groups over time. It is used in analysis of variance inferential tests (ANOVA).*
- Human subject** A living human being from whom data involving personally identifiable information is obtained in the context of a research study. The term *participant* is now preferred over the word *subject*, and outcome studies of social work often use the more accurate term *client*.
- Hypothesis** A tentative, testable assertion regarding the occurrence of certain behaviors or events; a prediction of study outcomes. It is used to determine how independent and dependent variables can be tested and written in either null or directional form. It is based on literature, theory, or observation of a phenomena. It forms the basis of experiments designed to establish plausibility, association, prediction, or causality.*
- Independent variable** The variable that affects or is presumed to affect the dependent variable under study and is included in the research design so that its effect can be determined. This is sometimes called the *experimental*, *manipulated*, or *treatment variable*, or in outcome studies, the *treatment* or *intervention*.*
- Inferential logic** The process of drawing conclusions from a research study using the principles of logic, specifically those pertaining to inductive and deductive reasoning.
- Institutional Review Board (IRB)** A committee designated to approve, monitor, and review biomedical and behavioral science involving humans, with the aim of protecting the rights and welfare of the research participants. It is a federal requirement in universities, large organizations, hospitals, and so on.*
- Intention-to-treat analysis** A method of evaluation for randomized trials in which all participants randomly assigned to one of the treatments are analyzed together, regardless of whether they completed or received that treatment, in order to preserve randomization.* This approach may also be used in with non-randomized quasi-experimental studies.
- Interrupted time series design** Longitudinal research in which ongoing repeated measurements of the outcomes are made and treatment is introduced at some point, while measurements continue as before.*
- Interval data** Variables scaled using a system that produces rank ordering and equal distances. Interval data lack an absolute zero point.
- Instrument change** Occurs when changes in obtained measures are due to the instrument calibration or changes in observers, judges, or interviewers (e.g., greater sensitivity with practice, or less observer attentiveness after repeated observations). This may be a threat to the internal validity of the findings of a study.*

JARS Journal Article Reporting Standards, guidelines contained within the *Publication Manual of the American Psychological Association* pertaining to the writing up and analysis of research reports. Many journals now require submitted manuscripts to be in compliance with the JARS.

Maturation The possibility that results are due to changes that occur in participants as a direct result of the passage of time, human developmental processes, or fatigue, and that may effect their performance on the dependent variable. This may be a threat to the internal validity of the findings of a study.*

Meta-analysis A systematic review that uses quantitative methods of published research interventions and studies to synthesize and summarize the results of a large number of research studies on one particular topic. This allows aggregate claims about interventions and their effects to be made and offers empirical suggestions about best practices or interventions. The unit of analysis in meta-analysis is the effect size found in different studies.*

Multiple posttreatment assessments The process of repeatedly and formally assessing a study's outcome measures at several (not just one) points in time following exposure to an intervention. These can be used to help establish the long-term effects of any treatment.

Multiple pretreatment assessments The process of repeatedly and formally assessing a study's outcome measures at several (not just one) points in time before clients receive treatment. These can be used to help establish any trends in the data (e.g., are the clients getting better, worse, or is functioning stable) prior to treatment.

Multiple treatment interference The carryover or delayed effects of prior experimental treatments when individuals receive two or more experimental treatments in succession. This may be a threat to the internal validity of the findings of a study.*

n The number of people in a sample.

N The number of people in a population.

Nonparametric tests A body of statistical tests used when the data represent a nominal (categorical) or ordinal level scale, or when the assumptions required for parametric tests cannot be met. This class of tests do not hold the assumptions of normality.*

Objectivity A presumed lack of bias or prejudice.*

One-group posttest-only design A pre-experimental design involving one group that is given a test after treatment is given. It attempts, therefore, to evaluate a program's outcomes when no available comparison group and no pretest data are available (or needed, as in a client satisfaction study).*

One group pretest–posttest design A pre-experimental design involving one group that is pretested, exposed to a form of treatment, and then posttested.*

Ordinal data Assigning numbers to variables, presenting the rank ordering (first, second, third, etc.) of the entities measured. This is the second level of measurement, one up from the first, nominal, or categorical.*

- Outcome measures** Specific standardized or nonstandardized benchmarks used to assess whether the intervention or program resulted in any changes.* Also known as *dependent variables*.
- Outcomes research** Research to measure practice or program effectiveness. Such studies examine what has changed as a result of the intervention being offered.*
- Parametric tests** Inferential statistical tests based upon the assumption that the data are normally distributed.
- Passage of time** Changes in client functioning that may occur during the natural course of events, unrelated to treatment. Because many client problems/issues naturally wax and wane over time, time itself may be a threat to a study's internal validity.
- Pearson correlation** A common index of correlation appropriate when the data represent either interval level or ratio scales. It takes into account each and every pair of scores and produces a coefficient (r) between 0.00 and plus or minus 1.00. A positive r indicates that, as one variable goes up or down, so does the other. Negative or inverse r indicates that as one variable goes up, the other goes down.*
- Placebo influences** An inactive treatment or procedure, literally meaning "I do nothing." The placebo effect (usually a positive or beneficial response) is attributable to the participant's or experimenter's expectation that the treatment will have an effect.*
- Posttest-only control/comparison group design** A research design involving at least two groups of participants. One group receives a treatment, the other receives no treatment, placebo treatment, or an alternative treatment condition. This design is quasi-experimental if the groups are formed naturally. It is an experiment if the groups are formed using random assignment.*
- Pre-experimental research design** A research design that involves studying only a single group of participants, either posttreatment only, or pre- and posttreatment. No control or comparison groups are used.*
- Pretest–posttest no-treatment control group design** A study wherein one group of clients is assessed, receives an intervention, and is reassessed. Their results are compared to a comparable group of clients who were assessed, not treated, and then reassessed. It is an attempt to control various threats to internal validity, such as passage of time, concurrent history, and regression to the mean.
- Pretest–posttest alternative treatment comparison design** A study wherein one group of clients is assessed, receives an intervention, and is reassessed. Their results are compared to a comparable group of clients who were assessed, who then receive an alternative treatment (e.g., treatment as usual, placebo care), and are then reassessed. It is an attempt to control various threats to internal validity, such as passage of time, concurrent history, regression to the mean, placebo influences, and expectancy bias.

Quasi-experimental design A type of research design in which the treatment and control or comparison groups are not created using random assignment procedures. It does involve the manipulation of an independent variable and the specification of a test hypothesis.*

Random assignment A method analogous to tossing a coin to assign clients to treatment groups. The experimental treatment is assigned if the coin lands on heads, and a conventional, control, or placebo treatment is given if the coin lands on tails.*

Random selection A sample selected in such a way that every member of the population has an equal chance of being selected.*

Randomized controlled trial (RCT) An outcome study wherein participants are randomly allocated to an experimental group or a control or comparison group and followed over time on the variables or outcomes of interest. RCTs are capable of high levels of internal validity.*

Ratio data The highest measurement scale that, in addition to being an interval scale, also has an absolute zero in the scale.*

Regression to the mean A statistical phenomenon that can make natural variation in repeated data look like real change. It happens when unusually large or small measures tend to be followed by measurements that are closer to the mean. This may be a threat to the internal validity of the findings of a study.*

Replication Conducting a measurement, experiment, or study again; the second instance may be a repetition of the original study using different participants. If the study is repeated and produces the same findings, this enhances the validity and generalizability of the findings.*

Selection bias This occurs with differential selection of participants for comparison groups. Score differences, consequently, can be attributed to pretreatment differences among groups. This may be a threat to the internal validity of a study.*

Single-system research design A study conducted with one person, family, group, or system to explore the results of an intervention targeting specific outcomes. Typically, repeated measures of client functioning are taken prior to intervention, during intervention, and perhaps after the intervention is discontinued. These designs are a form of *idiographic research*—studies involving small numbers of participants—as opposed to nomothetic research involving large numbers of participants. Also known as single-participant, single-subject, or $N = 1$ research.*

Social desirability bias The tendency of people to answer questions in ways that are typically acceptable in a particular culture. This will generally take the form of over-reporting good behavior and underreporting bad behavior.*

STROBE An explicit set of methodological standards for quasi-experimental studies. The acronym stands for *STrengthening the Reporting of OBservational*

studies in Epidemiology. Many journals now require that submitted manuscripts be consistent with these standards.

Switching replications design A outcome study wherein one group of clients receives an intervention and a second, comparable group does not receive the intervention. The outcomes are then assessed. The second group then receives the same intervention, and is assessed to see if they responded to the intervention in a way similar to the first treatment group.

Systematic review A review of clearly formulated questions that uses systematic and explicit methods to identify, select, and critically appraise relevant research and to collect and analyze data from the studies that are included in the review.*

t test A parametric inferential statistical test used to examine the differences between the means of two groups (an independent samples *t* test) or the mean values on some outcome measure obtained from the same group on two occasions (e.g., before and after treatment), a paired sample *t* test.

Therapist bias/allegiance effects A possible confound in an outcome study wherein one or more therapists delivering experimental or other interventions possesses more or less adherence to each of these treatments. If this differential allegiance is reflected in how they interact with clients, this may bias the outcomes of the study. This is a possible threat to the internal validity of an outcome study.

Treatment fidelity Specific checks placed in a study to confirm that the manipulation of the independent variable occurred as planned. It includes such things as treatment definitions specified, implementer training, treatment manuals written, supervision of treatment agents, sampling for consistency, proper utilization of data collection strategies, and so on.*

Type I error A conclusion that a treatment or intervention works when it actually does not. The risk of a type I error is often called *alpha*. In a statistical test, it describes the chance of rejecting the null hypothesis when it is in fact true. It is also called a false positive.*

Type II error A conclusion that there is no evidence a treatment works when it actually does work. The risk of a type II error is often called *beta*. In a statistical test, it describes the change of not rejecting the null hypothesis when it is in fact false. The risk of a type II error decreases as the number of participants in a study increases. It is also called a false negative.*

References

- Abel, E. M., & Greco, M. (2008). A preliminary evaluation of an abstinence-oriented empowerment program for public school youth. *Research on Social Work Practice, 18*, 223–230.
- Albright, D. L., & Thyer, B. A. (2010). A test of the validity of the LCSW examination: *Quis custodiet ipsos custodes? Social Work Research, 34*, 229–234.
- American Psychological Association. (2009). *Publication Manual of the American Psychological Association* (6th ed.). Washington, DC: Author.
- Armour, S., & Haynie, D. L. (2007). Adolescent sexual debut and later delinquency. *Journal of Youth and Adolescence, 36*, 141–152.
- Bales, K. (1996). Lives and labours in the emergence of organised social research, 1886–1907. *Journal of Historical Sociology, 9*, 113–138.
- Barlow, D. H. (2010). Negative effects from psychological treatments: A perspective. *American Psychologist, 65*, 13–20.
- Bausell, R. B. (2007). *Snake oil science: The truth about complementary and alternative medicine*. New York: Oxford University Press.
- Berk, R. A., Sorenson, S. B., Wiebe, D. J., & Upchurch, D. M. (2003). The legalization of abortion and subsequent youth homicide: A time series analysis. *Analysis of Social Issues and Public Policy, 3*, 45–64.
- Biglan, A., Ary, D., & Waagenaar, A. C. (2000). The value of interrupted time-series experiments for community intervention research. *Prevention Science, 1*, 31–49.
- Blenkner, M. (1962). Control-groups and the placebo-effect in evaluative research. *Social Work, 7*, 52–58.
- Booth, C. (1902–1903). *Life and labour of the people of London*. London: Macmillan.
- Bordnick, P. S., Elkins, R. L., Orr, T. E., Walters, P., & Thyer, B. A. (2004). Evaluating the relative effectiveness of three aversion therapies designed to reduce craving among cocaine abusers. *Behavioral Interventions, 19*, 1–24.

188 References

- Bowen, G. L., & Farkas, G. (1991). Application of time-series designs to the evaluation of social services program initiatives: The recycling fund concept. *Social Work Research and Abstracts*, 27(3), 9–15.
- Brandell, J. R., & Varkas, T. (2010). Narrative case studies. In B. A. Thyer (Ed.). *Handbook of social work research methods* (2nd ed., pp. 376–396). Thousand Oaks, CA: Sage Publications.
- Buttall, F. P. (2002). Exploring levels of moral development among sex offenders participating in community-based treatment. *Journal of Offender Rehabilitation*, 34(4), 85–95.
- Campbell, D. T., & Stanley, J. C. (1963). *Experimental and quasi-experimental designs for research*. Chicago: Rand-McNally.
- Capp, H., Thyer, B. A., & Bordnick, P. S. (1997). Evaluating improvement over the course of adult psychiatric hospitalization. *Social Work in Health Care*, 25, 55–66.
- Carrigan, M. H., & Levis, D. J. (1999). The contributions of eye movements to the efficacy of brief exposure treatment for reducing fear of public speaking. *Journal of Anxiety Disorders*, 13, 101–118.
- Carrillo, D. F., Gallant, J. P., & Thyer, B. A. (1993). Training M.S.W. students in interviewing skills. *Arête*, 18(1), 12–19.
- Carrillo, D. F., & Thyer, B. A. (1994). Advanced standing and two-year program M.S.W. students: An empirical investigation of foundation interviewing skills. *Journal of Social Work Education*, 30, 278–288.
- Chandra, A., Martno, S. C., Collins, R. L., Elliott, M. N., Berry, S. H., Kanouse, D. E., & Mui, A. (2008). Does watching sex on television predict teen pregnancy? Findings from a national Longitudinal Survey of Youth. *Pediatrics*, 122, 1047–1054.
- Chapin, F. S. (1917). The experimental method and sociology. *The Scientific Monthly*, 4, 133–144.
- Chapin, F. S. (1949). The experimental method in the study of human relations. *The Scientific Monthly*, 68, 132–139.
- Chen, S. Y., Jordan, C., & Thompson, S. (2006). The effect of cognitive behavioral therapy (CBT) on depression: The role of problem-solving appraisal. *Research on Social Work Practice*, 16, 500–510.
- Chu, Y. H., Frongillo, E. A., Jones, S. J., & Kaye, G. L. (2009). Improving patrons' meal selections through the use of point-of-selection nutritional labels. *American Journal of Public Health*, 99, 2001–2005.
- Clapp, J. D., Lange, J. E., Russell, C., Shillington, A., & Voas, R. B. (2003). A failed norms social marketing campaign. *Journal of Studies on Alcohol*, 64, 409–414.
- Cnaan, R., & Tripodi, S. (2010). Randomized controlled experiments. In B. A. Thyer (Ed.). *Handbook of social work research methods* (2nd ed., pp. 205–220). Thousand Oaks, CA: Sage.

- Cohen, J. (1994). The earth is round ($p < 0.05$). *American Psychologist*, *49*, 997–1003.
- Colosetti, S. D., & Thyer, B. A. (2000). The relative effectiveness of EMDR versus relaxation training with battered women prisoners. *Behavior Modification*, *24*, 719–739.
- Cook, T. D., & Campbell, D. T. (1979). *Quasi-experimentation: Design and analysis for field settings*. Boston: Houghton-Mifflin.
- Cook, T. D., & Shadish, W. R. (1994). Social experiments: Some developments over the past 15 years. *Annual Review of Psychology*, *45*, 545–580.
- Copp, H. L., Bordnick, P. S., Traylor, A. C., & Thyer, B. A. (2007). Evaluating wraparound services for seriously emotionally disturbed youth: Pilot study outcomes in Georgia. *Adolescence*, *42*, 723–732.
- Council on Social Work Education. (2008). *Educational Policy and Accreditation Standards*. Alexandria, VA: Author. Available at <http://www.cswe.org/File.aspx?id=13780>.
- Crolley, J., Roys, D., Thyer, B. A., & Bordnick, P. S. (1998). Evaluating outpatient behavior therapy for sex offenders: A pretest-posttest study. *Behavior Modification*, *22*, 485–501.
- Crow, S. J., Mitchell, J. E., Roerig, J. D., & Steffen, K. (2009). What potential role is there for medication treatment in Anorexia Nervosa? *International Journal of Eating Disorders*, *42*, 1–9.
- Dattalo, P. (2007). *Determining sample size: Balancing power, precision, and practicality*. New York: Oxford University Press.
- Davis, L. V. (1994). Rejoinder to Dr. Marsh. In W. Hudson & P. S. Nurius (Eds.), *Controversial issues in social work research* (pp. 73–74). Boston, MA: Allyn & Bacon.
- Davis, R. M. (2008). British American Tobacco ghost-wrote reports on tobacco advertising bans by the Advertising Association and J. J. Boddewyn. *Tobacco Control*, *17*, 211–214.
- DeAngelis, C. D., & Fontana, P. B. (2008). Guest authorship, mortality reporting, and integrity in Rofecoxib studies reply. *JAMA*, *300*, 905–906.
- Demicheli, V., Jefferson, T., Rivetti, A., & Price, D. (2008). *Vaccines for measles, mumps and rubella in children*. *Cochrane Database of Systematic Reviews*, *4*, Art. No. CD004407. DOI: 10.1002/14651858.CD004407.pub2.
- de Schmidt, G. A., & Gorey, K. M. (1997). Unpublished social work research: Systematic replication of a recent meta-analysis of published intervention effectiveness research. *Social Work Research*, *21*, 58–62.
- DeWalt, D. A., Davis, T. C., Wallace, A. S., Seligman, H. K., Bryant-Shilliday, B., Arnold, C. L., Freburger, J., & Shillinger, D. (2009). Goal setting in diabetes self-management: Taking the baby steps to success. *Patient Education and Counseling*, *77*, 218–223.

190 References

- Devine, E. T. (1908). Results of the Pittsburgh Survey. *Proceedings of the American Sociological Society*, 3, 85–92.
- Diehl, D., & Frey, A. (2008). Evaluating a community-school model of social work practice. *School Social Work Journal*, 32(2), 1–20.
- Dillon, S. (22 January 2009). Study sees an Obama effect as lifting Black test takers. *The New York Times*, A15. Accessed May 10, 2010, from <http://www.nytimes.com/2009/01/23/education/23gap.html>.
- DiNitto, D. (1983). Time-series analysis: An application to social welfare policy. *Journal of Applied Behavioral Science*, 19, 507–518.
- DiNitto, D. M., McDaniel, R. R., Ruefli, T. W., & Thomas, J. B. (1986). The use of ordinal time-series analysis in assessing policy inputs and impacts. *Journal of Applied Behavioral Science*, 22, 77–93.
- Donohue, B., & Thyer, B. A. (1992). Should the GRE be used as an admissions requirement by schools of social work? *Journal of Teaching in Social Work*, 6, 33–40.
- DuBois, W. E. B. (1899). *The Philadelphia Negro*. Philadelphia, PA: University of Pennsylvania.
- Duncan, T. E., & Duncan, S. C. (2004a). A latent growth curve modeling approach to pooled interrupted time series data. *Journal of Psychopathology and Behavioral Assessment*, 26, 271–278.
- Duncan, T. E., & Duncan, S. C. (2004b). An introduction to latent growth curve modeling. *Behavior Therapy*, 35, 333–363.
- Epstein, W. M. (1990). Rational claims to effectiveness in social work's critical literature. *The Social Science Journal*, 27, 129–145.
- Fischer, J. (1973). Is casework effective? A review. *Social Work*, 18, 5–20.
- Fischer, J. (1976). *The effectiveness of social casework*. Springfield, IL: Charles C. Thomas.
- Freud, S. (1962/1894). Obsessions and compulsions. In J. Strachey (Ed.). *The standard edition of the complete psychological works of Sigmund Freud* (Vol. 3, p. 81). London: Hogarth Press.
- Geron, S. M., & Stetekee, G. S. (2010). Applying for research grants. In B. A. Thyer (Eds.). *Handbook of social work research methods* (2nd ed., pp. 619–630). Thousand Oaks, CA: Sage Publications, Inc.
- Gordon, M. (1973). The social survey movement and sociology in the United States. *Social Problems*, 21, 284–298.
- Gorey, K. M., Thyer, B. A., & Pawluck, D. E. (1998). Differential effectiveness of prevalent social work practice models. *Social Work*, 43, 269–278.
- Grembowski, D., & Milgrom, P. M. (2000). Increasing access to dental care for Medicaid preschool children: The Access to Baby and Child Dentistry (ABCD) program. *Public Health Reports*, 115, 448–459.

- Grenier, A. M., & Gorey, K. M. (1998). The effectiveness of social work with older people and their families: A meta-analysis of conference proceedings. *Social Work Research, 22*, 60–64.
- Grey, A. L., & Dermody, H. E. (1972). Reports of casework failure. *Social Casework, 53*, 534–543.
- Halpern-Meekin, S., & Tach, L. (2008). Heterogeneity in two-parent families and adolescent well-being. *Journal of Marriage and Family, 70*, 435–451.
- Harrison, D. F., & Thyer, B. A. (1988). Doctoral research on social work practice. *Journal of Social Work Education, 24*, 107–114.
- Herbert, J. D., Sharp, I. R., & Guidano, B. A. (2002). Separating fact from fiction in the etiology and treatment of autism: A scientific review of the evidence. *The Scientific Review of Mental Health Practice, 1*, 23–43.
- Herron, W. G., & Sitkowski, S. (1986). Effect of fees on psychotherapy: What is the evidence? *Professional Psychology: Research and Practice, 17*, 347–351.
- Hogarty, G. E. (1989). Meta-analysis of the effects of practice with the chronically mentally ill: A critique and reappraisal of the literature. *Social Work, 34*, 363–373.
- Holden, G., Thyer, B. A., Baer, J., Delva, J., Dulmus, C. N., & Shanks, T. W. (2008). Suggestions to improve social work journal editorial and peer-review processes: The San Antonio response to the Miami Statement. *Research on Social Work Practice, 18*, 66–71.
- Holder, H. D., & Wagenaar, A. C. (1994). Mandated server training and reduced alcohol-involved traffic crashes: A time-series analysis of the Oregon experience. *Accident Analysis & Prevention, 26*, 89–97.
- Holland, J. F. (1997). Clinical trials in cancer. *Clinical Cancer Research, 3*, 2585–2586.
- Holosko, M. J. (2006). A suggested author's checklist for submitting manuscripts to *Research on Social Work Practice*. *Research on Social Work Practice, 16*, 449–454.
- Holosko, M. J. (2010). What types of designs are we using in social work research and evaluation? *Research on Social Work Practice, 20*, 665–673.
- Holosko, M. J., Thyer, B. A., & Danner, J. E. (2009). Ethical guidelines for designing and conducting evaluations of social work practice. *Journal of Evidence-based Social Work, 6*, 1–13.
- Hudson, W. W., Thyer, B. A., & Stocks, J. T. (1985). Assessing the important of experimental outcomes. *Journal of Social Service Research, 8*(4), 87–98.
- Hyun, M. K., Lee, M. S., Kang, K., & Choi, S. M. (2008). Body acupuncture for Nicotine Withdrawal Symptoms: A randomized placebo-controlled trial. *Complementary and Alternative Medicine, 7*, 233–238.
- Jainchill, N., Hawke, J., & Messina, M. (2005). Post-treatment outcomes among adjudicated adolescent males and females in modified therapeutic community treatment. *Substance Use & Misuse, 40*, 975–996.

- Johnson, P. A., Thyer, B. A., Daniels, M., Anderson, R., & Bordnick, P. S. (1996). Is the school social worker examination valid? *Arete*, 20(2), 1–5.
- Johnson, Y. M., & Stadel, V. L. (2007). Completion of advance directives: Do social work preadmission interviews make a difference? *Research on Social Work Practice*, 17, 686–696.
- Jones, C. D., Chancey, R., Lowe, L. A., & Risler, E. A. (2010). Residential treatment for sexually abusive youth: An assessment of treatment outcomes. *Research on Social Work Practice*, 20, 172–182.
- Keller, D. P., Schut, L. J., Puddy, R. W., Williams, L., Stephens, R. L., McKeon, R., & Lubell, K. (2009). Tennessee lives count: Statewide gatekeeper training for youth suicide prevention. *Professional Psychology: Research and Practice*, 40, 126–133.
- Klosko, J. S., Barlow, D. H., Tassinari, R., & Czerny, J. A. (1990). A comparison of alprazolam and behavior therapy in treatment of panic disorder. *Journal of Consulting and Clinical Psychology*, 58, 77–84.
- Knaevelsrud, C., & Maercker, A. (2007). Internet-based treatment for PTSD reduces distress and facilitates the development of a strong therapeutic alliance: A randomized controlled trial. *BMC Psychiatry*, 7, 13.
- Knaevelsrud, C., & Maercker, A. (2010). Long-term effects of an internet-based treatment for posttraumatic stress. *Cognitive Behaviour Therapy*, 39, 72–77.
- Kurdyak, P., Cairney, J., Sarnocinska-Hart, Callahan, R. C., & Strike, C. (2008). The impact of a smoking cessation policy on visits to a psychiatric emergency department. *Canadian Review of Psychiatry*, 53, 779–782.
- Lachin, J. M. (2000). Statistical considerations in the intent-to-treat principle. *Controlled Clinical Trials*, 21, 167–189.
- Larsen, J., & Hepworth, D. H. (1982). Skill development of helping skills in undergraduate social work education: Model and evaluation. *Journal of Education for Social Work*, 18, 66–73.
- LeCroy, C. W., & Krysik, J. (2007). Understanding and interpreting effect size measures. *Social Work Research*, 31, 243–248.
- Levinson, D. R. (2010). *Most Medicaid children in nine states are not receiving all required preventive screening services*. Washington, DC: Department of Health and Human Services, Office of the Inspector General. Accessed May 25, 2010, from <http://oig.hhs.gov/oei/reports/oei-05-08-00520.pdf>.
- Ligon, J., & Thyer, B. A. (2000). Client and family satisfaction with brief community mental health, substance abuse, and mobile crisis services in an urban setting. *Crisis Intervention*, 6, 93–99.
- Lilienfeld, S. O. (2007). Psychological treatments that cause harm. *Perspectives on Psychological Science*, 2, 53–70.
- Littell, J. H., Corcoran, J., & Pillai, V. (2008). *Systematic reviews and meta-analyses*. New York: Oxford University Press.

- Lowry, F. (Ed.). (1939). *Readings in social casework: 1920–1938*. New York: Columbia University Press.
- Macdonald, M. E. (1953). Some essentials in the evaluation of social casework. *Journal of Psychiatric Social Work*, 22(3), 135–137.
- Martin, V. (2005). The consequences of parental divorce on the life course outcomes of Canadian children. *Canadian Studies in Population*, 32, 29–51.
- Martino, S. C., Elliott, M. N., Collins, R. L., Kanouse, D. E., & Berry, S. H. (2008). Virginity pledges among the willing: Delays in first intercourse and consistency in condom use. *Journal of Adolescent Health*, 43, 341–348.
- McDowall, D., McCleary, R., Meddinger, E. E., & Hay, A. J. (1980). *Interrupted time series analysis*. Thousand Oaks, CA: Sage.
- Michielutte, R., Shelton, B., Paskett, E. D., Tatum, C. M., & Velez, R. (2000). Use of an interrupted time-series design to evaluate a cancer screening program. *Health Education Research*, 15, 615–623.
- Monnickendam, M., & Elliot, E. J. (1997). Effects of a practice-centered, cognitive-oriented computer course on computer attitudes: Implications for course content. *Social Work and Social Sciences Review*, 6, 175–185.
- Moore, A., & McQuay, H. (2006). *Bandolier's little book of making sense of the medical evidence*. New York: Oxford University Press.
- Myers, L. L., & Rittner, B. (2001). Adult psychosocial functioning of children raised in an orphanage. *Residential Treatment of Children and Youth*, 18(4), 3–21.
- National Association of Social Workers. (2008). *Code of Ethics*. Washington, DC: NASW Press.
- Nezu, A. M., & Nezu, C. M. (2008). *Evidence-based outcome research: A practical guide to conducting randomized controlled trials for psychosocial interventions*. New York: Oxford University Press.
- Novak, I., Cusick, A., & Lowe, K. (2007). A pilot study on the impact of occupational therapy home programming for young children with cerebral palsy. *American Journal of Occupational Therapy*, 61, 463–468.
- Newsome, W. S., Anderson-Butcher, D., Fink, J., Hall, L., & Huffer, J. (2008). The impact of school social work services on student absenteeism and risk factors related to school truancy. *School Social Work Journal*, 32(2), 21–38.
- Novella, S. (2011). What is acupuncture? *The Skeptical Inquirer*, 35(4), 28–29.
- Ottenbacher, K. J. (1997). Designing and interpreting clinical studies. In M. Fuhrer (Ed.), *Assessing medical rehabilitation practices* (pp. 233–256). Baltimore, MD: Paul H. Brookes.
- Pabian, W., Thyer, B. A., Straka, E., & Boyle, P. (2000). Do the families of children with developmental disabilities obtain recommended services? A follow-up study. *Journal of Human Behavior in the Social Environment*, 3(1), 45–58.

- Palmgreen, P., Lorch, E. P., Stephenson, M. T., Hoyle, R. H., & Donohew, L. (2007). Effects of the Office of National Drug Control Policy's marijuana initiative campaign on high-sensation-seeking adolescents. *American Journal of Public Health, 97*, 1644–1649.
- Parrish, D. E., & Rubin, A. (2011). An effective model for continuing education training in evidence-based practice. *Research on Social Work Practice, 21*, 77–87.
- Perry, R. E. (2006). Do social workers make better child welfare workers than non-social workers? *Research on Social Work Practice, 16*, 392–405.
- Pharis, M. E. (1976). Ten reasons why I am not bothered by outcome studies which claim to show psychotherapy is ineffective. *Clinical Social Work Journal, 4*, 58–61.
- Pignotti, M. (2005). Thought Field Therapy Voice Technology vs. random meridian point sequences: A single-blind controlled experiment. *The Scientific Review of Mental Health Practice, 4*(1), 72–81.
- Pignotti, M. & Thyer, B. A. (2009). Why randomized clinical trials are important and necessary to social work practice. In H-W. Otto, A. Polutta, & H. Ziegler (Eds.), *Evidence-based practice: Modernizing the knowledge base of social work* (pp. 99–109). Farmington Hills, MI/Opladen, Germany: Barbara Budrich Publishers.
- Randall, E. J., & Thyer, B. A. (1994). A preliminary test of the validity of the LCSW examination. *Clinical Social Work Journal, 22*, 223–227.
- Raphael, B. (1986). *When disaster strikes: How individuals and communities cope with disaster*. New York: Basic Books.
- Raskin, M., Johnson, G., & Rondestvedt, J. W. (1973). Chronic anxiety treated by feedback-induced muscle relaxation: A pilot study. *Archives of General Psychiatry, 28*, 263–267.
- Reid, W. J., & Hanrahan, P. (1982). Recent evaluations of social work: Grounds for optimism. *Social Work, 27*, 328–340.
- Reid, W. J., & Shyne, A. W. (1969). *Brief and extended casework*. New York: Columbia University Press.
- Residents of Hull House. (1895). *Hull-House maps and papers; A presentation of nationalities and wages in a congested district of Chicago, together with comments and essays on problems growing out of the social conditions*. New York: Crowell.
- Reynolds, A. J., Temple, J. A., Robertson, D. L., & Mann, E. A. (2001). Long-term effects of an early childhood intervention on educational achievement and juvenile arrest: A 15-year follow-up of low-income children in public schools. *JAMA, 285*, 2339–2346.
- Richards, K. V., & Thyer, B. A. (2011). Does Individual Development Account participation help the poor? A review. *Research on Social Work Practice, 21*, 348–362.

- Rogers, C. R. (1933). A good foster home: Its achievements and limitations. *Mental Hygiene*, 17, 21–40.
- Rosenbaum, J. E. (2008). Patient teenagers? A comparison of the sexual behavior of virginity pledgers and matched nonpledgers. *Pediatrics*, 123, e110–e120.
- Rosenthal, D., & Frank, J. D. (1956). Psychotherapy and the placebo effect. *Psychological Bulletin*, 53, 294–302.
- Royse, D., Thyer, B. A., & Padgett, D. (2010). *Program evaluation: An introduction* (5th ed.). Belmont, CA: Brooks/Cole/Cengage.
- Rubin, A., & Babbie, E. R. (2008). *Research methods for social work* (6th ed.). Belmont, CA: Thomson.
- Rubin, A., & Parrish, D. (2007). Problematic phrases in the conclusions of published outcome studies: Implications for evidence-based practice. *Research on Social Work Practice*, 17, 334–347.
- Salloum, A. (2008). Group therapy for children after homicide and violence: A pilot study. *Research on Social Work Practice*, 18, 198–211.
- Sanderson, S., Tatt, I. D., & Higgins, J. P. T. (2007). Tools for assessing quality and susceptibility to bias in observational studies in epidemiology: a systematic review and annotated bibliography. *International Journal of Epidemiology*, 36, 666–676.
- Sandler, J. C., Freeman, N. J., & Socia, K. M. (2008). Does a watched pot boil? A time-series analysis of New York state's Sex Offender Registration and Notification law. *Psychology, Public Policy, and Law*, 14, 284–302.
- Schilling, R., Baer, J. C., Barth, R., Fraser, M., Herman, D., Holden, G. et al. (2005). Peer review and publication standards in social work journals: The Miami Statement. *Social Work Research*, 29, 119–121.
- Schissel, B. (1996). Law reform and social change: A time-series analysis of sexual assault in Canada. *Journal of Criminal Justice*, 24, 123–138.
- Schneider, D. J., May, G., Carithers, T., Coyle, K., Potter, S., Endahl, J., Robin, L., McKenna, M., Debrot, K., & Seymour, J. (2006). Evaluation of a fruit and vegetable distribution program—Mississippi, 2004–05 school year. *Journal of the American Medical Association*, 296, 1833–1834.
- Segal, S. P. (1972). Research on the outcome of social work therapeutic interventions: A review of the literature. *Journal of Health and Social Behavior*, 13, 3–17.
- Seekins, T., Fawcett, S. B., Cohen, S. H., Elder, J. P., Jason, L. A., Schnelle, J. F., & Winett, R. A. (1988). Experimental evaluation of public policy: The case of state legislation for child passenger safety. *Journal of Applied Behavior Analysis*, 21, 233–243.
- Shadish, W. R. (2011). Randomized controlled studies and alternative designs in outcome studies: Challenges and opportunities. *Research on Social Work Practice*, 21, xxx–xxx.

- Shadish, W. R., Clark, M. H., & Steiner, P. M. (2008). Can nonrandomized experiments yield accurate answers? A randomized experimental comparing random to nonrandom assignment. *Journal of the American Statistical Association*, *103*, 1334–1343.
- Shadish, W. R., Cook, T. D., & Campbell, D. T. (2002). *Experimental and quasi-experimental designs for generalized causal inference*. New York: Houghton Mifflin.
- Shadish, W. R., Galindo, R., Wong, V., Steiner, P. M., & Cook, T. D. (2011). A randomized experiment comparing random and cutoff-based assignment. *Psychological Methods*, *16*, 179–191.
- Shadish, W. R., Matt, G. E., Navarro, A. M., & Phillips, G. (2000). The effects of psychological therapies under clinically representative conditions: A meta-analysis. *Journal of Consulting and Clinical Psychology*, *126*, 512–529.
- Shadish, W. R., & Ragsdale, K. (1996). Random versus nonrandom assignment in controlled experiments: Do you get the same answer? *Journal of Consulting and Clinical Psychology*, *64*, 1290–1305.
- Shipton, B., & Spain, A. (1981). Implications of payment of fees for psychotherapy. *Psychotherapy: Theory, Research and Practice*, *18*, 68–73.
- Smeeth, L., Cook, C., Fombonne, E., Heavey, L., Rodrigues, L. C., Smith, P. G., & Hall, A. J. (2004). MMR vaccination and pervasive developmental disorders: A case-control study. *The Lancet*, *364*, 963–969.
- Smith, G. C. S., & Pell, J. P. (2003). Parachute use to prevent death and major trauma related to gravitational challenge: Systematic review of randomized controlled trials. *BMJ*, *327*, 1459–1461.
- Smith, G. S., Thyer, B. A., Clements, C., & Kropf, N. P. (1997). An evaluation of coalition building training for aging and developmental disabilities service providers. *Educational Gerontology*, *23*, 105–114.
- Solomon, P., Cavanaugh, M. M., & Draine, J. (2009). *Randomized controlled trials*. New York: Oxford University Press.
- Spinelli, M. (1997). Interpersonal psychotherapy for depressed antepartum women: A pilot study. *American Journal of Psychiatry*, *154*, 1028–1030.
- Staff. (2010). Should protocols for observational research be registered? *The Lancet*, *375*, 348.
- Stewart, R. E., & Chambless, D. L. (2009). Cognitive-behavioral therapy for adult anxiety disorders in clinical practice: A meta-analysis of effectiveness studies. *Journal of Consulting and Clinical Psychology*, *77*, 595–606.
- Sze, W. C., Keller, R. S., & Keller, D. B. (1979). A comparative study of two different teaching and curricular arrangements in human behavior and social environments. *Journal of Education for Social Work*, *15*, 103–108.
- Thomas, E. J. (1975). Uses of research methods in interpersonal practice. In N. A. Polansky (Ed.), *Social work research* (revised edition, pp. 254–283). Chicago: University of Chicago Press.

- Thomlison, R. J. (1984). Something works: Evidence from practice effectiveness studies. *Social Work*, 29, 51–56.
- Thyer, B. A. (1987). Contingency contracting to promote automobile safety belt use by students (letter). *The Behavior Therapist*, 10, 150, 160.
- Thyer, B. A. (1988). Teaching without testing: A preliminary report of an innovative technique for social work education. *Innovative Higher Education*, 13, 47–53.
- Thyer, B. A. (1991). Guidelines for evaluating outcome studies on social work practice. *Research on Social Work Practice*, 1, 76–91.
- Thyer, B. A. (2002). How to write up a social work outcome study for publication. *Journal of Social Work Research and Evaluation*, 3, 215–224.
- Thyer, B. A. (2008). *Preparing research articles*. New York: Oxford University Press.
- Thyer, B. A. (Ed.) (2010). *Handbook of social work research methods*. Thousand Oaks, CA: Sage Publications, Inc.
- Thyer, B. A. (2011). LCSW examination pass rates: Implications for social work education. *Clinical Social Work Journal*, 39, 296–300. DOI: 10.1007/s10615-009-0253-x.
- Thyer, B. A., Jackson-White, G., Sutphen, R., & Carrillo, D. F. (1992). Structured study questions as a social work teaching method. *Innovative Higher Education*, 16, 235–245.
- Thyer, B. A., & Myers, L. L. (2007). *A social worker's guide to evaluating practice outcomes*. Alexandria, VA: Council on Social Work Education.
- Thyer, B. A., Myers, L. L., & Nugent, W. R. (2011). Do regular social work faculty earn better student course evaluations than adjunct faculty or doctoral students? *Journal of Teaching in Social Work*, 31, 365–377.
- Thyer, B. A., Sowers-Hoag, K. M., & Love, J. P. (1986). The influence of field instructor-student gender combinations on student perceptions of field instruction quality. *Arete*, 11(2), 25–30.
- Thyer, B. A., & Vodde, R. (1994). Is the ACSW examination valid? *Clinical Social Work Journal*, 22, 105–122.
- Thyer, B. A., Vonk, M. E., & Tandy, C. C. (1996). Are advanced standing and two-year MSW program students equivalently prepared? An empirical investigation. *Arete*, 20(2), 42–46.
- Tripodi, T., & Harrington, J. (1979). Uses of time-series designs for formative program evaluation. *Journal of Social Service Research*, 3(1), 67–78.
- Turner, E. H., Matthews, A., Linardos, E., Tell, R. A., & Rosenthal, R. (2008). Selective publication of antidepressant trials and its influence on apparent efficacy. *The New England Journal of Medicine*, 358, 252–260.
- Videka-Sherman, L. (1988). Meta-analysis of research on social work practice in mental health. *Social Work*, 33, 325–338.

- Viggiani, P. A., Reid, W. J., & Bailey-Dempsey, C. (2002). Social worker-teacher collaboration in the classroom: Help for elementary students at risk. *Research on Social Work Practice, 12*, 604–620.
- Vingilis, E., McLeod, A. L., Seeley, J., Mann, R., Voas, R., & Compton, C. (2006). The impact of Ontario's extended drinking hours on cross-border cities of Windsor and Detroit. *Accident Analysis and Prevention, 38*, 63–70.
- Von Elm, E., Altman, D. G., Egger, M., Pocock, S. J., Gotsche, P. C., & Vandembrouche, J. P. (2007). The strengthening and Reporting of Observational Studies in Epidemiology (STROBE) statement: Guidelines for reporting observational studies. *Epidemiology, 18*, 800–804.
- Vonk, E. M., & Thyer, B. A. (1999). Evaluating the effectiveness of short-term treatment at a university counseling center. *Journal of Clinical Psychology, 55*, 1095–1106.
- Vonk, M. E., Zucrow, E., & Thyer, B. A. (1996). Female MSW students' satisfaction with practicum supervision: The effect of supervisor gender. *Journal of Social Work Education, 32*, 415–419.
- Waber, R. L., Shiv, B., Camon, Z., & Ariely, D. (2008). Commercial features of placebo and therapeutic efficacy. *JAMA, 299*, 1016–1017.
- Wagenaar, A. C., & Maldonado-Molina, M. M. (2007). Effects of drivers license suspension policies on alcohol-related crash involvement: Long-term follow-up in forty-six states. *Alcoholism: Clinical and Experimental Research, 31*, 1399–1406.
- Wagenaar, A. C., Maldonado-Molina, M. M., Erickson, D. J., Ma, L., Tobler, A. L., & Komro, K. A. (2007). General deterrence effects of U.S. statutory DUI fine and jail penalties: Long-term follow-up in 32 states. *Accident Analysis and Prevention, 39*, 982–994.
- Wagenaar, A. C., & Toomey, T. L. (2002). Effects of minimum drinking age laws: Review and analyses of the literature from 1960 to 2000. *Journal of Studies on Alcohol, Supplement No. 14*, 206–225.
- Welner, A., Welner, Z., & Fishman, R. (1979). Psychiatric adolescent inpatients: Eight- to ten-year follow-up. *Archives of General Psychiatry, 36*, 698–700.
- Whitt-Glover, M. C., Hogan, P., Lang, W., & Heil, D. P. (2008). Pilot study of a faith-based physical activity program among sedentary Blacks. *Preventing Chronic Disease, 5*(2), A51.
- Wicke, T., & Silver, R. C. (2009). A community responds to collective trauma: An ecological analysis of the James Byrd murder in Jasper, Texas. *American Journal of Community Psychology, 44*, 233–248.
- Wilson, J. M., Wallace, L. S., & DeVoe, J. E. (2009). Are state Medicaid application enrollment forms readable? *Journal of Health Care for the Poor and Underserved, 20*, 423–431.

- Wolf, D. B., & Abell, N. B. (2003). Examining the effects of meditation techniques on psychosocial functioning. *Research on Social Work Practice, 13*, 27–42.
- Wood, W. G. (1982). Do fees help heal? *Journal of Clinical Psychology, 38*, 669–673.
- Yegidis, B., Weinbach, R. W., & Myers, L. L. (2011). *Research methods for social workers* (7th ed.). New York: Allyn and Bacon.
- Yoken, C., & Berman, J. S. (1984). Does paying a fee alter the effectiveness of treatment? *Journal of Consulting and Clinical Psychology, 52*, 254–260.
- Ziman, J. (1978). *Reliable knowledge: An exploration of the grounds for belief in science*. Cambridge: Cambridge University Press.
- Zite, N. B., Philipson, S. J., & Wallace, L. S. (2007). Consent to Sterilization section of the Medicaid-Title XIX form: Is it understandable? *Contraception, 75*, 256–260.

