

Tilburg University

Methods and data analysis of comparative research

van de Vijver, F.J.R.; Leung, K.

Published in:
Handbook of cross-cultural psychology, 2nd ed.

Publication date:
1997

[Link to publication in Tilburg University Research Portal](#)

Citation for published version (APA):
van de Vijver, F. J. R., & Leung, K. (1997). Methods and data analysis of comparative research. In J. W. Berry, Y. H. Poortinga, & J. Pandey (Eds.), *Handbook of cross-cultural psychology, 2nd ed.* (pp. 257-300). (2nd. ed.; No. vol. 1). Allyn & Bacon.

General rights

Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

**Handbook of Cross-Cultural Psychology
Second Edition**

Edited by John W. Berry, Ype H. Poortinga, Janak Pandey,
Pierre R. Dasen, T. S. Saraswathi, Marshall H. Segall,
and Cigdem Kagitçibasi

VOLUME 1

Theory and Method

VOLUME 2

Basic Processes and Human Development

VOLUME 3

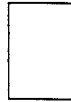
Social Behavior and Applications

性相近 習相遠

**Basic human nature is similar at birth;
Different habits make us seem remote.**

From the *San Zi Jing*

Second Edition



Handbook of Cross-Cultural Psychology

VOLUME 1
THEORY AND METHOD

Edited by

John W. Berry
Queen's University, Canada

Ype H. Poortinga
Tilburg University, The Netherlands

Janak Pandey
University of Allahabad, India

Allyn and Bacon

Boston • London • Toronto • Sydney • Tokyo • Singapore

7

METHODS AND DATA ANALYSIS OF COMPARATIVE RESEARCH

FONS VAN DE VIJVER
Tilburg University
The Netherlands

KWOK LEUNG
Chinese University of Hong Kong
Hong Kong

Contents

Introduction	259
Specific Issues in Comparative Methodology and Data Analysis	259
Equivalence	261
Methods	262
Sampling of Cultures	262
Sampling of Subjects	264
Procedure	264
Instrument Translation	266
Administration	267
Design	269
Data Analysis	271
Preliminary Analyses	271
Establishing Scalar Equivalence	279
Statistical Tests of Cross-Cultural Differences: Introduction	280
Statistical Tests of Cross-Cultural Differences:	
Level-Oriented Techniques	281
Statistical Tests of Cross-Cultural Differences:	
Structure-Oriented Techniques	283
Four Common Types of Comparative Studies	287
Methods and Analysis of Four Common Types of	
Comparative Studies	289
Conclusion	294
References	294

Introduction

The major goal of this chapter is to provide a comprehensive overview of the methodological issues encountered in cross-cultural research. Since the reviews in the first edition of the *Handbook* on testing and assessment by Irvine and Carroll (1980) and on experimentation by Brown and Sechrest (1980), many developments have taken place. In our presentation, we focus on data sets that are comparative in nature. Most studies of this type involve data from at least two cultural groups, but some studies are monocultural. In such studies, previous work must provide data and results before meaningful cross-cultural comparisons to be made. Monocultural studies commonly conducted by ethnographers and anthropologists that do not touch upon cross-cultural comparison will fall outside of the scope of our review.

We see the process of conducting cross-cultural research as composed of three important steps. First, the research questions must be explicitly stated. Second, a method that is appropriate to the research questions raised should be selected. Method is defined here as the design, sampling, administration, and instrumentation involved in the collection of data. Finally, the appropriate data analysis should be chosen in light of the research questions raised and the method chosen. We consider these three steps as intertwined, and they should be considered simultaneously prior to data collection. This three-step framework is used in organizing the materials that follow.

The first section of the chapter describes specific issues of cross-cultural research, such as quasi-experimentation. The second section describes in more detail the methodological aspects of cross-cultural studies. The third section deals with the analysis of cross-cultural data. The fourth section reviews the main issues in the methodology and analysis of four common types of cross-cultural studies. Conclusions are drawn in the final section.

Specific Issues in Comparative Methodology and Data Analysis

Before the methods and analyses of cross-cultural studies can be discussed, the applicability of "true experiments" (Campbell & Stanley, 1966) and the associated statistical framework to these studies—the Neyman–Pearson theory—should be explored in order to highlight their special characteristics.

The classical Neyman–Pearson theory provides the most commonly applied statistical framework in testing intergroup differences in psychology. The framework is appropriate for analyzing data from experiments with experimental and control groups. The two groups are considered to be equal, except for the manipulation that is present in the experimental group and absent in the control group ("all other things being equal," as it is often called). The theoretical question the researcher wants to examine concerns the presence of a difference in the dependent measures between the experimental and control groups. This is tested by a *t* test or analysis of variance. The researcher chooses *a priori* a probability

that is considered appropriate, usually .05 or .01, for concluding whether or not there are differences between the experimental and control groups. The framework has been developed as a tool to analyze data collected in experimental settings and to reduce the risk of making false inferences.

The framework has turned out to work well mainly in so-called true experiments (Campbell & Stanley, 1966), in which subjects are randomly assigned to different experimental conditions. The "all other things being equal" argument in general does not apply to studies in which subjects are not assigned randomly to experimental treatments. Group membership, a major experimental treatment in cross-cultural studies, is predetermined and cannot be randomly assigned. When the cultural differences between the groups of subjects involved in a cross-cultural study are extensive, it does not make much sense to assume the validity of the "all other things being equal" argument and to compare the groups as if the data were collected in a true experiment. In cross-cultural studies, the application of the Neyman–Pearson framework can yield misleading results (Poortinga & Malpass, 1986). For instance, when cognitive tests are presented to Western literate and non-Western illiterate subjects, the educational and cultural differences between the two groups tend to be so massive that a test of the null hypothesis of no intergroup differences in performance is inadequate. Quite likely, every item will show a significant difference between the two groups.

Furthermore, the interpretation of such a test is equivocal. In the experimental paradigm the interpretation of the difference in the dependent measures is simple. The treatment, typically well defined, such as a drug that has been administered, has produced the score difference between the experimental and control groups. In a similar vein, the differences in the cognitive tests between the literates and illiterates can be attributed to the treatment "culture." However, the attribution does not convey much meaning and dodges the question of a proper interpretation of the score differences. Culture is too global a concept to be used as a meaningful independent variable in the interpretation. In comparison with the experimental branches of psychology, cross-cultural psychology should be more sensitive to the interpretability of findings. Whereas the task often ends for the experimental psychologist with the observation of a significant difference because the observation will typically confirm or falsify a hypothesis, the task of the cross-cultural psychologist is certainly not complete with the observation of significant intergroup differences.

A crucial problem in quasi-experiments (in which there is no random assignment of subjects) is the ruling out of rival hypotheses. This issue has been extensively discussed, and the consensus is that culture must be "unpacked" (e.g., Whiting, 1976; Poortinga, Van de Vijver, Joe, & Van de Koppel, 1989). The use of culture as an explanatory variable is not satisfactory, and culture must be decomposed into a set of psychologically meaningful constructs, which are then used to explain the cultural differences observed (e.g., Leung, 1989; Poortinga & Van de Vijver, 1987). When cultural differences on a dependent variable are documented, it is almost impossible to pin down which aspect of culture is responsible for the observed differences, in the absence of additional data. A feasible strategy is to

identify the most likely variables that may account for the expected cultural differences and measure these variables in the study. A number of analytical procedures, which will be described later, can be employed to identify which aspect is indeed the most plausible explanation for the cultural differences observed. An adequate cross-cultural study must have built-in elements in its design to rule out plausible rival hypotheses (cf. Cook & Campbell, 1979).

Equivalence

Equivalence is a major concern in cross-cultural research; meaningful cross-cultural comparisons can only be made if the data from different cultures are comparable. Equivalence has been discussed extensively, and several types have been identified (e.g., Berry, 1969; Poortinga, 1971, 1989; Van de Vijver & Poortinga, 1982). Because terms used in the literature to describe equivalence are often unclear and confusing, we propose that three types of equivalence be distinguished: *structural*, *measurement unit*, and *scalar*. Cross-cultural researchers are often interested in *structural equivalence*, which refers to the similarity of psychometric properties of data sets from different cultures. Specifically, psychometric properties are often taken to refer to correlations of the items of an instrument (*instrument* is used in this chapter for any measurement device such as tests, questionnaires, and observational scales) or to correlations of an instrument with external measures. Multidimensional scaling, factor analysis, and the analysis of covariance structures (structural equations) are commonly employed to study structural equivalence. Thus, if equal factor structures are obtained in various cultural groups, it can be concluded that the psychological constructs underlying the instrument are identical. However, structural equivalence does not imply that both the origin and the measurement unit of the instrument are identical. Structural equivalence is primarily based on similarity in correlations across a variety of cultures and correlations are not affected by linear transformations of the variables. For example, if the scores of all persons in one cultural group are multiplied by a positive constant, the correlations remain unaffected and the factor loadings will also remain the same. Therefore, similar factor loadings can arise from scales with different origin and measurement units.

The second and third types of equivalence are concerned with measurement equivalence. When the scores of two cultural groups are compared, it is possible that the unit of measurement is identical, but that the scales do not have a common origin. This will be called *measurement unit equivalence*. Temperature scales in degrees of Celsius and Kelvin show this kind of equivalence. It has been argued that some intelligence tests can be validly applied within but not across cultural groups due to different origins of the scale in the cultural groups. In the case of measurement unit equivalence, differences between two scores (e.g., the scores between two classmates or the scores of an individual at two measurement occasions) can be compared both within and across cultures, while the scores themselves can only be compared within cultures.

If it can be ascertained that scores show not only an identical unit of measurement, but also a common origin, *scalar equivalence* or *full score comparability* is said to have been obtained. Scalar equivalence allows the comparison of the scores obtained, both within and across cultural groups. Examples are such variables as weight and height. For psychological measurements it is often difficult to establish scalar equivalence. In general it is easier to disprove than to prove scalar equivalence.

In the cross-cultural literature, the term *metric equivalence* has often been introduced to refer to the case when two or more data sets from different cultures exhibit similar psychometric properties (Berry, 1969). Within this framework, *subsystem validation* refers to the case when independent and dependent variables show the same relationship within cultures and across cultures (e.g., Roberts & Sutton-Smith, 1962). *Scalar equivalence* refers to the case in which scores from different cultures have a similar origin and unit of measurement (e.g., Poortinga, 1971).

We find some of this terminology imprecise. The term *metric* in metric equivalence denotes the unit of measurement in common usage in the psychometric literature, and does not denote structural equivalence nor a common origin of the scores, both of which are implied in the current usage of the term. Thus, we propose that this term be abandoned.

The first subtype of metric equivalence, subsystem validation, is actually a special case of structural equivalence, and can be subsumed under structural equivalence. The second subtype, scalar equivalence, is defined in the same way as in our scheme, and should be retained.

Methods

Sampling of Cultures

The selection of cultures in a cross-cultural study is often central to its scope for evaluating the hypotheses proposed. Three types of sampling procedures for the selection of cultures are commonly found in the literature. First, *convenience sampling* is often adopted in cross-cultural studies. Researchers select a culture simply because they may be from that culture, are acquainted with collaborators from that culture, or happen to be spending a sabbatical leave in that culture. The choice of culture is haphazard, driven by convenience, and not related to the theoretical questions raised. Very often, these studies adopt a "let's look and see" approach and do not develop any *a priori* predictions about cultural differences. When cultural differences are found, post hoc explanations are often developed to explain the differences.

The second approach is *systematic sampling*, in which cultures are selected in a systematic, theory-guided fashion. Usually, cultures are selected because they represent different values on a theoretical continuum. The classic study by Berry (1967) provides an excellent example of this approach. Two groups were studied,

one agricultural and one hunting. It was hypothesized that agricultural societies impose stronger pressure on conformity, and hence will lead to field dependence. Hunting societies encourage their members to be autonomous and hence are conducive to field independence. These two groups were selected systematically to evaluate this hypothesis. Another example of this approach is provided by Leung, Au, Fernandez-Dols, and Iwawaki (1992). In their study, four cultures were selected, namely, Spain, Japan, Canada, and the Netherlands. Japan and Spain tend to collectivistic, whereas Canada and the Netherlands tend to be individualistic (Hofstede, 1980). The comparison of these two groups will reveal the impact of individualism–collectivism. On the other hand, Spain and the Netherlands tend to be feminine, whereas Japan and Canada tend to be masculine (Hofstede, 1980). The comparison of Spain and the Netherlands with Japan and Canada will reveal the effects of cultural masculinity and femininity. An interesting feature of this study is that in both types of comparison, each group is composed of a Western and an Eastern culture. If differences are found between the two groups, the possibility that the differences are due to East–West differences can be ruled out.

We believe that in the systematic approach, bicultural comparisons are adequate only if there is a compelling theoretical framework in which the results can be interpreted, as is the case in Berry's (1967) study. When a study is exploratory, or when the theoretical framework guiding the study is rudimentary, the number of cultures in a study should be preferably larger than two. Campbell (1986) argued that the number of rival explanations is greatly reduced when the number of cultures involved in evaluating a hypothesis increases (cf. Leung et al.'s, 1992, study mentioned above).

In order to maximize the effectiveness of the systematic approach, cultures that are far apart on the theoretical dimension upon which they vary should be selected. This approach will maximize the chance to detect cultural differences. However, if only two cultures are selected that are highly dissimilar, they are likely to vary in other dimensions as well, and numerous alternative interpretations have to be ruled out. The problem does not arise when more than two cultures are studied; the larger the number of cultures selected, the fewer the alternative interpretations will be possible.

The third approach is *random sampling*. In this approach, a large number of cultures are randomly sampled, usually for evaluating a universal structure or a pan-cultural theory. Truly random samples are basically nonexistent in the literature, as no one has the resources to select a large number of cultures on a random basis for a single study. However, several studies have tried to follow this approach, and their sample may eventually begin to approximate a random sample (usually not of all groups but of all literate groups). For instance, Schwartz (1992, 1994) has sampled 36 cultures to evaluate the structure of human values. He basically included any cultural group in which he could find a collaborator to participate in the project. Buss et al. (1990) also followed a similar approach in sampling 37 cultures in their study of mate selection. Peterson et al. (1995) have surveyed managers from more than 20 countries on event management issues.

Sampling of Subjects

In order to make valid cross-cultural comparisons, the subjects from different cultural groups must be similar in terms of relevant background characteristics. Otherwise, it is hard to conclude whether the cultural differences observed are due to cultural differences or sample-specific differences. If we compare a group of illiterate subjects from one culture to a group of highly educated subjects from another culture, the differences observed are likely to be explainable in terms of educational differences rather than differences in some other aspect of their cultures. One approach to overcome this problem is to match the samples in terms of demographic characteristics so that sample differences can be ruled out as alternative explanations for observed cultural differences. For instance, college students from different cultures are often compared, and it is usually assumed that college students from different cultures are similar in their demographical characteristics. In a similar vein, Hofstede (1980, 1983) reduced the influence of unwanted intergroup differences by studying subjects from a single multinational organization from 53 countries. Schwartz (1992, 1994) sampled secondary school teachers from various countries to maximize the comparability of his subjects.

It is sometimes impossible to match samples from different cultures because of practical reasons, or because there are sharp cross-cultural differences in the demographic background of subjects. An adequate approach is then to measure the major demographic variables and treat them as covariates in the subsequent data analysis. For instance, in a study comparing the delinquent behaviors of adolescents in the United States, Australia, and Hong Kong, it was found that there were substantial differences in the father's educational standing in the three cultures (Feldman, Rosenthal, Mont-Reynaud, Leung, & Lau, 1991). The educational standing of the fathers of the Hong Kong subjects was significantly lower than that of the fathers of the Australian and American subjects. To overcome this problem, an analysis of covariance was used to compare cultural means partialling out the influence of father's educational standing.

It is unfortunate that many cross-cultural studies tend to ignore sample differences and fail to assess the impact of such differences. As the results are confounded by sample differences, it is difficult to provide an unambiguous interpretation.

Procedure

In this section we will review issues related to the procedural aspects of a cross-cultural study: the selection and evaluation of the adequacy of a measurement instrument, its translation, and its administration.

In an early stage of a project the question has to be raised whether the same instrument can be applied in all cultural groups. In the case of an already existing measurement instrument, its appropriateness in an intercultural context has to be judged. This amounts to answering the question whether the operationalizations

chosen in the instrument will be adequate in all cultural groups studied. Are the measurement operations specified in the instrument an adequate representation of the psychological domain that is to be covered? Embretson (1983) has introduced the concept of construct representation. The concept refers to the coverage of the psychological domain. Do the measurement operations specified in the instrument represent an adequate and sufficient sample of the behavioral manifestations of the psychological construct that is measured by the instrument? Any answer to this question requires knowledge of the cultural context in which the instrument will be applied.

The outcome of the decision process can take three forms: to *apply* the instrument, to *adapt* it, or to *assemble* a new version. In the first alternative the instrument or a translated version will be used without any modification. If the construct is not fully covered in the new group, the instrument can be adapted by rephrasing, adding, or replacing items that measure the missing aspects. If the researcher finds the original instrument entirely inadequate, a new instrument has to be assembled.

The decision whether to apply or adapt an existing instrument or to assemble a new one has both theoretical and practical implications. We propose to make the application of the same instrument the default choice. The advantages of this choice are (1.) the possibility to compare research results with other results reported in the literature, (2.) the possibility to maintain scalar equivalence (which is not achievable if results of newly assembled instruments are compared), and (3.) the small amount of money and effort that is required to administer an existing instrument as compared to the development and establishment of the psychometric properties of a new or adapted instrument. However, the direct application of an existing instrument may not always be the best choice. If an instrument does not cover important aspects of the psychological construct under study or if it shows a clear ethnocentric bias, adaptation or the assemblage of a new instrument would be a better choice. The decision may be seen as involving a cost-benefit analysis, with time and money as the costs and construct representation as the benefit.

There are numerous examples of *application* in the literature. For instance, Hofstede's (1980, 1983) classic study involves a value questionnaire that was administered in over 10 languages in 53 countries. The use of the Minnesota Multiphasic Personality Inventory (MMPI) in China provides a good example to illustrate the process of *adaptation*. When the items of the MMPI were tested in China, it was found that some items were meaningless in the Chinese context, and these items had to be modified (Cheung, 1989). However, most of the original items in the MMPI were retained, and it was actually possible to interpret the Chinese results in light of the American norms. The case of *assembling* a new instrument is rare in the literature, but two examples can be cited. Church (1987) argued that Western personality instruments are unable to capture many of the indigenous personality constructs of the Filipino culture. In light of these difficulties, he proposed a number of directions for the construction of a new personality instrument for the Filipino culture. In a similar vein, Cheung et al. (1996) have

argued that adaptation of Western personality instruments is inadequate in capturing all the major dimensions of personality in the Chinese culture. They started from scratch and created a personality instrument, called the Chinese Personality Assessment Inventory (CPAI), for the Chinese people. This instrument contains several indigenous personality dimensions, such as "face" and "harmony," as well as many items that are particularly meaningful in the Chinese context.

Instrument Translation

In the case of the *application* and the *adaptation* the instrument has to be translated. The translation-backtranslation method is probably the best known method for instrument translations (e.g., Brislin, 1980; Hambleton, 1993, 1994). An instrument is translated from one language to another and then backtranslated to the original language by an independent translator. This method often provides adequate results, but sometimes it produces a stilted language that reproduces the original language version well, but is not easily readable and comprehensible. This is particularly the case when test items contain local idioms that, almost by definition, are difficult to translate. Backtranslations can provide researchers who lack proficiency in the target language control of the adequacy of the translation. However, it is noteworthy that in the field of professional translations the procedure is almost never utilized (Wilss, 1982). Professional translations are commonly produced and checked by teams of competent bilinguals; hence, instead of relying on backtranslations, these teams utilize judgmental methods to assess the accuracy of the translation.

Werner and Campbell (1970) have proposed to decenter instruments that are used in a cross-cultural context—to adjust both the original and the translated versions simultaneously. The aim in decentering is not the verbatim reproduction of the original text but the enhancement of the naturalness and readability of the original and translated version.

Brislin, Lonner, and Thorndike (1973) have generated a useful set of guidelines to ensure good translatability (cf. Brislin, 1980, p. 432):

1. Use short, simple sentences in order to minimize the cognitive load of the instrument; a simple item-per-item check whether the phrasing can be simplified can lead to considerable improvement in translatability.
2. Employ the active rather than the passive voice.
3. Repeat nouns instead of using pronouns (which in some languages may be difficult to translate).
4. Do not use metaphors and colloquialisms, which are usually not well translatable.
5. Avoid the subjunctive mood (e.g., verb forms with "could" and "would").
6. Add sentences when key concepts are communicated. Reword these phrases to provide redundancy.
7. Avoid adverbs and prepositions telling "where" and "when," such as beyond and upper.

8. Avoid possessive forms where possible.
9. Use specific words, such as chickens and pigs, rather than general terms, such as livestock.
10. Avoid words indicating vagueness, such as probably and frequently.
11. Use wording familiar to translators where possible.
12. Avoid sentences with two different verbs that suggest different actions.

Various techniques have been proposed to check the accuracy of translations. An overview has been presented by Hambleton (1993, 1994). A distinction can be made between judgmental and empirical methods. Judgmental evidence of translation equivalence usually amounts to the application of a translation-backtranslation design. An assessment of the accuracy of the translation by a set of competent bilinguals is an alternative way to assess accuracy. Hambleton proposes three designs to study the accuracy of translations: (1.) bilinguals take the source and target versions of the test; (2.) source language monolinguals take the original and backtranslated versions, and (3.) monolinguals in both languages take the test. The latter is by far the most frequently applied design. Various psychometric techniques are available to evaluate the equivalence of the items in the source and target languages. These are known as item bias or *differential item functioning* techniques and will be discussed later.

Administration

Four areas will be distinguished in the following overview of issues related to a proper administration of instruments in a cross-cultural study (cf. Van de Vijver & Poortinga, 1991, 1992): the personal characteristics of the tester (or interviewer), interactions between the tester and the examinees, response procedures, and the stimuli of the instrument. In general, it will be difficult or even impossible to generate an exhaustive list of the problems that may arise in the administrative aspects of cross-cultural research. However, an overview of the common problems may sensitize the reader to the kinds of problems that can be encountered.

The presence of a tester, experimenter, or interviewer can be a threat to the validity of the results, particularly when this person has a different cultural background from the subjects in the sample. The potential influence has been recognized in observational studies of mother-child interactions (Super, 1981). In intelligence testing, the influence of racial differences between the tester and the examinee has been studied systematically (Jensen, 1980). Overall, the influence tends to be small, though the results are not consistent. In many cross-cultural studies the cultural distance between the tester or interviewer and the subjects will be considerably larger than in the American studies reviewed by Jensen. No systematic study has been undertaken of tester effects in settings more representative of cross-cultural settings.

A second area to be considered is the interaction between the tester and the respondent. In many research designs there is verbal communication between the two, and various problems may occur as a result of such communication. In

some cases the choice of the language used may be problematic. For instance, when Reuning and Wortley (1973) administered a variety of cognitive tests to the Bushmen, Kalahari desert dwellers, they faced the problem that their subjects had a highly heterogeneous linguistic background. Because it would have been difficult to hire and train an interpreter for each vernacular, they chose to minimize the verbal exchange in the testing procedure.

The reduction of verbal communication is not always possible because verbal exchange is essential in surveys and psychological testing. If the researcher decides to administer the instruments with the help of one or more interpreters, the potential influence of the interpreters should be evaluated, even when they are carefully trained. An assessment of the interpreter's influence usually requires that a group of respondents be interviewed by two interpreters. The results obtained by these interviewers are then compared with the help of an index of agreement. The choice of this index depends, among other things, on the nature of the data gathered. Cohen's kappa or its weighted version can be used in the case of nominal or ordinal data (Cicchetti, Showalter, & McCarthy, 1990; Cohen, 1960), and an intraclass correlation (Shrout & Fleiss, 1979) or Cronbach's alpha (e.g., Winer, 1971) in the case of interval data.

The third area involves response procedures. Subjects may be unfamiliar with a certain response procedure. For instance, the Porteus' Maze Test, a paper-and-pencil test, has been administered to groups of subjects who had never used a pencil before. Not surprisingly, their scores were very low (cf. Van de Vijver & Poortinga, 1991). If subjects are unfamiliar with a response procedure, it is important to reserve time for familiarizing the subjects with the procedure as part of the test introduction. In the area of personality and social psychology, Likert scales are often applied. Particularly among groups having little experience with this response format, the use of verbal descriptions of the response alternatives instead of numbers might be preferred.

A good example of the impact of response procedures can be found in the work of Serpell (1979). He administered a pattern-copying task to children in the United Kingdom and Zambia. The children's copying skills were assessed using two response media: pencil-drawing and iron-wire modelling, a popular pastime among Zambian boys. It was found that the British children scored higher than the Zambian children on the pencil-drawing task while the Zambian children reached higher scores on the iron-wire modelling task.

In some cases no empirical evidence may be available to judge the accuracy of a response procedure. A pilot study could then be carried out in which potentially useful response procedures are compared in a monotrait-multimethod matrix, in which several response procedures for measuring the same construct are examined. The correspondence of the results across the response procedures indicates the validity of the procedures.

Stimulus-related aspects are by far the most extensively studied area of procedural problems in cross-cultural research. Stimulus familiarity is the most often mentioned source of invalid intergroup score differences in the literature (e.g., Irvine & Carroll, 1980). A study by Deregowski and Serpell (1971) illustrates the

importance of stimulus familiarity. Scottish and Zambian children were asked to sort miniature models of animals and motor vehicles in one experimental condition and their photographs in another one. No intergroup differences were found for the actual models whereas in the sorting of photographs, the Scottish children obtained higher scores than the Zambian children.

In the past, various attempts have been made to adapt the stimuli of cognitive tests in such a way that intergroup differences caused by stimulus familiarity would be eliminated. Both the culture-free and culture-fair test movements were intended to serve this purpose. Even though the original ideas of the movements have been long abandoned and it is widely acknowledged that such tests cannot be constructed (Frijda & Jahoda, 1966), the concern for stimulus familiarity is still widely shared. Stimuli differ in terms of their cultural entrenchment. Simple geometrical stimuli such as squares, circles, and triangles are often used as stimuli in cognitive tests because their cultural loading is assumed to be limited though certainly not absent.

In the area of personality and social psychology, stimulus familiarity also plays an important role. Items of personality scales frequently use complex words or expressions. Effort should be made to use simple, unambiguous stimuli and to avoid the undesirable introduction of verbal abilities, such as vocabulary and text comprehension skills, as sources of individual differences.

Design

A distinction will be made between the design of structure-oriented and level-oriented studies in cross-cultural psychology. Structure-oriented studies examine relationships among variables and attempt to identify similarities and differences in these relationships across cultures. For example, is the structure of intelligence universal? Level-oriented studies, on the other hand, focus on differences in the magnitude of variables across cultures. For example, are members of culture A more individualistic than members of culture B?

The design of structure-oriented studies is often straightforward: it replicates the design of the original study. The design of level-oriented studies tends to be more complicated, and an adequate choice of research variables and design is needed to enhance the interpretability of the findings obtained. There is at least one important issue common to all level studies: Which covariates should be included? It was argued before that the Neyman–Pearson framework assumes a random assignment of individuals to treatments and that cross-cultural studies can never adopt a truly experimental design. Cultural groups differ in many respects, only some of which are of interest in a particular study. All these group differences can in principle explain observed score differences. An important aid in the reduction of the number of rival explanations are covariates. Covariates can be helpful in the interpretation of cross-cultural score differences in two ways. First, they can be used to validate the interpretation of the cross-cultural differences as hypothesized by the experimenter. For instance, if individualism–collectivism is assumed to be related to a psychological phenomenon, say inter-

group hostility, individuals from individualistic and collectivistic countries could be included in the study. In addition to an intergroup hostility measure, a test of individualism–collectivism should be administered to all individuals. These scores could then be used in an analysis of covariance, in which cultural groups are the independent variable, the hostility measure the dependent measure, and the individualism–collectivism score the covariate. The covariate is used to validate the cross-cultural differences postulated by the theory. Earley (1989) has evaluated the effect of individualism–collectivism on social loafing with this approach.

Second, covariates can also be used to check the effects of nuisance variables. The inclusion of such covariates will control for cultural differences that influence the behavior in question, but that are not specified by the theory. For instance, if men and women differ in the level of hostility and if the student groups in the two cultures in the previous example have a different male–female ratio, gender could be used as a covariate, because the observed cross-cultural differences could be due to the difference of gender composition of the two groups as well as to cross-cultural differences in intergroup hostility. The covariate is not meant here to provide an explanation of the cross-cultural differences, but to control for nuisance variables. Covariance analysis as discussed in textbooks is almost always exclusively concerned with the elimination of the impact of nuisance variables. The conclusions of an analysis of covariance can be misleading if the assumption of parallel regression lines within each cultural group is violated (cf. Lord, 1967). A simple statistical test of the equality of regression coefficients in two cultural groups is described in Cohen and Cohen (1983: chapters 10 and 12) and Pedhazur (1982, chapter 12).

Covariates can be based on aggregate rather than individual measures as the previous examples could suggest. In a study of intergroup differences in some cognitive test, educational quality could be assessed. Such a measure located at the class or even cultural level can be used as a covariate at the individual level, meaning that all subjects of a class or school will get the same score on the variable.

We strongly encourage the use of covariates because they provide an effective way to confirm a particular interpretation of intergroup differences and to falsify alternative interpretations. Yet, the limitations of methodological and statistical procedures should be acknowledged. Statistical techniques can help to evaluate the impact of contextual variables, but will not provide information on which covariates to choose. For example, intergroup differences in cognitive test performance might be assumed to be related to educational quality or to Westernization, to mention a few possibilities. Methodological and statistical considerations cannot dictate the choice. All that can be asked from methodology and statistics is a set of tools to enable the evaluation of the accuracy of the choice, or, in case both sets of variables have been measured, the evaluation of their relative importance.

Leung and Zhang (1995) have concluded that many studies have been exported from the West to non-Western countries, and some of the issues examined in these studies are of little relevance to the local culture. It is entirely possible that results obtained in many of these studies are shaped by the cultural back-

ground of the researchers, and that different results may be obtained if a different cultural vantage point is taken in the design of these studies. Two approaches may be adopted to design a culturally balanced study, in which no single culture will dominate the research questions explored and bias the results obtained. First, a *decentered* approach can be adopted, in which a culturally diverse perspective is taken in the conceptualization and design of a study. For instance, when Schwartz (1992) tested his pan-cultural model of value structure, he encouraged researchers from different cultures to add culture-specific value items to his pan-cultural set. Smith and Peterson (1988) have taken into account the influence of culture in their formulation of a theory of leadership behavior and their empirical test of the theory (Peterson et al., 1995).

The second approach is the *convergence* approach. The basic idea is to design a study that is as culturally distant as possible from existing studies and to see if the results obtained overlap with existing results. If the new results overlap with existing results, it can be concluded that the cultural origin of existing studies have not biased the results obtained. If different results are obtained, however, the possibility that the cultural origin of existing studies has biased the results must be further investigated. The best examples to illustrate this approach are provided by Bond and his colleagues. The Chinese Culture Connection (1987) designed a value survey based entirely on Chinese values and administered it in 22 countries. It was found that three factors showed overlap with factors identified by Hofstede (1980), whose results were based on a Western instrument. A new factor emerged, termed Confucian work dynamism, which correlated highly with economic growth. In the realm of person perception, Yang and Bond (1990) administered a set of emic Chinese descriptors together with a set of imported American descriptors to a group of Taiwanese subjects. Of the five Chinese factors identified, only four were adequately explained by the American factors, and one factor was uniquely Chinese.

Data Analysis

In this section we will first describe bias, followed by a description of psychometric techniques to detect differential item functioning as a special case of bias. In the last part of the section we will describe the most common statistical techniques for analyzing cross-cultural data sets.

Preliminary Analyses

Prior to the data analysis that addresses the central research question or hypothesis, preliminary analyses will often be required. If a psychological instrument is used, its psychometric properties should be established, in particular its reliability. In most cross-cultural studies this seems to be routine practice. It is surprising that tests of intergroup differences in reliability are almost never carried out even though the observation of dissimilar reliability coefficients can provide valuable

clues about measurement accuracy and hence, the appropriateness of an instrument for cross-cultural comparison. Procedures to test the equality of independent alpha coefficients have been described by Kraemer (1981) and Hakstian and Whalen (1976).

The interpretation of intergroup differences can be seen as an attribution process. Two kinds of attributions can be envisaged. Observed intergroup differences may be valid, and members of group A have on average more of a particular propensity such as anxiety, intelligence, or collectivism than members of group B. The observed differences may also be due to bias (measurement problems). For instance, the items used may be affected by intergroup differences in stimulus familiarity or social desirability, which have produced the cultural differences observed.

A distinction can be made between three types of bias. The first is called *construct bias*. This kind of bias occurs when the psychological construct is not identical across cultural groups. Construct bias implies that the theoretical construct is not or is inadequately represented in the instrument. In Embretson's (1983) terms, construct bias refers to a poor construct representation. An example can be found in the area of intelligence. Everyday conceptions of intelligence, mainly in non-Western cultures, have been found to differ from the conception underlying intelligence tests (Serpell, 1993; Sternberg, 1985; Super, 1983). Everyday conceptions of intelligence tend to be broader than scientific theories. In addition to reasoning and factual knowledge that are shared in both conceptions, "social intelligence" is also included in everyday conceptions. "Social intelligence" involves social skills, obedience, and knowing one's role in the family, class, and peer group. A Western intelligence test will therefore show construct bias in many non-Western contexts. Culture-bound syndromes, such as amok, that are studied in ethnopsychiatry provide another example (Draguns, 1989; Harkness & Super, 1990). In the area of personality the Chinese concept of "filial piety" can be mentioned; filial piety refers to taking care of one's parents, conforming to their requests, and treating them well. The Chinese concept is much broader than the Western concept of being a good son or daughter (Ho, in press). A direct comparison of these two will result in construct bias.

It was argued before that a cross-cultural researcher may choose to apply or adapt an existing instrument, or assemble a new one. In the terminology of this section, the decision should be based on whether construct bias is present in the instrument. The assessment of construct bias should be based on knowledge about the cultural groups. If an instrument has been applied in several cultural groups with the same instrument and no additional data are available, statistical tests alone will not lead to a full understanding of the nature of the construct bias present. A proper assessment of construct bias should be based on research conducted in each cultural group, exploring whether the implicit definitions of the concept of the test are consistent across the cultural groups. Examples of this approach can be found in the work of Serpell (1993), Sternberg (1985), and Super (1983).

The second kind of bias is called *method bias*. If method bias occurs, the psychological construct is well represented by the instrument but the assessment

procedure introduces unwanted intergroup differences. Empirical studies that reveal method bias are Deregowski and Serpell's (1971) sorting task of miniature models and pictures of animals and motor vehicles, described earlier and Serpell's (1979) study of pattern copying using a paper-and-pencil format and iron-wire models.

Method bias can be examined by monotrait-multimethod matrices or triangulation. In this approach, a psychological construct is investigated using a systematic variation of methods. If the cross-cultural differences observed are similar across methods, method bias is unlikely. Method bias is said to occur if the intergroup differences vary across the methods. An analysis of covariance structures is often used in this situation, as will be illustrated later on.

A specific way to study method bias involves the repeated administration of the same instrument. Test-retest studies of cognitive tests have often shown score increases that are larger in non-Western groups than in Western groups (Kendall, Verster, & Von Mollendorf, 1988; Van de Vijver, Daal, & Van Zonneveld, 1986). A significant improvement in one group at the second occasion, or a gain pattern that is differential across groups, undermines the validity of the first test administration.

The third kind of bias is the most investigated. It was originally called *item bias* and is now better known as *differential item functioning*. Whereas construct bias and method bias involve the appropriateness of the whole instrument, differential item functioning occurs at the item level. Item bias refers to anomalies in the instrument at the item level caused by poor translation or inappropriate items in a particular context. A widely accepted definition of differential functioning has been proposed in the area of ability testing. An item is said to show item bias if persons from different cultural groups with an equal ability do not have the same probability of giving a correct answer. Individuals with an equal ability or attitude from different cultural groups should, apart from chance fluctuations, show the same average score for items of an unbiased instrument. From a psychometric point of view, the assessment of this kind of bias is best developed. A multitude of psychometric techniques have been proposed to test the presence of item bias. We will not describe them in detail. Rather, we shall briefly describe and illustrate two of them, followed by the presentation of a taxonomy of the techniques.

Historically speaking, analysis of variance was probably the first technique that has been applied to study differential item functioning (Cleary & Hilton, 1968). We shall discuss here a slightly modified procedure. Suppose that a test for authoritarianism of 30 five-point Likert-scale items has been administered in two cultural groups of 200 persons each. If we are interested in the presence of differential item functioning, the first step is to divide the subjects into score level groups. Individuals with an equal score are assumed to have an equal level of authoritarianism, and subjects with the same score are grouped together. Because the scores on the Likert scale range from one to five, the total score can vary from a minimum of 30 to a maximum of 150. The split of the score distribution into score levels should be based on the score of all cultural groups together; the same

cutoff scores should be applied to all cultural groups. Theoretically speaking, there can be 121 score level groups in this case (from 30 to 150, including both ends). In practice, a much smaller number will be used as the number of subjects will be unevenly distributed across the score levels (Clauser, Mazor, & Hambleton, 1994). Quite often, an attempt is made to choose the cutoff scores in such a way that the number of subjects in each group is approximately the same. Score level will be one of the independent variables in our data analysis; the other one will be the cultural group. Differential item functioning is tested in a set of analyses of variance, one per item, with culture and score level as independent variables and the item score as dependent variable.

Following Mellenbergh (1982), we shall make a distinction between two types of item bias: uniform and nonuniform. Figure 7-1 presents the curves which depict the average score of two groups on a particular item, technically called empirical item characteristic curves (Allen & Yen, 1979). When the curves more or less coincide, there is no bias (Figure 7-1a). When the curves are more or less parallel without coinciding, there is uniform bias (Figure 7-1b). When the curves are not parallel, the items are said to show a nonuniform bias (Figure 7-1c). In this case, the difference in the average test score will depend on the score level. For instance, for low authoritarian subjects, the item is endorsed more strongly in one culture, while for high authoritarian subjects, the item is endorsed more strongly in the other culture. A combination of both types of bias is presented in Figure 7-1d. In terms of the analysis of variance, an item is said to be uniformly biased when the main effect of culture is significant. In this case subjects from one cultural group have a consistently higher score than individuals with the same underlying propensity from another cultural group. A significant interaction of level and culture indicates the presence of nonuniform bias.

Item bias analyses can be carried out in an iterative or a noniterative way. In the latter case the analyses of variance are carried out for all items and the presumably biased items (i.e., all items with a significant main effect for culture and/or a significant interaction between culture and level) are removed simultaneously. Intergroup score comparisons are carried out on the reduced item set. In an iterative procedure the elimination proceeds on an item-by-item basis. In the first step, all items are considered. The item with the largest bias component (i.e., the smallest probability in the computer output) is then removed if the component is significant. The whole procedure is then repeated for the reduced set of items until no more bias components are significant. An attractive feature of iterative procedures is that the total score is updated in each iterative step, which allows for a finer detection of bias. It might well be that after the removal of a few items the meaning of the total score changes somewhat and this change can result in the removal of different items than in the case of a noniterative procedure. However, iterative procedures are cumbersome because after the removal of an item new cutoff scores for the score levels have to be calculated.

The removal of biased items does not inevitably lead to the elimination of intergroup differences in the average scores (Poortinga & Van der Flier, 1988). Items can be biased or unbiased, irrespective of the presence (or absence) of inter-

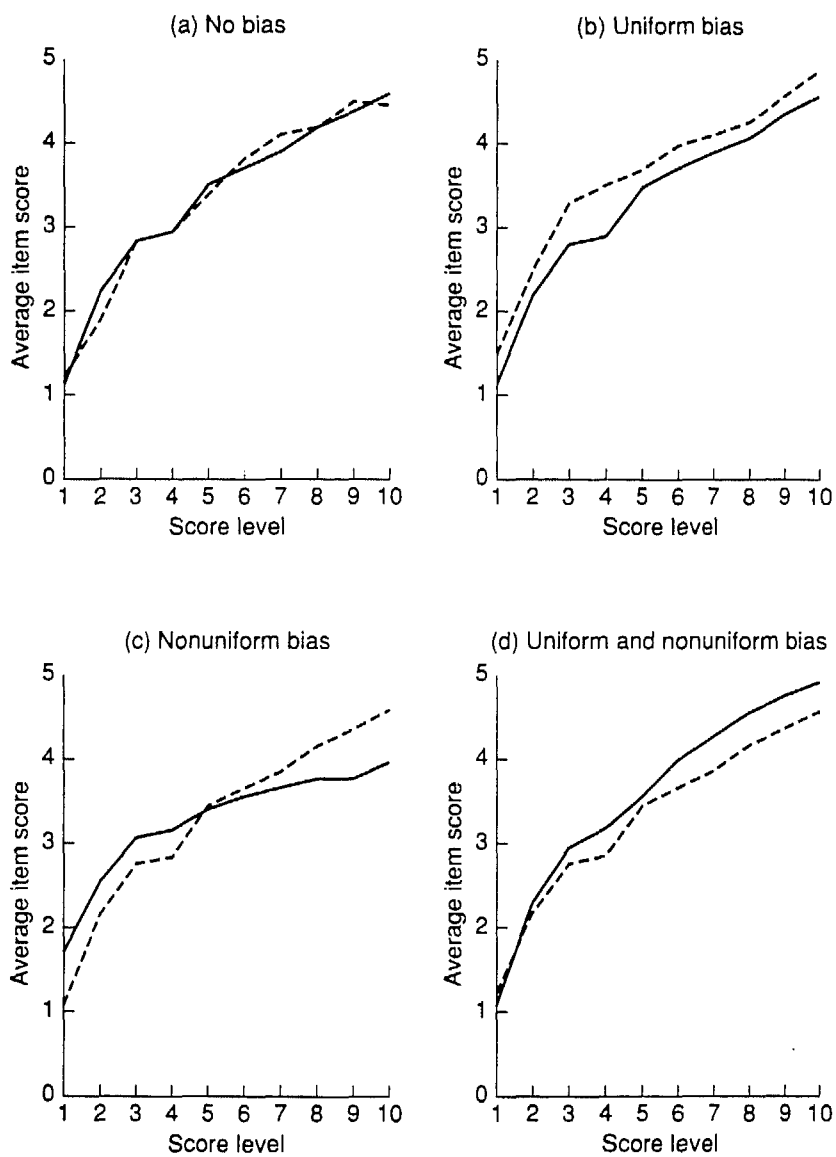


FIGURE 7-1 The average performance of the two cultural groups on an item that shows (a) no bias (b) uniform bias (c) non-uniform bias, and (d) both uniform and nonuniform bias (hypothetical example).

group differences. After all, item bias analysis does not test whether there are overall intergroup differences in total score, that is, whether individuals of group A would have a higher propensity on X than individuals of group B. Rather, item bias analyses test whether there are intergroup differences per score level, that is, whether individuals from group A with a particular attitude level have the same average score on a particular item as people from group B with the same attitude level. The item bias analysis uses an analysis of variance with level and cultural group as independent variables and a particular item score as dependent variable. In contrast, an analysis of variance testing the presence of overall intergroup differences treats culture as the independent variable and the item or total test score as the dependent variable.

The most popular technique to test differential item functioning today is the Mantel–Haenszel statistic (Holland & Thayer, 1988). The statistic is closely related to item response theory (e.g., Hambleton & Swaminathan, 1985). More specifically, the Mantel–Haenszel statistic tests whether a single Rasch model, a model from item response theory, fits the data in each group. The rationale behind the Mantel–Haenszel statistic and the analysis of variance approach explained earlier is similar. The major difference is that the Mantel–Haenszel procedure works with dichotomous data whereas the analysis of variance is based on interval data.

Item response theory represents a more general approach for assessing differential item functioning (e.g., Hambleton & Swaminathan, 1985; Hulin, 1987). This model assumes that an unbiased item evokes a similar response from respondents that are similar in their standing on a latent trait regardless of their cultural backgrounds. In the general form, this model links item responses to latent traits by means of a logistic curve specified by three parameters. The first parameter is concerned with the discrimination capability of the item; the second parameter is concerned with the difficulty level of the item; and the third is concerned with the extent to which guessing is involved in responding to the item. In specific applications, two-parameter models, which exclude the guessing parameter, are often employed for modelling attitudinal data. To detect biased items, item characteristic curves, which relate the probability of making a certain response to standing on a latent trait, are examined. Items are equivalent across two cultures if their item characteristic curves are similar across these cultures. Differential item functioning is present when parameters differ significantly across cultural groups. Item response theory has been applied in cross-cultural research on self-concept (Leung & Drasgow, 1985), job satisfaction (Candell & Hulin, 1987), intelligence (Ellis, 1989; Van de Vijver, 1988), and attitudes toward mental health (Ellis & Kimmel, 1992).

The standard procedure for the application of item response theory is as follows:

1. Item response theory assumes that a scale is unidimensional, and the unidimensionality of the scale must be checked. If the scale is multidimensional, each unidimensional subscale must be examined separately.
2. An item response theory model with the appropriate number of parameters is selected to fit the data in each culture.

3. The parameters identified for each cultural group are equated on the same metric through an iterative linking procedure.
4. Biased items are detected and eliminated with the aid of item characteristic curves and a chi-square test. The parameters are equated again with the linking procedure with unbiased items only, and this procedure stops when no biased items are detected.
5. The biased items identified are eliminated from the scale before cross-cultural comparisons are made.

Item response theory has characteristics that make it appropriate for cross-cultural applications. First, the estimates of item parameters do not depend on the propensity level of the group studied. This is not the case in classical test theory in which the difficulty of an item, operationalized as item average, depends on the average ability level of the group. Similarly, the estimates of person parameters in item response theory are independent of the items of the instrument. Second, most models in item response theory allow for a fit test. The extent to which the empirical data can be taken to obey the theoretical model can be examined (e.g., Hambleton & Swaminathan, 1985; Lord, 1980; Van den Wollenberg, 1988).

The most important limitations of item response theory are twofold. The applicability of item response models may be reduced by the strict assumptions that have to be met, particularly in the Rasch model. Furthermore, large sample sizes are required to obtain stable estimates, particularly in the three-parameter model.

TABLE 7-1 Schematic overview of differential item functioning techniques (after Van de Vijver, 1994)

Sampling distribution	Model equation	
	Linear	Nonlinear
Unconditional procedures		
Unknown	Partial correlation index (Stricker, 1982)	Delta plots (Angoff, 1982)
Known	Analysis of variance (Cleary & Hilton, 1968)	
Conditional procedures		
Unknown	Standardized <i>p</i> -difference (Dorans & Kulick, 1986)	Item response theory (McCauley & Mendoza, 1985)
Known	Analysis of variance with score level as one of the variables	Mantel-Haenszel procedure (Holland & Thayer, 1988)

Three questions are relevant in the choice of a particular item bias statistic. First, what kind of measurement model should be used? Some techniques are based on a linear model such as an analysis of variance, while others, such as the Mantel-Haenszel statistic, are based on a nonlinear model. In general, interval-level data tend to be analyzed using linear models, while dichotomous data are often analyzed by item response theory, a nonlinear model. Second, is the technique conditional or unconditional? Most modern techniques are so-called conditional procedures. These techniques compare the scores of individuals across cultural groups per score level. Both of the previous examples are conditional. Until the eighties, unconditional procedures were more common, such as the comparison of item averages. It has been shown several times (e.g., Lord, 1977, 1980) that unconditional procedures can underestimate the number of biased items. Therefore, conditional procedures are to be preferred. The third question refers to the sampling distribution of the item bias statistic. In both our examples, the sampling distributions are known. This allows for a statistically rigorous test of the null hypothesis of no bias. Yet, various bias statistics that have been proposed have unknown sampling distributions, which makes a statistical evaluation of item bias questionable, whatever the intuitive appeal of the statistic (e.g., Stricker's, 1982, partial correlation index). A taxonomy of bias statistics on the basis of these three questions is presented in Table 7-1.

A perusal of the cross-cultural literature shows that differential item functioning techniques are infrequently applied in cross-cultural psychology. We find this disappointing; in many cases it should be standard practice to carry out an item bias analysis prior to the actual data analysis. Item bias techniques have been mostly applied to cognitive test scores, and much less so in the area of personality and social psychology. There is no good reason for the uneven distribution of the application of item bias techniques, unless one would want to maintain that items in personality questionnaires are of a much higher quality and less open to bias than are cognitive test items.

Two general findings emerge from the application of differential item functioning techniques in cross-cultural psychology. First, item bias may be psychometrically well defined and operationalized, but it may be difficult to grasp its psychological meaning. In current applications, it is not uncommon to find that item bias is reported but no sensible explanation can be provided for the bias (Scheuneman, 1987; Van de Vijver, 1994). Furthermore, item bias indices are not stable in cross-validation studies. Retests with the same instrument may show other items to be biased. The common difficulties encountered in empirical applications of item bias techniques, such as inadequate stability and interpretability, may reduce the attractiveness of these procedures. Still, if we start to routinely apply item bias techniques to cross-cultural data, we may build up a body of knowledge about item quality from a cross-cultural perspective.

Second, some item bias studies have shown a substantial proportion of items to be biased, sometimes more than half of the items. In such a case the item bias analysis seems to point to a serious lack of validity of the instrument. A prudent approach would then be to refrain from intergroup comparisons.

Establishing Scalar Equivalence

Techniques based on correlations such as factor analysis have been proposed and used to test scalar equivalence. For instance, Eysenck and his coworkers concluded that scalar equivalence can be assumed when the factor structures obtained with a measurement instrument in various cultural groups are similar. A similar argument has been put forward by Berry (1980). However, as argued before, similarity of correlations matrices or factor structures across cultural groups can only demonstrate structural equivalence, and does not speak to scalar equivalence. Structural equivalence imposes fewer restrictions on the data than scalar equivalence.

At least three approaches have been proposed to establish full score comparability in the literature. First, various authors assume but do not test full score comparability. If a test is administered in two cultural groups and the test scores are compared without any concern for comparability, full score comparability is implicitly assumed. An example comes from the literature on culture-free and culture-fair intelligence testing. Reports involving these somewhat obsolete instruments hardly involve statistical tests of full score comparability (e.g., Anastasi, 1976; Cattell, 1940; Cattell & Cattell, 1963). In our view, researchers should attempt to provide evidence for full score comparability of their instruments.

The second and third approaches are internal validation procedures. The procedures are called internal because the data used to validate equivalence are derived from the instrument itself. The second approach involves intra-cultural techniques in which empirical data are compared to theoretical expectations for each culture. It is possible to formulate hypotheses about the order of difficulty or endorsement rate of items in some instruments. For instance, items of tests of arithmetic abilities can often be ordered by the complexity of the arithmetical operation required. Operations requiring the manipulation of one-digit numbers will be easier than operations requiring two-digit numbers; additions and subtractions will be easier than multiplications and divisions. Strong evidence against scalar equivalence is obtained if theoretical expectations are not borne out. As a second example, the use of fit tests in applications of item response theory can be mentioned (e.g., Hambleton & Swaminathan, 1985; Lord, 1980; Van den Wollenberg, 1988). A good fit within each group provides initial evidence for scalar equivalence. Intracultural validation techniques provide necessary though insufficient evidence for the presence of scalar equivalence.

The third approach can be called *cross-cultural validation*. The best known example is the work on item bias, or differential item functioning (Berk, 1982; Holland & Wainer, 1993). Various psychometric techniques have been developed which scrutinize consequences of the lack of bias at the item level (cf. the description of item bias before).

A special case of the monotrait-multimethod approach, described earlier for the examination of method bias, is the use of multiple measures to capture the same construct. Triangulation, as this procedure is often called, can provide some insight in scalar equivalence, especially when the statistical techniques described in the previous section do not apply, such as in the case of single-item measures

(e.g., measures in Piagetian psychology, social behavior). Triangulation amounts to utilizing multiple measures, as diverse as possible, to measure the construct. If convergent results are obtained with different measures, bias is not likely to have produced the results. For instance, Hess, Chang, and McDevitt (1987) found that in comparison with American mothers, Chinese mothers were more likely to attribute the academic performance of their children to effort. Consistent with this result, Chinese children were also more likely to attribute their academic performance to their own effort than were American children. The convergence between the children and mothers has strengthened the validity of the cultural difference observed. In contrast, Serpell's (1979) study of Zambian and Scottish children's copying skills using iron-wire models and pencil-drawing is an example of nonconvergent operations. It should be pointed out that although multiple measures can assess the confounding influence of bias, it does not guarantee scalar equivalence even when convergence is obtained. The equality of the origin and the unit of the measurement scale is not directly assessed in triangulation.

Statistical Tests of Cross-Cultural Differences: Introduction

The statistical techniques described in the previous section examine the cross-cultural applicability of research instruments and the validity of the use of these instruments in cross-cultural comparisons. In this section, we will describe statistical tests that are applied after the adequacy of the psychometric characteristics and the absence of bias have been established. A distinction between structure- and level-oriented techniques will be made in our presentation. Because of space limitation, we will only provide a brief overview of the statistical techniques.

Prior to any statistical analyses, it should be decided whether the data need to be standardized, and if so, which standardization procedure is to be used (e.g., Hofstede, 1980; Leung & Bond, 1989). Culture-level analyses can yield strikingly dissimilar results for standardized and nonstandardized data sets. Standardization is usually defined as the computation of z scores ($z = (X - M)/S$, in which X is the score to be standardized, M is the mean and S is the standard deviation). Standardization is defined here more generally and refers both to z scores and to transformations to other deviance scores such as X/S and $X - M$. The aim of standardization is the reduction or elimination of unwanted intergroup differences such as those due to response sets. If scores are standardized per cultural group, intergroup differences in means, standard deviations, or both are eliminated. Such a procedure requires justification, because intergroup differences in average scores may not be exclusively due to response sets or other unwanted sources but may reflect valid differences. The justification is usually based on the presumed equality of averages across cultures. For instance, Schwartz (1992), who has transformed raw scores to deviations from the mean in his value survey, argues that the average importance score that people give to all the value items in his inventory should be similar across individuals, because his instrument represents a comprehensive set of human values. If such a reasoning cannot be justified, analyses based

on the original as well as the standardized data should be conducted and the results obtained compared.

Statistical Tests of Cross-Cultural Differences: Level-Oriented Techniques

The most frequently reported statistical tests of *level* differences are the *t* test and analysis of variance (e.g., Glass & Hopkins, 1984; Hays, 1994). The most commonly tested null hypothesis specifies that there are no intergroup differences. In a *t* test, the cultural group is the independent variable and the score on a psychological instrument is the dependent variable. The popularity of the *t* test, in cross-cultural psychology as well as elsewhere, is undoubtedly attributed to its simplicity, availability (in computer packages), and robustness against violations of assumptions. The same holds for the analysis of variance, which is carried out when data of more than two cultural groups are studied. The major interest tends to be in the main effect for culture, which, assuming that the effect is significant, indicates that at least one culture has an average on the dependent variable different from the other cultures. More complex designs, so-called factorial designs, are often reported in cross-cultural research. These are designs in which, in addition to culture, one or more independent variables, such as gender or age, are included. The inclusion of such additional variables, say gender, is particularly relevant when the reaction patterns of men and women are expected to differ across the cultures studied (e.g., the male-female differences on the dependent variable are more pronounced in one culture). These differences in reaction patterns will come out in an analysis of variance as a significant interaction of culture and gender.

Regression analysis is often used in level-oriented studies. Regression analysis evaluates the influence of one or more independent variables on a dependent variable in terms of the amount of variance in the dependent variable that the independent variables can explain. Regression coefficients express the degree of relationship between the independent and the dependent variables. The squared multiple correlation, another relevant statistic of the regression analysis, is the amount of variance explained by the independent variables, which gives an overall evaluation of the success of the independent variables in predicting variation in the dependent variable. In cross-cultural studies, level-oriented hypotheses involve a test of whether the intercept of a regression equation is similar across different cultural groups (Cohen & Cohen, 1983; Pedhazur, 1982).

Regression analysis can be carried out on raw or standardized scores (mean of zero and unit variance). Standardization affects the size of the coefficients, but leaves the significance level unaffected. In practice it has become more common to report standardized regression coefficients because they are independent of the measurement units of the independent variables.

The choice between an analysis of variance (or *t* test or *z* test) and a regression analysis mainly depends on the measurement level of the independent variables. Nominal and ordinal data are often analyzed in an analysis of variance,

whereas predictors based on interval data are usually analyzed with a regression model. Yet, the choice is more a matter of convenience than of principle, given the close link between the two. Cohen and Cohen (1983) and Pedhazur (1982) describe how an analysis of variance can be seen as a regression analysis; the independent variables of an analysis of variance are the predictors of a regression analysis. The significance tests of the regression coefficients in the regression analysis (which are *t* tests) yield similar results as the significance tests of the analysis of variance (*F* test). Specifically, the squared *t* values of the regression statistics are equal to the *F* ratios of the analysis of variance.

Regression analyses can be used to identify relationships that hold both within and across cultures, which are actually structure-oriented issues. This approach is illustrated with two variables (*X* and *Y*). The first step in this technique is to obtain a pan-cultural regression equation of *Y* on *X*, in which data from all cultures are included. In the second step, culture is added as a dummy variable, and another regression analysis including *X*, the dummy variable, and the interaction of *X* and the dummy variables as predictors, is carried out. The multiple correlations of the two equations are then tested for equality. If there is no significant difference between these two multiple correlations, it is concluded that the relationship holds within each culture and across all cultures. In other words, the regression weights and the intercepts of the equation are similar in all cultures. A significant difference of the two multiple correlations points to the presence of intergroup differences on the dependent variable not explained by *X*. For an elaboration of this approach, see Leung (1987) and Poortinga and Van de Vijver (1987).

When more than one dependent variable is involved, covariance structure analysis employing models such as LISREL and EQS becomes appropriate. The basic idea is to test whether the interrelationships of the variables are similar in different cultures. For instance, confirmatory factor analysis involving two or more cultures is frequently conducted to see if a factor structure is similar in different cultures (e.g., De Groot, Koot, & Verhulst, 1994; Leung et al., 1992; Marsh & Byrne, 1993; Sachs, 1992; Watkins, 1989). The reader is referred to Poon, Chan, Lee, and Leung (1993) for a method to compare the equivalence of a covariance structure in a large number of cultures.

Multilevel analysis is another procedure that can be used to address both level- and structure-oriented issues (Bock, 1989; Bryk & Raudenbush, 1992; Goldstein, 1987; Lee, 1990). Even though these procedures have not been applied in cross-cultural research, their potential value is obvious. At least two levels of analysis are possible in cross-cultural research (e.g., Leung, 1989; Leung & Bond, 1989). In the culture-level approach, culture is the unit of analysis, and the results obtained are characterizations of cultures, but not individuals. The classic study on values by Hofstede (1980) is based on this approach. There is no assumption with regard to whether relationships found across cultures will hold within each of the cultures included in the analysis. Culture-level analyses can guard against the ecological fallacy, the incorrect application of culture-level characteristics to individuals. When a culture is known to be individualistic, the mean score on a scale of individualism will be higher than in a collectivistic culture, but it does not mean

that each person has a high score on individualism. Furthermore, cross-level inferences can be fallacious because of a difference in meaning of constructs at individual and cultural levels. Gender at the individual level can have two values, male and female; an aggregated gender score at the cultural level refers to the proportion of males and females in a group, which is quite a different concept.

In the individual-level approach, the individual is the unit of analysis, and this is the dominant approach in cross-cultural psychology. The relationships between variables at individual and cultural levels need not be equal (cf. Ostroff, 1993). Yet, it is more elegant theoretically to demonstrate their equality. An example can be found in "subsystem validation," in which "hypotheses are examined both intraculturally and cross-culturally, so that explanatory variables may be tested at two levels" (Berry & Dasen, 1974, p. 19). The objective of this approach is to establish that the relationships among a set of variables hold within a culture as well as across cultures. For instance, in the classic study by Segall, Campbell, and Herskovits (1966), it was found that when a culture is associated with a more "carpentered" environment (more corners formed of intersecting planes perpendicular to each other), people from this culture are more susceptible to geometric illusions. This finding explains why one cultural group is more susceptible to geometric illusions than another cultural group. Their findings also imply that if a person is exposed to a more carpentered environment, he or she is more susceptible to geometric illusions. Thus, their findings are able to explain cultural differences as well as individual differences in susceptibility to geometric illusions.

Statistical Tests of Cross-Cultural Differences: Structure-Oriented Techniques

Cross-cultural psychologists are often interested in a comparison of the *structure* underlying the data rather than a direct comparison of the observed variables as discussed in the previous section. For instance, much research has been devoted to the question of whether the structure of intelligence is universal. Do the same cognitive processes contribute to test performance in different cultural groups? These questions have probably their intellectual roots in the notion of the psychic unity of humankind (Tylor, 1871), which can be interpreted as the idea that the structure behind human behavior is universal.

While multivariate statistical techniques are often applied in structure-oriented analysis, ANOVA or *t* tests are commonly applied when the independent variables are discrete. The focus is here to evaluate whether the differences of the dependent variable across the various levels of the independent variable are similar or different in each culture. For instance, Buss (1989) applied *t* tests to examine mate preferences of males and females in 33 cultures. He confirmed hypothesized gender differences in mate preferences in the cultures studied.

The similarity of psychological structure has been studied mostly by means of factor analysis (Harman, 1976; McDonald, 1985). Like regression analysis, factor analysis postulates that an observed score is a weighted sum of a usually limited set of contributing factors. Unlike in regression analysis, however, the

contributing factors are not observable in factor analysis. An observed score, for example, an intelligence test score of a person, is a weighted sum of unobservable factor scores, such as reasoning ability, perceptual speed, and memory, which in turn are determined by subtests of the intelligent test. Based on the intercorrelations of the subtests, factor analysis determines the score of each person on the factors and the correlations of the subtests with the factors, the so-called factor loadings. We will not dwell upon the classical problems of factor analysis here, which include the determination of the number of factors and the rotation problem, because these problems are inherent to factor analysis and not unique to its cross-cultural applications.

When factor analysis is applied to cross-cultural data, one major question to be considered is the (lack of) similarity of the factor analytic solution across the groups. The question amounts to a check on the equality of the factor loadings. Do the instruments (tests, items, and observational measures) have the same correlations with the factors in each cultural group? The equality of the factor loadings is sometimes visually checked, which is a questionable practice as more powerful procedures exist.

Such a procedure starts with a so-called target rotation (e.g., McDonald, 1985). Factor analytic solutions can be freely rotated (the rotation problem). This subjectivity is usually "solved" by applying a rotation procedure such as Varimax. However, independently obtained factor loadings (no matter whether they are rotated by Varimax or any other rotation procedure) may be more similar than a visual inspection may suggest. The factor loading matrices may be rotated to each other in order to maximize their agreement. This is a legitimate procedure because of the arbitrariness of factor analytic solutions. In a target rotation the axes are rotated in such a way that the agreement between the sets of factor loadings is optimized. One of the groups is arbitrarily chosen as the target to which the factor loading matrices of the other groups will be rotated.

After having rotated the factor loadings, their similarity can be evaluated in a factor by factor comparison by means of a coefficient of agreement. The most often used coefficient of agreement has been developed by Tucker (1951); it has become known as *Tucker's coefficient of agreement* and also as *proportionality coefficient* (Zegers & Ten Berge, 1985). The coefficient is comparable to a correlation coefficient, the only difference is that, unlike a correlation coefficient, the coefficient of agreement is sensitive to a constant that is added to one of the variables. As an alternative to the coefficient of agreement, the identity coefficient can be proposed that is sensitive to any linear transformation (Zegers & Ten Berge, 1985). The coefficient is defined as

$$e_{xy} = \frac{2 \sum x_i y_i}{\sum x_i^2 + \sum y_i^2},$$

in which x_i and y_i represent the factor loadings in the two groups. As a rule of thumb, values of the identity coefficient lower than .90 are taken to point to a lack of agreement and values higher than .95 are seen as evidence for the similarity of the factor matrices.

Other procedures have also been developed to evaluate the agreement of factor analytic solutions across groups. Thus, equivalence of the Eysenck personality scales (e.g., Eaves, Eysenck, & Martin, 1989; Eysenck & Eysenck, 1983) has often been studied employing a procedure proposed by Kaiser, Hunka, and Bianchini (1971). There has been some debate as to whether procedures such as those proposed by Kaiser et al. (1971) or by Tucker (1951) are sufficiently powerful to detect item bias. Using simulated data, the critics (e.g., Bijnen, Van der Net, & Poortinga, 1986; Van de Vijver & Poortinga, 1994) have shown that values well over .90 can be obtained when in fact there are items with dissimilar loadings across groups. So, caution is required because these agreement indices sometimes do not reflect the influence of nonequivalent items.

Multidimensional scaling procedures have also been employed to compare the structure of cross-cultural data sets. An example can be found in the work of Schwartz (e.g., Schwartz, 1992, 1994; see also Endler, Lobel, Parker, & Schmitz, 1991; Russell, Lewicka, & Niit, 1989). Multidimensional scaling attempts to reproduce the distances between stimuli (such as test or item scores or behavioral measures) in a small number of dimensions. To compare multidimensional scaling solutions obtained in different groups, the technique has the same rotational problem as factor analysis. Distances between stimuli are not affected by (orthogonal) rotations of the axes. Consequently, configurations of such analyses as obtained in different cultural groups have an arbitrary spatial orientation. Target rotations have to be applied prior to an evaluation of the agreement of the solutions; when such rotations are not carried out, the agreement will be underestimated. The procedure of carrying out target rotations and computing an index of agreement of the solutions described above for factor analysis also applies here. No empirical applications of target rotations following a multidimensional scaling procedure are known to the authors.

Cluster analysis is another technique that aims to reduce a large set of correlations or distances to a smaller number of dimensions or factors. However, although this technique is suitable for cross-cultural research, it has not been used much in the cross-cultural literature.

The final techniques to be discussed here have also been mentioned in the discussion of level-oriented techniques, namely the analysis of covariance structures or linear structure models, which analyzes the covariance matrix of a set of measures (Bollen, 1989; Byrne, 1989). Two common computer packages for covariance modelling, namely LISREL (Byrne, 1989; Jöreskog & Sörbom, 1993) and EQS (Bentler, 1992; Byrne, 1994), can be used to analyze multigroup data. Several applications of linear structure models in cross-cultural research can be envisaged, including confirmatory factor analysis and the analysis of multitrait-multimethod (or monotrait-multitrait) matrices for assessing method bias. Confirmatory factor analysis is a versatile tool for testing cross-cultural differences in covariance structures. Compared to the classical factor analytic procedures described above, confirmatory factor analysis allows for the test of a large set of hierarchically linked hypotheses of cross-cultural invariance. The analyses usually consist of two series. The first analysis tests whether the covariance matrix of the measures is the same for all cultural groups. Fit tests play an essential role in

covariance structure analysis. Unfortunately, fit tests in LISREL and EQS are not easy to interpret. The overall goodness-of-fit index is a chi-square distributed variable that is known to be sensitive to sample size; in large samples small inter-group differences in the covariance matrix will yield a significant value. Various fit measures have been developed that are less dependent on sample size (Bollen & Long, 1993).

If the null hypothesis is not rejected, it is highly likely that the psychological structure underlying the performance is identical across cultural groups. If the hypothesis of equal covariance matrices has to be rejected (which is usually the case), the second series of hypothesis testing will start. The second series consists of a set of hierarchically nested models that successively increase the number of equality constraints across groups. The choice of the models is free, but the order specified here (as well as in various other sources; e.g., Jöreskog & Sörbom, 1993; Vandenberg & Self, 1993; Van de Vijver & Harsveld, 1994) usually follows a theoretically relevant sequence. The first analysis of the second series specifies an equal number of factors in each group. The specification of the number of factors should be based on preliminary analyses or on earlier research findings. If the hypothesis of an equal number of factors has to be rejected, an exploratory factor analysis can be carried out to investigate the reason for the lack of fit; for instance, factors may have split up or merged in various cultural groups. If this is the case, the analyses may proceed with different numbers of factors across groups. A total lack of correspondence of the factors, after the possibilities of split and merged factors have been explored, would point to a small overlap in the psychological meaning of the instrument across the groups. Such a lack of correspondence can be expected for statistical reasons when many correlations in one or more groups are close to zero. The input variables of the factor analysis (item scores, test scores, etc.) are then largely independent of each other and the use of any multivariate technique should be strongly questioned.

The second step of the second series will test whether the matrix of factor loadings can be considered equal across the cultural groups. A set of factor loadings have to be left free in the first group; the values in the other groups are constrained to have the same values as the factor loadings in the first group. This will again yield various fit indices, among which are incremental fit indices. Because this model is subsumed in the previous model (stating an equal number of factors), the difference in the chi-square fit indices of the two models can be interpreted meaningfully: the difference in chi-square values follows a chi-square distribution with a number of degrees of freedom that is equal to the difference of the number of degrees of freedom of the two models. Acceptance of the hypothesis points to the equality of the composition of the latent factors across the groups. Rejection of the hypothesis provides evidence that the psychological structure underlying the data is dissimilar across the cultural groups. A better fit may be found when only a subset of the factor loadings are set equal to each other across cultural groups. If this is the case, subtle differences in the psychological structure have been observed. If an acceptable fit of a model specifying equal factor loadings is found, constraints can be added, such as equality of the covariance

matrices of the latent factors in all groups. The study of the fit of hierarchically nested models provides a flexible tool to analyze covariance structures. It can be used to detect smaller or larger differences in psychological meaning of measurement instruments across cultural groups. For an example of confirmatory factor analysis in cross-cultural psychology, the reader is referred to Watkins (1989).

Covariance structure analysis can also be employed for causal modelling. A set of variables, either consisting of observed variables or of a combination of observed and latent (i.e., unobservable) variables, are assumed to have *a priori* specified antecedent-consequence relationships. The model can be based on theoretical expectations or an earlier exploratory study. In some cases the exploratory study and the test of the causal model are derived from random splits of the same sample; the model is then developed on half of the data and cross-validated in the other half. An example of a cross-cultural application of a causal model can be found in Van Haaften and Van de Vijver (in press). The number of applications of causal modelling in the cross-cultural literature is limited (cf., however, Little, Oettingen, & Baltes, 1995; Little, Oettingen, Stetsenko, & Baltes, 1995), but given its flexibility and usefulness, its use is recommended.

Covariance structure analysis can also be used to assess method bias, either in a test-retest design or in a monotrait-multimethod approach. In the latter case all methods that are employed (observed variables) load on the same latent factor(s). LISREL and EQS allow for a test of the similarity of the loadings of the methods across the two groups.

Finally, it should be pointed out that the distinction between level- and structure-oriented techniques is not strict in some statistical techniques. In regression, multilevel, and covariance modelling techniques, the differentiation between level- and structure-oriented questions is quite subtle. For instance, suppose that educational achievement is predicted on the basis of a set of aptitude tests in two different cultures and equality of regression lines is tested. Similarity of regression coefficients involves structural relationships whereas equality of the intercept would refer to level-oriented relationships. The same applies to multi-level models that can tackle both level- and structure-oriented questions. In empirical applications of covariance modelling there tends to be an emphasis on structural relationships. However, the models are sufficiently flexible to deal with intergroup differences in averages as well. In sum, the designation of regression and multilevel models as level-oriented and covariance modelling as structure-oriented is more inspired by their common usage than by theoretical characteristics of these models. They could as well be seen as hybrid models.

Four Common Types of Comparative Studies

In the remainder of the chapter we shall make a distinction between four types of cross-cultural studies, depending on whether the orientation is exploratory or hypothesis testing, and on whether or not contextual factors are considered (see Table 7-2).

TABLE 7-2 Common types of comparative studies

Consideration of contextual factors	Orientation	
	No	Hypothesis testing
Yes	Generalizability	Psychological differences
	Theory-driven	External validation

The first two types emphasize hypothesis testing. The first kind of studies, *generalizability studies*, attempts to establish the generalizability of research findings obtained in one, typically Western, group to other Western or non-Western groups. In general, these studies make little or no reference to local cultural elements.

In the second type, called *theory-driven studies*, cultural factors are part of the theoretical framework. Cultural variation is deliberately sought as a validation of the model, and specific *a priori* predictions are proposed and tested. The framework is tested by sampling various cultures that differ on some focal dimension. Theory-driven studies test a theory about a particular relationship between cultural variables and a psychological outcome. Contextual elements are crucial in this type of studies.

Hypothesis testing receives little emphasis in the following two types of cross-cultural research. The first type, *psychological differences studies*, is probably most common in the literature. A measurement instrument is applied in at least two cultures and the researcher is interested in whether there are any differences in averages, standard deviations, reliabilities, or other psychometric properties of the instrument across the cultural groups. Usually, the original instrument has been applied before in a Western context, and an application of the instrument in another cultural group is thought to provide an interesting extension. There is often no compelling theory about the nature of the cross-cultural differences to be expected. Contextual factors are typically not included in the design, and post hoc explanations are invoked to interpret the cross-cultural differences observed.

The last type of cross-cultural research refers to what has been called *external validation*, which attempts to explore the meaning and causes of cross-cultural differences with the aid of contextual factors. In this type of studies, specific *a priori* hypotheses are absent and usually a large set of contextual variables are included in an exploratory manner. Only a few statistical techniques have been applied in external validation studies. Regression analysis is the most frequently applied technique, which assesses the effectiveness of independent variables in explaining cross-cultural variations in the dependent variable. This kind of validation does not address structural or scalar equivalence, but aims at providing evidence for a particular interpretation of cross-cultural differences.

Poortinga and Van de Vijver (1987) have outlined a general procedure for external validation with the inclusion of covariates. The procedure presupposes

that data (tests, observational instruments, interviews, surveys, etc.) are collected in at least two cultural groups. Data should also be collected on additional variables, termed *context variables*, that are likely to be able to explain cross-cultural differences that may be obtained. The data analysis starts with an analysis of variance to test the null hypothesis of no cultural differences. In the next analysis context variables are introduced; they are used as covariates in an analysis of covariance or as predictors in a regression analysis. In terms of an analysis of variance, the main effect of culture is tested twice; the first analysis tests group differences before correction for context variables; the second analysis tests intergroup differences in residual scores after correction. Let us call the corresponding F ratios F_1 and F_2 , respectively. Significant F values point to intergroup differences. A comparison of the significance of F_1 and F_2 can yield various possibilities. If F_1 is not significant, there are no intergroup differences to be explained (even though there is still a remote possibility that the introduction of context variables could reveal significant intergroup differences). Context variables will play a central role when F_1 is significant. Introduction of context variables can give rise to three possibilities. First, context variables may be unrelated to the dependent variable, in which case intergroup differences cannot be accounted for by these context variables. Second, context variables can be related to the dependent variable and intergroup differences on the dependent variable become smaller after correction, but they are still significant. In this case context variables provide a partial explanation of intergroup differences. Third, when F_2 is not significant, intergroup differences can be accounted for entirely by the context variables.

It should be pointed out that internal and external validation procedures have different goals. Internal validation aims at establishing the cross-cultural equivalence of the data. The key question is to ascertain whether the scores of individuals in all cultural groups can be directly compared. In external validation procedures, scalar equivalence is assumed, and the research goal is to shed light on the meaning and interpretation of the cross-cultural differences. In other words, internal validation procedures attempt to detect and remove culturally biased items, whereas external validation procedures attempt to explore the causes of cross-cultural differences observed.

Methods and Analysis of Four Common Types of Comparative Studies

In this section, the four types of cross-cultural studies—generalizability, theory-driven, psychological differences, and external validation—are examined with regard to the following issues: sampling of cultures and of subjects, procedure, design, analysis, and major strengths and weaknesses. See Table 7-3 for a summary. Examples from the literature will be described.

In *generalizability studies* a theory, a correlational or causal relationship, or an instrument derived from a theory is tested in another cultural context. The goal of the study is to establish the generalizability of the theory, the relationship, or

TABLE 7-3 Methods and analysis for the four common types of comparative studies

Type of study	Sampling of cultures	Design	Major analysis	Major strength	Major weakness
Generalizability	convenience	replication of original study or new study	structure techniques (e.g., correlations, factor analysis, analysis of covariance structures)	study of equivalence	no contextual variables included
Psychological differences	systematic or convenience sampling	replication of original study or new study	both level (e.g., <i>t</i> test and ANCOVA) and structure techniques	"openmindedness" about cross-cultural differences	ambiguous interpretation
Theory-driven	systematic (maximize contrast on focal variable)	new study; covariates may be included	both level and structure techniques	study of relationship of cultural factors and behavior	lack of attention to alternative interpretations
External validation	systematic	measures at different levels of aggregation; covariates included	level techniques	focus on interpretation of cross-cultural differences	choice of covariate variables may be meaningless

the instrument. The cultures are often chosen on the basis of convenience sampling. Two different subject sampling schemes can be applied: random or matched sampling. Generalizability will be high when the original results are replicated and the subjects are sampled randomly. However, a lack of replicability cannot be interpreted unambiguously in random sampling. Negative findings could be due to cultural differences or to a lack of equivalence of the samples. A new data set using matched samples is then required to establish a more unambiguous interpretation. The procedure of the study usually follows the procedure of the original one; in some cases stimuli may be added to enhance the appropriateness of the instrument for the local context. The design, too, is a replication of the original one.

For replications, data analysis will consist of two parts: the first part will be identical to the analysis of the original study. Second, because the establishment of generalizability is the aim of the study, an assessment of the similarity of the original and new results is required. Factor analyses, followed by target rotation and the computation of an index of agreement or multigroup analyses of covariance structures, are to be preferred over a more informal evaluation of the similarity of the outcome. Compared to studies in which results can be contrasted with those obtained in previous studies conducted in other cultures, studies that are conducted simultaneously in a number of cultures will employ more exploratory data analyses for identifying cultural similarities and differences in the results.

The major strength of generalizability studies is their ability to test the equivalence of the results across cultures. As prior data are available with which new data sets can be compared, various hypotheses about cross-cultural differences and similarities can be investigated. A weakness of generalizability studies is that they often fail to include contextual variables. If cultural differences are found, it is often not at all clear how these should be interpreted. Furthermore, bias analyses are infrequently carried out in these studies. Thus, it is too common to take unexpected differences in item scores at face value (instead of carrying out an item bias analysis).

Most examples of generalizability studies in the literature involve studies of applications of an instrument, derived from a theory. Schwartz (1992) has collected data from various countries on the universality of the structure of human values. Irvine (1979) and Vernon (1969, 1979) have compared the structure of intelligence across cultures. A study of the choice of conflict resolution procedures (Leung, 1987) is an example of a cross-cultural study of a causal relationship. Amir and Sharon's (1987) replication of Western social-psychological studies in Israel among high school and university students is another good example of a generalizability study. Finally, there are many attempts recently to validate the big-five personality factors in a variety of cultures (e.g., McCrae & Costa, 1985; McCrae & John, 1992).

In *theory-driven studies*, cultures are often systematically sampled in order to maximize their contrast on some focal variable. The sampling of subjects requires scrutiny. The cultures in the sample will often differ in many more respects than

those intended and of interest to the researcher. As the matching of the groups on all relevant ambient variables cannot be achieved, contextual measures should be added to enhance the interpretability of the findings. The measurement instruments should assess various other variables on which the cultures differ and which could obscure the cross-cultural differences being studied. The experimental procedure used is often similar across cultures. Because theory-driven studies are usually level-oriented studies, data analysis usually involves analysis of variance or covariance. In the latter case the context variables are the covariates.

The most important strength of theory-driven studies is the explicit postulation of the relationship between cultural factors and the focal behavior, which is often considered the main goal of cross-cultural psychology (e.g., Berry, Poortinga, Segall, & Dasen, 1992). The major weakness of theory-driven studies is their lack of attention to item bias and alternative explanations for the cross-cultural differences observed.

An example of a theory-driven study without covariates is Berry's (1976; Berry et al., 1986) study of the cognitive styles of hunters and food gatherers. Cultural variations in perceptual styles, educational patterns, and societal structures are all hypothesized to be interrelated and to be functionally related to the food gathering patterns of a cultural group. An example with a covariate is Earley's (1989) study in which American subjects were found to show more social loafing (the phenomenon that people work less when they are in a group than when they have to do the same task individually) than Chinese subjects. In the study, subjects' individualism–collectivism score was measured as a covariate. After controlling for cross-cultural differences in individualism–collectivism in an analysis of covariance, the cross-cultural differences in social loafing disappeared. The covariance analysis provided strong evidence for the role of individualism–collectivism in explaining cross-cultural differences in social loafing.

Studies of *psychological differences* involve the application of a measurement instrument such as a test, an interview scheme, or an observation scale, in a new cultural context. The purpose is to explore cross-cultural differences either in the magnitude recorded by the instrument or in the structure underlying the instrument. Many articles in the *Journal of Cross-Cultural Psychology* fall into this category. For instance, Guida and Ludlow (1989) compared the test anxiety of American and Chilean school children and found that for upper and middle-class subjects, American subjects reported a lower level of test anxiety than Chilean subjects. Two post hoc explanations were then given to explain this finding, none of which was tested in the study.

Two schemes for sampling cultures are employed: systematic and convenience sampling. The subjects can be chosen freely, and as usual, a choice has to be made between matched or random sampling. The procedure will typically amount to the administration of a translated instrument in a new culture. If the instrument has been applied before, the design of the previous study will usually be replicated. Covariates are typically not included in this type of study. The statistical analysis can be based on either level- or structure-oriented techniques. Quite often both aspects are combined; evidence is first provided for the similarity of

psychometric properties (e.g., reliability analysis, factor analysis and target rotations, or analysis of covariance structures), followed by an analysis of variance or *t* test at the item level.

The strength of psychological differences studies is their "open-mindedness" about the presence or absence of cross-cultural differences, a useful strategy to explore cross-cultural differences. When no cross-cultural differences are observed, it is quite likely that neither bias nor intergroup differences exist. The weaknesses of the studies are rather severe. The occurrence of bias is usually not explored. Also, because of the absence of context variables in these studies, the interpretation of the cross-cultural differences observed is not self-evident. It is often difficult to evaluate post hoc interpretations put forward to explain the observed cross-cultural differences. Finally, "fishing" may occur (Cook & Campbell, 1979). It is common that a large number of statistical tests are conducted to test the null hypothesis of no cultural differences. Such multiple testing procedures ("fishing" for significance) can easily lead to the false rejection of the null hypothesis, and hence to incorrect conclusions about the occurrence of cross-cultural differences. Various simple remedies have been proposed in the literature, such as post hoc procedures in analysis of variance or Bonferroni procedures (e.g., Glass & Hopkins, 1984; Hays, 1994). These procedures control for Type I errors when a large number of statistical tests are performed.

External validation studies start from observed cross-cultural differences. These studies aim to identify an appropriate interpretation of the differences. In some cases, external validation is based on previous studies (either generalizability or psychological differences studies) in which cross-cultural differences are reported, while in other cases the observation of cross-cultural differences and external validation are combined in one study. In both cases the choice of cultures, subjects, procedure, and design are straightforward. External validation studies usually involve survey data or secondary data (i.e., data derived from other sources, such as information on national income). External validation studies may be based on various aggregation levels (cf. the section on multilevel modelling). Most frequently reported are the individual level (e.g., when a test of individualism-collectivism is administered and is used as a covariate), an intermediate level (e.g., family and school), and the cultural level (e.g., gross national product, population density). Culture-level data can be derived from various sources such as the Human Relations Area Files (HRAF files; Barry, 1980), other cross-cultural research, and yearbooks of national and international organizations such as OECD, WHO, and UNICEF.

External validation studies are either exploratory in clarifying sources of cross-cultural differences. Analysis of covariance, regression, and causal modelling are the major statistical techniques for studying external validation.

The strength of this approach is its focus on the interpretation of cross-cultural differences, an often neglected issue in cross-cultural psychology. In principle, external validation provides a refutable framework for interpretation. However, the choice of variables and the level of analysis may be arbitrary or meaningless from a psychological point of view. As an example, the distance from a

country's capital to the equator has been found to be a good predictor of various psychological test scores, for example, of cognitive tests. It is obvious that the statistical result does not convey much information about the psychological variables underlying the performance differences.

Examples of this approach can be found in the work of Bond (1991) and Williams and Best (1982). These authors first demonstrated cross-cultural differences (in health measures in Bond's study and in gender stereotypes in Williams and Best's studies). They then related the differences to a wide variety of culture-level measures, such as values, GNP, and per capita expenditure on education and health. The results obtained allow them to interpret the cross-cultural differences observed in terms of these external variables.

Conclusion

The research question or hypothesis, method, and data analysis of cross-cultural studies are closely related. Only properly chosen methods and data analytical procedures will permit an unbiased evaluation of proposed theoretical formulations. In cross-cultural psychology, the interpretation of the meaning of research findings is crucial but evasive. Many interpretations can usually be generated to explain a cross-cultural difference, and it is often difficult to assess their validity. The best approach is to formulate a number of rival hypotheses on an *a priori* basis and design studies that are able to rule out inappropriate explanations. In our opinion, knowledge in cross-cultural psychology accumulates at an unnecessarily slow pace primarily because many cross-cultural researchers rely heavily on post hoc theorizing. This chapter is meant to encourage cross-cultural researchers to place more emphasis on methods and data analysis to improve the effectiveness of their studies. It is also hoped that the chapter will dispel the myth that methodological and statistical sophistication is an obstacle or a distraction in the research enterprise. Quite the contrary, proper methods and data analytical procedures can help clarify conceptual ambiguities, disentangle the influence of confounding variables, and rule out invalid interpretations of cross-cultural differences. Berry (1980) has stated clearly that "Cross-cultural psychology is defined primarily by its *method*" (p. 1, italic in original). Researchers who are committed to cross-cultural research should take methodological issues seriously. This chapter may facilitate cross-cultural researchers to take full advantage of the methodological and statistical procedures available for shaping their work and contributing to the advancement of the field.

References

- Allen, M. A., & Yen, W. M. (1979). *Introduction to measurement theory*. Monterey, CA : Brooks/Cole.
- Amir, Y., & Sharon, I. (1987). Are social psychological laws cross-culturally valid? *Journal of Cross-Cultural Psychology*, 18, 383-470.
- Anastasi, A. (1976). *Psychological testing* (4th ed.). New York: Macmillan.

- Angoff, W. H. (1982). Use of difficulty and discrimination indices for detecting item bias. In R. A. Berk (Ed.), *Handbook of methods for detecting bias* (pp. 96–116). Baltimore, MD: Johns Hopkins University Press.
- Barry, H. (1980). Descriptions and uses of the Human Relations Area Files. In H. C. Triandis & J. W. Berry (Eds.), *Handbook of cross-cultural psychology* (Vol. 2, pp. 445–478). Boston: Allyn and Bacon.
- Bentler, P. M. (1992). *EQS structural equation program manual*. Los Angeles: BMDP Statistical Software.
- Berk, R. A. (Ed.). (1982). *Handbook of methods for detecting item bias*. Baltimore: Johns Hopkins University Press.
- Berry, J. W. (1967). Independence and conformity in subsistence-level societies. *Journal of Personality and Social Psychology*, 7, 415–418.
- Berry, J. W. (1969). On cross-cultural comparability. *International Journal of Psychology*, 4, 119–128.
- Berry, J. W. (1976). *Human ecology and cognitive style: Comparative studies in cultural and psychological adaptation*. Beverly Hills, CA: Sage.
- Berry, J. W. (1980). Introduction to methodology. In H. C. Triandis & J. W. Berry (Eds.), *Handbook of cross-cultural psychology* (Vol. 1, pp. 1–28). Boston: Allyn and Bacon.
- Berry, J. W., & Dasen, P. R. (Eds.). (1974). *Culture and cognition: Readings in cross-cultural psychology*. London: Methuen.
- Berry, J. W., Poortinga, Y. H., Segall, M. H., & Dasen, P. R., (1992). *Cross-cultural psychology. Research and applications*. Cambridge: Cambridge University Press.
- Berry, J. W., Van de Koppel, J. M. H., Sénéchal, C., Annis, R. C., Bahuchet, S., Cavalli-Sforza, L. L., & Witkin, H. A. (1986). *On the edge of the forest: Cultural adaptation and cognitive development in Central Africa*. Lisse: Swets & Zeitlinger.
- Bijnen, E. J., Van der Net, T. Z. J., & Poortinga, Y. H. (1986). On cross-cultural comparative studies with the Eysenck Personality Questionnaire. *Journal of Cross-Cultural Psychology*, 17, 3–16.
- Bock, D. (1989). *Multilevel analysis of educational data*. New York: Academic Press.
- Bollen, K. J. (1989). *Structural equations with latent variables*. New York: Wiley.
- Bollen, K. J., & Long, J. S. (Eds.). (1993). *Testing structural equation models*. Newbury Park, CA: Sage.
- Bond, M. H. (1991). Chinese values and health: A cross-cultural examination. *Psychology and Health*, 5, 137–152.
- Brislin, R. W. (1980). Translation and content analysis of oral and written material. In H. C. Triandis & J. W. Berry (Eds.), *Handbook of cross-cultural psychology* (Vol. 1, pp. 389–444). Boston: Allyn and Bacon.
- Brislin, R. W., Lonner, W. J., & Thorndike, R. (1973). *Cross-cultural research methods*. New York: Wiley.
- Brown, E. D., & Sechrest, L. (1980). Experiments in cross-cultural research. In H. C. Triandis & J. W. Berry (Eds.), *Handbook of cross-cultural psychology* (Vol. 2, pp. 297–318). Boston: Allyn and Bacon.
- Bryk, A. S. & Raudenbush, W. (1992). *Hierarchical linear models: Applications and data analysis methods*. Newbury Park, CA: Sage.
- Buss, D. (1989). Sex differences in human mate preferences: Evolutionary hypotheses tested in 37 cultures. *Behavioral and Brain Sciences*, 12, 1–49.
- Buss, D. M., Abbott, M., Angleitner, A., Asherian, A., Biaggio, A., Blanco-Villasenor, A., Bruchon-Schweitzer, M., Ch'u, H., Czapinski, J., De Raad, B., Ekehammar, B., El Lohamy, N., Fioravanti, M., Georgas, J., Gerde, P., Guttman, R., Hazan, F., Iwawaki, S., Janakiramaiah, N., Khosrokhani, F., Kreitner, S., Lachenicht, L., Lee, M., Liik, M., Little, B., Mika, S., Moadel-Shadid, M., Moane, G., Montero, M., Mundy-Castle, A. C., Niit, T., Nsenduluka, E., Pienkowski, R., Pirttila-Blackman, A-M., Ponce de Leon, J., Rousseau, J., Runco, M. A., Safir, M. P., Samuels, C., Sanitioso, R., Serpell, R., Smid, N., Spencer, C., Tadinac, M., Todorova, E. N., Troland, K., Van den Brande, L., Van Heck, G., Van Langenhove, L., & Yang, K-S. (1990). International preferences in selecting mates. A study of 37 cultures. *Journal of Cross-Cultural Psychology*, 21, 5–47.
- Byrne, B. M. (1989). *A primer of LISREL: Basic ap-*

- lications and programming for confirmatory factor analytic models. New York: Springer.
- Byrne, B. M. (1994). *Structural equation modelling with EQS and EQS/Windows: Basic concepts, applications, and programming*. Thousand Oaks, CA: Sage.
- Campbell, D. T. (1986). Science's social system of validity-enhancing collective believe change and the problems of the social sciences. In D. W. Fiske & R. A. Shweder (Eds.), *Metatheory in social science* (pp. 108–135). Chicago: University of Chicago Press.
- Campbell, D. T., & Stanley, J. C. (1966). *Experimental and quasi-experimental designs for research*. Chicago: Rand McNally.
- Candell, G. L., & Hulin, C. L. (1987). Cross-language and cross-cultural comparisons in scale translations: Independent sources of information about item nonequivalence. *Journal of Cross-Cultural Psychology*, 17, 417–440.
- Cattell, R. B. (1940). A culture-free intelligence test, I. *Journal of Educational Psychology*, 31, 176–199.
- Cattell, R. B., & Cattell, A. K. S. (1963). *Culture Fair Intelligence Test*. Champaign, IL: Institute for Personality and Ability Testing.
- Cheung, F. M. (1989). A review on the clinical applications of the Chinese MMPI. *Psychological Assessment*, 3, 230–237.
- Cheung, F. M., Leung, K., Fan, R. M., Song, W. Z., Zhang, J. X., & Chang, J. P. (1996). Development of the Chinese Personality Assessment Inventory (CPAI). *Journal of Cross-Cultural Psychology*, 27, 181–199.
- Chinese Culture Connection (1987). Chinese values and the search for culture-free dimensions of culture. *Journal of Cross-Cultural Psychology*, 18, 143–164.
- Church, T. A. (1987). Personality research in a non-Western setting: The Philippines. *Psychological Bulletin*, 102, 272–292.
- Cicchetti, D. V., Showalter, D., & McCarthy, P. (1990). A computer program for calculating subject-by-subject kappa or weighted kappa coefficients. *Educational and Psychological Measurement*, 50, 153–158.
- Clauser, B. E., Mazor, K. M., & Hambleton, R. K. (1994). The effects of score group width on the Mantel-Haenszel procedure. *Journal of Educational Measurement*, 31, 67–78.
- Cleary, T. A., & Hilton, T. L. (1968). An investigation of item bias. *Educational and Psychological Measurement*, 28, 61–75.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and Psychological Measurement*, 20, 37–46.
- Cohen, J., & Cohen, P. (1983). *Applied multiple regression/correlation analysis* (2nd ed.). Hillsdale, NJ: Erlbaum.
- Cook, T. D., & Campbell, D. T. (1979). *Quasi-experimentation: Design and analysis issues for field settings*. Chicago: Rand McNally.
- De Groot, A., Koot, H. M., & Verhulst, F. C. (1994). Cross-cultural generalizability of the Child Behavior Checklist cross-informant syndromes. *Psychological Assessment*, 6, 225–230.
- Deregowski, J. B., & Serpell, R. (1971). Performance on a sorting task: A cross-cultural experiment. *International Journal of Psychology*, 6, 273–281.
- Dorans, N. J., & Kulick, E. (1986). Demonstrating the utility of the standardization approach to assessing unexpected differential item performance on the scholastic aptitude test. *Journal of Educational Measurement*, 23, 355–368.
- Draguns, J. (1989). Normal and abnormal behavior in cross-cultural perspective: Specifying the nature of their relationship. *Nebraska Symposium on Motivation*, 37, 235–277. Lincoln, NE: University of Nebraska Press.
- Earley, C. (1989). Social loafing and collectivism: A comparison of the United States and the People's Republic of China. *Administrative Science Quarterly*, 34, 565–581.
- Eaves, L. J., Eysenck, H. J., & Martin, N. G. (1989). *Genes, culture and personality: An empirical approach*. London: Academic Press.
- Ellis, B. B. (1989). Differential item functioning: Implications for test translations. *Journal of Applied Psychology*, 74, 912–921.
- Ellis, B. B., & Kimmel, H. D. (1992). Identification of unique cultural response patterns by means of item response theory. *Journal of Applied Psychology*, 77, 177–184.
- Embretson, S. E. (1983). Construct validity: Construct representation versus nomothetic span. *Psychological Bulletin*, 93, 179–197.
- Endler, N. S., Lobel, T., Parker, J. D., & Schmitz, P. (1991). Multidimensionality of state and

- trait anxiety: A cross-cultural study comparing American, Canadian, Israeli and German young adults. *Anxiety Research*, 3, 257-272.
- Eysenck, H. J., & Eysenck, S. B. J. (1983). Recent advances in the cross-cultural study of personality. In J. N. Butcher & C. D. Spielberger (Eds.), *Advances in personality assessment* (Vol. 2, pp. 41-69). Hillsdale, NJ: Erlbaum.
- Feldman, S. S., Rosenthal, D. A., Mont-Reynaud, R., Leung, K., & Lau, S. (1991). Ain't misbehavin': Adolescent values and family environments as correlates of misconduct in Australia, Hong Kong, and the United States. *Journal of Research on Adolescence*, 1, 109-134.
- Frijda, N., & Jahoda, G. (1966). On the scope and methods of cross-cultural research. *International Journal of Psychology*, 1, 109-127.
- Glass, G. V., & Hopkins, K. D. (1984). *Statistical methods in education and psychology* (2nd ed.). Englewood Cliffs, NJ: Prentice-Hall.
- Goldstein, H. (1987). *Multilevel models in educational and social research*. New York: Oxford University Press.
- Guida, F. V., & Ludlow, L. H. (1989). A cross-cultural study of test anxiety. *Journal of Cross-Cultural Psychology*, 20, 178-190.
- Hakstian, A. R., & Whalen, T. E. (1976). A *k*-sample significance test for independent alpha coefficients. *Psychometrika*, 41, 219-231.
- Hambleton, R. K. (1993). Translating achievement tests for use in cross-national studies. *European Journal of Psychological Assessment*, 9, 57-68.
- Hambleton, R. K. (1994). Guidelines for adapting educational and psychological tests: A progress report. *European Journal of Psychological Assessment* (*Bulletin of the International Test Commission*), 10, 229-244.
- Hambleton, R. K., & Swaminathan, H. (1985). *Item response theory: Principles and applications*. Dordrecht: Kluwer-Nijhoff.
- Harkness, S., & Super, C. M. (1990). Culture and psychopathology. In M. Lewis & S. M. Miller (Eds.), *Handbook of developmental psychopathology. Perspectives in developmental psychology* (pp. 41-52). New York: Plenum Press.
- Harman, H. H. (1976). *Modern factor analysis* (3rd rev. ed.). Chicago: University of Chicago Press.
- Hays, W. L. (1994). *Statistics* (5th ed.). Orlando, FL: Harcourt Brace & Company.
- Hess, R. D., Chang, C. M., & McDevitt, T. M. (1987). Cultural variations in family beliefs about children's performance in mathematics: Comparisons among People's Republic of China, Chinese-American, and Caucasian-American families. *Journal of Educational Psychology*, 79, 179-188.
- Ho, D. Y. F. (in press). Filial piety and its psychological consequences. In M. H. Bond (Ed.), *Handbook of Chinese Psychology*. Hong Kong: Oxford University Press.
- Hofstede, G. (1980). *Culture's consequences: International differences in work-related values*. Beverly Hills, CA: Sage.
- Hofstede, G. (1983). Dimensions of national cultures in fifty countries and three regions. In J. B. Derogowski, S. Dziurawiec, & R. C. Annis (Eds.), *Explications in cross-cultural psychology* (pp. 335-355). Lisse: Swets & Zeitlinger.
- Holland, P. W., & Thayer, D. T. (1988). Differential item performance and the Mantel-Haenszel procedure. In H. Wainer & H. I. Braun (Eds.), *Test validity* (pp. 129-145). Hillsdale, NJ: Erlbaum.
- Holland, P. W., & Wainer, H. (Eds.). (1993). *Differential item functioning*. Hillsdale, NJ: Erlbaum.
- Hulin, C. L. (1987). Psychometric theory of item and scale translations: Equivalence across languages. *Journal of Cross-Cultural Psychology*, 18, 115-142.
- Irvine, S. H. (1979). The place of factor analysis in cross-cultural methodology and its contribution to cognitive theory. In L. Eckensberger, W. Lonner, & Y. H. Poortinga (Eds.), *Cross-cultural contributions to psychology* (pp. 300-341). Lisse: Swets & Zeitlinger.
- Irvine, S. H., & Carroll, W. K. (1980). Testing and assessment across cultures. In H. C. Triandis & J. W. Berry (Eds.), *Handbook of cross-cultural psychology* (Vol. 2, pp. 181-244). Boston: Allyn and Bacon.
- Jensen, A. R. (1980). *Bias in mental testing*. New York: Free Press.
- Jöreskog, K. G., & Sörbom, D. (1993) *LISREL 8*. Chicago: Scientific Software International.

- Kaiser, H. F., Hunka, S., & Bianchini, J. (1971). Relating factors between studies based upon different individuals. *Multivariate Behavioral Research*, 6, 409-422.
- Kendall, I. M., Verster, M. A., & Von Mollendorf, J. W. (1988). Test performance of blacks in South Africa. In S. H. Irvine & J. W. Berry (Eds.), *Human abilities in cultural context* (pp. 299-339). Cambridge: Cambridge University Press.
- Kraemer, H. C. (1981). Extension of Feldt's approach to testing homogeneity of coefficients of reliability. *Psychometrika*, 46, 41-45.
- Lee, S. Y. (1990). Multilevel analysis of structural equation models. *Biometrika*, 77, 763-772.
- Leung, K. (1987). Some determinants of reactions to procedural models for conflict resolution. *Journal of Personality and Social Psychology*, 53, 898-908.
- Leung, K. (1989). Cross-cultural differences: Individual-level vs. culture-level analysis. *International Journal of Psychology*, 24, 703-719.
- Leung, K., Au, Y., Fernandez-Dols, J. M., & Iwawaki, S. (1992). Preference for methods of conflict processing in two collectivist cultures. *International Journal of Psychology*, 27, 195-209.
- Leung, K., & Bond, M. H. (1989). On the empirical identification of dimensions for cross-cultural comparison. *Journal of Cross-Cultural Psychology*, 20, 133-151.
- Leung, K., & Drasgow, F. (1985). Relation between self-esteem and delinquent behavior in three ethnic groups: An application of item response theory. *Journal of Cross-Cultural Psychology*, 17, 151-167.
- Leung, K., & Zhang, J. X. (1995). Systemic consideration: Factors facilitating and impeding the development of psychology in developing countries. *International Journal of Psychology*, 30, 693-706.
- Little, T. D., Oettingen, G., & Baltes, P. B. (1995). *The revised Control, Agency, and Means- Ends Interview (CAMI)*. Berlin: Max-Planck-Institut für Bildungsforschung.
- Little, T. D., Oettingen, G., Stetsenko, A., & Baltes, P. B. (1995). Children's action-control beliefs and school performance: How do American children compare with German and Russian children? *Journal of Personality and Social Psychology*, 69, 686-700.
- Lord, F. M. (1967). A paradox in the interpretation of group comparisons. *Psychological Bulletin*, 68, 304-305.
- Lord, F. M. (1977). A study of item bias, using Item Characteristic Curve Theory. In Y. H. Poortinga (Ed.), *Basic problems in cross-cultural psychology* (pp. 19-29). Lisse: Swets & Zeitlinger.
- Lord, F. M. (1980). *Applications of item response theory to practical testing problems*. Hillsdale, NJ: Erlbaum.
- Marsh, H. W., & Byrne, B. M. (1993). Confirmatory factor analysis of multigroup-multimethod self-concept data: Between-group and within-group invariance constraints. *Multivariate Behavioral Research*, 28, 313-349.
- McCrae, R. R., & Costa, P. T. (1985). Updating Norman's "adequacy taxonomy": Intelligence and personality dimensions in natural language and in questionnaires. *Journal of Personality and Social Psychology*, 49, 710-721.
- McCrae, R. R., & John, O. P. (1992). An introduction to the five-factor model and its applications. *Journal of Personality*, 60, 175-215.
- McCauley, C. D., & Mendoza, J. (1985). A simulation study of item bias using a two-parameter item response model. *Applied Psychological Measurement*, 9, 389-400.
- McDonald, R. P. (1985). *Factor analysis and related methods*. Hillsdale, NJ: Erlbaum.
- Mellenbergh, G. J. (1982). Contingency table models for assessing item bias. *Journal of Educational Statistics*, 7, 105-118.
- Ostroff, C. (1993). Comparing correlations based on individual-level and aggregated data. *Journal of Applied Psychology*, 78, 569-582.
- Pedhazur, E. J. (1982). *Multiple regression in behavioral research* (2nd ed.). New York: Holt, Rinehart, & Winston.
- Peterson M. F., Smith, P. B., Akande, A., Ayestaran, S., Bochner, S., Callan, V., Cho, N. G., Jesuino, J. C., D'Amorim, M., Francois, P., Hofmann, K., Koopman, P. L., Leung, K., Lim, T. K., Mortazavi, S., Munene, J., Radford, M., Ropo, A., Savage, G., Setiadi, B., Sinha, T. N., Sorenson, R., & Viedge, C. (1995). Role con-

- flict, ambiguity, and overload: A 21-nation study. *Academy of Management Journal*, 38, 429-452.
- Poon, W. Y., Chan, W., Lee, S. Y., & Leung, K. (1993). Preliminary analysis of multiple group structural equation modelling via cluster analysis. *Proceedings of the American Statistical Association 1993 Convention, Social Statistics Section*, 368-373.
- Poortinga, Y. H. (1971). Cross-cultural comparison of maximum performance tests: Some methodological aspects and some experiments with simple auditory and visual stimuli. *Psychologia Africana*, Monograph Supplement, No. 6.
- Poortinga, Y. H. (1989). Equivalence of cross-cultural data: An overview of basic issues. *International Journal of Psychology*, 24, 737-756.
- Poortinga, Y. H., & Malpass, R. S. (1986). Making inferences from cross-cultural data. In W. J. Lonner & J. W. Berry (Eds.), *Field methods in cross-cultural psychology* (pp. 17-46). Beverly Hills, CA: Sage.
- Poortinga, Y. H., & Van de Vijver, F. J. R. (1987). Explaining cross-cultural differences: Bias analysis and beyond. *Journal of Cross-Cultural Psychology*, 18, 259-282.
- Poortinga, Y.H., Van de Vijver, F. J. R., Joe, R. C., & Van de Koppel, J. M. H. (1987). Peeling the onion called culture: A synopsis. In C. Kagitcibasi (Ed.), *Growth and progress in cross-cultural psychology* (pp. 22-34). Lisse: Swets & Zeitlinger.
- Poortinga, Y. H., & Van der Flier, H. (1988). The meaning of item bias in ability tests. In S. H. Irvine & J. W. Berry (Eds.), *Human abilities in cultural context* (pp. 166-183). Cambridge: Cambridge University Press.
- Reuning, H., & Wortley, W. (1973). Psychological studies of the Bushmen. *Psychologia Africana*, Monograph Supplement, 7.
- Roberts, J., & Sutton-Smith, B. (1962). Child training and game involvement. *Ethology*, 1, 166-185.
- Russell, J. A., Lewicka, M., & Niit, T. (1989). A cross-cultural study of a circumplex model of affect. *Journal of Personality and Social Psychology*, 57, 848-856.
- Sachs, J. (1992). Covariance structure analysis of a test of moral orientation and moral judgment. *Educational and Psychological Measurement*, 52, 825-833.
- Scheuneman, J. D. (1987). An experimental, exploratory study of causes of bias in test items. *Journal of Educational Measurement*, 24, 97-118.
- Schwartz, S. H. (1992). Universals in the content and structure of values: Theoretical advances and empirical tests in 20 countries. In M. Zanna (Ed.), *Advances in experimental social psychology* (Vol. 25, pp. 1-65). Orlando, FL: Academic Press.
- Schwartz, S. H. (1994). Studying human values. In A. Bouvy, F. J. R. Van de Vijver, P. Boski, & P. Schmitz (Eds.), *Journeys into cross-cultural psychology* (pp. 239-254). Lisse: Swets & Zeitlinger.
- Segall, M. H., Campbell, D. T., & Herskovits, M. J. (1966). *The influence of culture on visual perception*. Indianapolis, IN: Bobbs-Merrill.
- Serpell, R. (1979). How specific are perceptual skills? *British Journal of Psychology*, 70, 365-380.
- Serpell, R. (1993). *The significance of schooling. Life-journeys in an African society*. Cambridge: Cambridge University Press.
- Shrout, P. E., & Fleiss, J. L. (1979). Intraclass correlations: Uses in assessing rater reliability. *Psychological Bulletin*, 86, 420-428.
- Smith, P. B., & Peterson, M. F. (1988). *Leadership, organizations and culture*. Beverly Hills, CA: Sage.
- Sternberg, R. J. (1985). Implicit theories of intelligence, creativity, and wisdom. *Journal of Personality and Social Psychology*, 49, 607-627.
- Stricker, L. J. (1982). Identifying test items that perform differentially in population subgroups: A partial correlation index. *Applied Psychological Measurement*, 6, 261-273.
- Super, C. M. (1981). Behavior development in infancy. In R. H. Munroe, R. L. Munroe, & B. B. Whiting (Eds.), *Handbook of cross-cultural human development* (pp. 181-270). New York: Garland SPTM Press.
- Super, C. M. (1983). Cultural variation in the meaning and uses of children's "intelligence." In J. B. Deregowski, S. Dziurawiec, & R. C. Annis (Eds.), *Expiscations in cross-cultural psychology* (pp. 199-212). Lisse: Swets & Zeitlinger.

- Tucker, L. R. (1951). *A method for synthesis of factor analysis studies* (Personnel Research Section Report No. 984). Washington, DC: Department of the Army.
- Tylor, E. B. (1871). *Primitive culture* (2 vols.). London: Murray.
- Van de Vijver, F. J. R. (1988). Systematizing item content in test design. In R. Langeheine & J. Rost (Eds.), *Latent trait and latent class models* (pp. 291-307). New York: Plenum.
- Van de Vijver, F. J. R. (1994). Item bias: Where psychology and methodology meet. In A. Bouvy, F. J. R. Van de Vijver, P. Boski, & P. Schmitz (Eds.), *Journeys into cross-cultural psychology* (pp. 111-126). Lisse: Swets & Zeitlinger.
- Van de Vijver, F. J. R., Daal, M., & Van Zonneveld, R. (1986). The trainability of abstract reasoning: A cross-cultural comparison. *International Journal of Psychology*, 21, 589-615.
- Van de Vijver, F. J. R., & Harsveld, M. (1994). The incomplete equivalence of the paper-and-pencil and computerized version of the General Aptitude Test Battery. *Journal of Applied Psychology*, 79, 852-859.
- Van de Vijver, F. J. R., & Poortinga, Y. H. (1982). Cross-cultural generalization and universality. *Journal of Cross-Cultural Psychology*, 13, 387-408.
- Van de Vijver, F. J. R., & Poortinga, Y. H. (1991). Testing across cultures. In R. K. Hambleton & J. Zaal (Eds.), *Advances in educational and psychological testing* (pp. 277-308). Dordrecht: Kluwer.
- Van de Vijver, F. J. R. & Poortinga, Y. H. (1992). Testing in culturally heterogeneous populations: When are cultural loadings undesirable? *European Journal of Psychological Assessment*, 8, 17-24.
- Van de Vijver, F. J. R., & Poortinga, Y. H. (1994). Methodological issues in cross-cultural studies on parental rearing behavior and psychopathology. In C. Perris, W. A. Arrindell, M. Eisemann (Eds.), *Parental rearing and psychopathology* (pp. 173-197). Chichester: Wiley.
- Van den Wollenberg, A. L. (1988). Testing a latent trait model. In R. Langeheine & J. Rost (Eds.), *Latent trait and latent class models* (pp. 31-50). New York: Plenum.
- Van Haaften, E. H., & Van de Vijver, F. J. R. (in press). Psychological consequences of environmental degradation, *Journal of Health Psychology*.
- Vandenberg, R. J., & Self, R. M. (1993). Assessing newcomers' changing commitments to the organization during the first 6 months of work. *Journal of Applied Psychology*, 78, 557-568.
- Vernon, P. E. (1969). *Intelligence and cultural environment*. London: Methuen.
- Vernon, P. E. (1979). *Intelligence: Heredity and environment*. San Francisco: Freeman.
- Watkins, D. (1989). The role of confirmatory factor analysis in cross-cultural research. *International Journal of Psychology*, 24, 685-702.
- Werner, O., & Campbell, D. T. (1970). Translating, working through interpreters, and the problem of decentering. In R. Naroll & R. Cohen (Eds.), *A handbook of cultural anthropology* (pp. 398-419). New York: American Museum of Natural History.
- Whiting, B. B. (1976). The problem of the packaged variable. In K. Riegel & J. Meacham (Eds.), *The developing individual in a changing world* (Vol. 1, pp. 303-309). The Hague: Mouton.
- Williams, J. E., & Best, D. L. (1982). *Measuring sex stereotypes: A thirty-nation study*. Beverly Hills, CA: Sage.
- Wilss, W. (1982). *The science of translation: Problems and methods*. Tuebingen: Narr.
- Winer, B. J. (1971). *Statistical principles in experimental design* (2nd ed.). New York: McGraw Hill.
- Yang, K. S., & Bond, M. H. (1990). Exploring implicit personality theories with indigenous or imported constructs: The Chinese case. *Journal of Personality and Social Psychology*, 58, 1087-1095.
- Zegers, F. E., & Ten Berge, J. M. F. (1985). A family of association coefficients for metric scales. *Psychometrika*, 50, 17-24.