

Mengyao Hu, 2016

[Introduction](#)

[Guidelines](#)

- [1. Exploring the variables.](#)
- [2. Simple and multiple linear regression models.](#)
- [3. Categorical data analysis.](#)
- [4. Multilevel models.](#)
- [5. Longitudinal analysis.](#)
- [6. Latent variable models.](#)
- [7. Differential item functioning.](#)
- [8. Machine learning.](#)
- [9. Incorporate complex survey data features.](#)
- [10. Introduction to Bayesian inference.](#)

[References](#)

[Appendices: A](#)

Introduction In recent years, the number and scope of multinational, multicultural, or multiregional surveys, which to as '3MC' surveys, has increased dramatically. With the increased availability of large datasets covering multiple countries such as the European Social Survey (ESS) and the Survey of Health, Ageing and Retirement in Europe (SHARE), many researchers have become engaged in analyzing these data. Not surprisingly, there has been increased interest in the development of statistical tests appropriate for cross-cultural survey data analysis. This chapter aims to provide a comprehensive introduction of different statistical methods, from basic statistics to advanced modeling approaches. However, that this chapter does not aim to teach statistics, but rather to provide an overview of what statistical tests are available when to apply them in 3MC research. We also provide links and references to each statistical method for those who like additional detail.

[↑ Back](#)

Guidelines

1. Exploring the variables.

1.1 The classification of variable types is important, because it will help to determine which statistical procedure be used. For example, when the dependent variable is continuous, a linear regression can be applied (see [Guideline 2.2](#)); when it is categorical (binary), a logistic model can be applied (see [Guideline 3.2](#)); when it is categorical (nominal or ordinal), multinomial, or ordinal, logistic regressions may be used (see [Guideline 3.3](#)). If, in latent variable models (see [Guideline 6](#)), the latent variable is continuous, confirmatory factor analysis (CFA) or item response theory (IRT) models can be used (see [Guidelines 6.1](#) and [6.5](#)). Figures 1 and 2 below list the choices of regression and latent variable measurement models as regard the variable types of the dependent and independent variables. Several commonly used variable types are listed as below:

- **Nominal variables:** variable values assigned to different groups. For example, respondent gender may be 'male' or 'female.'
- **Ordinal variables:** categorical variables with ordered categories. For example, 'agree,' 'neither agree nor disagree,' or 'disagree.'
- **Continuous variables:** variables which take on numerical values that measure something. "If a variable can take any value between two specified values, it is called a continuous variable; otherwise, it is called a discrete variable." Continuous variables are understood to have equal intervals between each adjacent pair of values in the distribution. Income is an example of a continuous variable.
- **Discrete (ratio) variable:** "a discrete variable can only take on a finite value, typically reflected as a whole number". The variables have an absolute '0' value. One example is the number of children a person has.

Figure 1: Variable types and choices of regression models.

(image)

Figure 2: Variable types and choices of latent variable measurement models (see [Guideline 6](#) for more details).

(image)

1.2 Distribution of variables:

1.2.1 **Graphical illustrations of distributions:** it is commonly recommended to look at graphical summaries of both continuous and categorical distributions before fitting any models. Details of the graphical options listed below can be found in this online statistics book: [Online Statistics Education: An Interactive Multimedia of Study](#).

- For categorical variables:
 - Bar graphs
 - Pie charts
- For continuous or discrete variables:
 - Stem and Leaf Plots
 - Histograms
 - Box plots
- For any type of variable:
 - Frequency distributions

In 3MC data analysis, to get a direct visual comparison, researchers can plot distributions by country or racial group.

1.2.2 **Numerical summaries of distributions:** a distribution can be summarized with various descriptive statistics. The mean and median capture the center of a distribution (central tendency), while the variance describes the distribution spread or variability (see [online book material](#)).

- **Mean:** the average of a number of values. It is calculated by adding up the values and dividing by the number of values added.
- **Median:** the “...median is the number separating the higher half of a data sample, a population, or a probability distribution, from the lower half” . For a highly skewed distribution, the median may be a more appropriate measure of central tendency than the mean. For example, the median is more widely used to characterize income, since potential outliers (e.g., those with very high incomes) have much more impact on the mean.
- **Variance:** a measure of the extent to which a set of numbers are 'spread out.'
- **Precision:** the reciprocal of the variance; most commonly seen in Bayesian analysis (see [Guideline 9](#)).

1.3 Suggested reading:

- Tests of the equality of two means (see [online material](#))
-
-

1.4 Potential uses in 3MC research:

- A good starting point for analysis is to look at the distributions of variables of interest and at graphical illustrations of the variables in each cultural group.
- One way of comparing survey estimates across various cultures is to directly compare mean estimates. A two-sample t-test can be used to evaluate the equality of two means (see [Guideline 1.3](#)). However, researchers need

aware that the observed mean differences are not necessarily equal to the latent construct mean differences ([Guideline 6](#)), and direct comparison using observed mean differences may lead to invalid results (see). In a factors irrelevant to the question content, such as response style differences in different cultures, may influence comparability across cultures. More advanced models (such as latent variable models) can be used to evaluate control for these factors.

[↑ Back](#)

2. Simple and multiple linear regression models.

2.1 A bivariate relationship is the relationship between two variables. For example, one may be interested in how height is associated with weight (i.e., whether those who are taller tend to weigh more). Basic information about bivariate relationships can be found [here](#).

- **Scatterplots:** before running any models, a scatterplot is essential to explore the associations (negative or positive) between variables.
- **Correlations between variables:** Pearson's correlation is the most commonly used method of evaluating the relationship between two variables. Refer to [this website](#) for more information.

2.2 Linear regression models can allow researchers to predict one variable using other variable(s). The dependent variable in linear regression models is a continuous variable. Basic information about simple linear and multiple regression models can be found [here](#).

- **ANOVA table:** in the output of regression model results, an analysis of variance (ANOVA) table is usually provided, consisting "of calculations that provide information about levels of variability within a regression and form a basis for tests of significance".
 - [Resource 1](#)
 - [Example of a regression model results output using Stata](#)
- **Dummy predictor variables:** as described by [Stata](#), a dummy variable or *indicator variable* is an artificial variable created to represent an attribute with two or more distinct categories or levels. If a categorical variable is added directly to the regression models without being specially specified, the software will treat it as continuous. If the differences between the categories (e.g., category 2 minus category 1) do not have any actual meaning, dummy variables are usually created in these situations to make sure that such categorical variables are correctly specified in the model. For example, in 3MC data analysis, to compare Country A to Country B on the level of the dependent variable, one can create a country dummy variable using one of the countries as a reference group and add it as an independent variable to the model. When multiple countries exist, one can use one of the countries as the reference category, and treat the variable as categorical in the model.
 - For information on dummy variables and how they are created and used, see [this website](#).
 - For information on regression models with categorical predictors using SAS, see [here](#).
- **Interactions of predictor variables:** sometimes a regression model is used to test whether the relationship between the dependent variable (DV) and one specific independent variable (IV) depends on another IV. To test this, an interaction term between the two IVs can be added to the model.
 - [Resource 1](#)
- **Transformations of variables:** when non-linearity is found for predictors, transformations may be considered to 'normalize' a variable which has a skewed distribution. For more detail, see [this website](#).
- **Lack of fit testing:** various techniques are available to test for the lack of fit in regression models, including graphical methods (e.g., $Q-Q$ plots) and numerical methods (e.g., R^2 and F -tests).
- **Model diagnostics:** techniques are available to test the appropriateness of the model and whether the model's assumptions hold.
 - [Resource 1](#)
 - [Resource 2](#) (using R)

- **Selecting reduced regression models** (variable selection): techniques for determining the model which contains the most appropriate independent variables, giving the maximum R^2 value.
 - [Resource 1](#) (Including SAS code)
 -

2.3 Suggested reading:

- Applied Statistical Analysis and Data Display: An Intermediate Course with Examples in S-PLUS, R, and SAS
- Statistical Methods, 8th ed
- The Little SAS Book, 4th ed

2.4 Potential uses in 3MC research:

- As in linear regression models, a country variable/indicator can be added to the regression model as a covariate (e.g.,).

[↑ Back](#)

3. Categorical data analysis.

3.1 **Analysis of two-way tables:** categorical data are often displayed in a two-way table. Sometimes, one or both variables are continuous. If so, the continuous variable(s) can be categorized into groups. A two-way table can be constructed using the new variables. Note that this approach may lead to a loss of information on the continuous variables. See also [online material](#).

3.1.1 The Pearson chi-square test evaluates whether the row and column variables in a two-way table are associated.

- [Resource 1](#)

3.1.2 Odds ratios (OR) and relative risks (RR) describe the proportions in contingency tables. See for a comprehensive introduction.

- [Resource 1](#)
-

3.1.3 Log-linear models are commonly used to model the cell counts of contingency tables such as two-way tables.

- [Resource 1](#)

3.2 **Logistic regression models:** these can be used when the dependent variable is a binary categorical variable. This technique allows researchers to model or predict the probability that an individual will fall into one specific category given other independent variables. Logistic regression is a type of generalized linear model where the logit function for selecting one category is expressed through a linear function of the predictors. Thus, as in other linear regression models, the predictors can include both continuous and categorical variables.

- [\(link\)](#)
- [Resource 2](#)

3.3 **Multinomial and ordinal logistic regressions:** when the DV is a nominal variable, a multinomial logistic regression model can be used. If the DV is an ordinal variable, an ordinal logistic regression can be used.

- [Resource 1](#)

3.4 Suggested reading:

-
-
-

3.5 Potential uses in 3MC research:

- To evaluate responses to a categorical variable across two different cultures, one can construct a two-way table using the categorical variable and the country indicator as the rows and columns. A Pearson chi-square test is used to evaluate whether the variable differs by cultures.
- As in logistic regression models, a country variable/indicator can be added to the logistic regression model as a covariate.

[↑ Back](#)

4. Multilevel models.

Multilevel models are usually used when there is a hierarchical structure, such as when sampling units are nested in geographical areas (e.g., cluster sampling) and when they are selected in longitudinal studies. Multilevel models are known as hierarchical linear models, mixed models, random effects models, and variance components models. The Center for Multilevel Modeling at the University of Bristol offers a [free online course](#) on multilevel modeling. Additional information on multilevel modeling can be found [here](#) and in [this book](#). When many cultural groups are present, a multilevel framework can be used, with country being treated as a random variable. Multilevel models with latent variables, such as multilevel structural equation models (MLSEM), can also be run, as discussed by [this book](#) and [this book](#). See [Guideline 6.4](#) for more information on SEM models.

4.1 Suggested reading:

-
-
-
-

4.2 Potential uses in 3MC research:

- Many 3MC studies have a multilevel data structure with respondents nested within countries. Recent research in multilevel cross-cultural research has emerged in last several decades. For more information, see [this book](#).

[↑ Back](#)

5. Longitudinal analysis.

Longitudinal data analysis refers to techniques used to evaluate data collected through repeated measures.

- [Resource 1](#)

5.1 **Modeling longitudinal/panel data:** in panel surveys, respondents are interviewed at multiple points in time, producing 'panel' or 'longitudinal' data. The first step in analyzing longitudinal data is to look at the descriptive statistics, then, select one of several possible methods of analysis. The traditional technique is the repeated measures analysis of variance (rmANOVA), although this has several limitations. More commonly used approaches include multilevel models and marginal models.

5.1.1 **Descriptive plots:** The 'spaghetti' plot “involves plotting a subject’s values for the repeated outcome n (vertical axis) vs. time (horizontal axis) and connecting the dots chronologically” . Plots can be created at the individual data level and the mean level. For binary outcomes, proportions can be used to generate them for different population groups. In 3MC studies, the plots can be generated for different cultural or country groups.

- [Resource 1](#)

5.1.2 Repeated measures analysis of variance (rmANOVA):

- For more information, see this [online example](#) on rmANOVA.
- The rmANOVA approach is not recommended due to the limitations as mentioned below:
 - Subjects missing any data will not be included in the analysis.
 - A limited number of covariance structures are allowed.
 - Time-varying covariates are not allowed.

5.1.3 **Multilevel models for longitudinal data:** multilevel models account for between respondent variance including random effects in the model, such as random slope and random intercept.

- [Resource 1](#)
- [Resource 2](#)
- (presentation slides available [here](#))

5.1.4 **Marginal modeling approaches:** If the between-subject variation is not of interest, the marginal modeling approach, where only the correlated error terms are included in the model, can be used—no random effects are added to the model.

-

5.2 Suggested reading:

-
-
-
-
-
-
-
-
-

5.3 Potential uses in 3MC research:

- A country variable/indicator can be added to the marginal models as a covariate, or it can be added in a multilevel model as a fixed effect.

[↑ Back](#)

6. Latent variable models.

Latent variable models include both observed variables (the data) and latent variables. A latent variable is unobservable which represents hypothetical constructs or factors . A latent variable can be measured by several observed variables. An example of latent variable provided by describes the construct of intelligence. As mentioned by , “there is no single, definitive measure of intelligence. Instead, researchers use different types of observed variables, such as tasks of verbal

reasoning or memory capacity, to assess various facets of intelligence.” Examples of such latent variables are usually measured in a measurement model which evaluates the relationship between latent variables and their indicators.

6.1 **Exploratory factor analysis (EFA) and confirmatory factor analysis (CFA):** these are two types of measurement models where latent variables are indicated by multiple observed variables. The difference between EFA and CFA is related to the existence of a hypothesis about the measurement model before doing the analysis. As mentioned by Bollen (1989), “CFA attempts to confirm hypotheses and uses path analysis diagrams to represent variables and factors, whereas EFA tries to uncover complex patterns by exploring the dataset and testing predictions.” As seen in Figure 3, since both indicators and the latent variables are all continuous, both EFA and CFA are based on linear functions. Figure 3 illustrates the differences between EFA and CFA. Since EFA is purely data-driven and may be arbitrary in nature, it is suggested by some literature to always use CFA, which is theory-driven. However, as mentioned by Bollen (1989), it is more appropriate to use EFA rather than CFA when the theory is not well established.

See [Appendix A](#) for a comprehensive overview on EFA and for CFA. The code for conducting EFA and CFA are included in [Appendix A](#).

Figure 3: Comparisons between EFA and CFA*.

(image)

*Adapted from *Exploratory and Confirmatory Factor Analysis presentation*.

Multigroup CFA (MCFA) is commonly used in 3MC research for measurement equivalence testing. The basic idea is to start with the same model, but allow the coefficients to differ by groups (assuming configural equivalence), then start introducing constraints in the model coefficients—such as to make them equal across the groups. The model fit of the previously run models can be compared. Among all the models, the parsimonious model with the best fit solution will be selected to evaluate the data. If the model reveals no violations of scalar equivalence, the country means can be compared directly. In a panel study, with data available at different time points, one can also evaluate measurement equivalence across cultures over time. See [Guideline 6.2](#) below for more information on measurement equivalence testing. For more information on MCFA, see:

- [Resource 1](#)
-
-

6.2 **Measurement equivalence in 3MC research:** as mentioned by Bollen (1989), “measurement equivalence implies that a measurement instrument used in different cultures measures the same construct.” There are different levels of measurement equivalence. Three of the most widely discussed levels are configural, metric, and scalar equivalence. These three levels are hierarchical, where the higher ones have higher requirements of equivalence and require the achievement of the lower ones.

Configural equivalence refers to similar construction of the latent variable. In other words, the same indicators are associated with the latent concepts in each culture. It does not require each culture to view the concept in the same way. For example, it allows the strength (i.e., loadings) to be different across cultures. *Metric equivalence* requires the same slope across cultures, which captures the associations between indicator and the latent variable. In other words, it implies “the equality of the measurement units or intervals of the scale on which the latent concept is measured across cultural groups”. *Scalar equivalence* implies that on the basis of equality of the measurement units, the country means of the latent variable also have the same origin across cultures. Under this equivalence level, the model achieves measurement equivalence, and researchers can compare the country scores directly.

In situations where full equivalence is difficult to achieve, researchers also evaluate the conditions under which different cultures achieve partial equivalence. An example of partial equivalence is when most of the indicators are equivalent across cultures, but only one has a different slope and threshold across cultures. One can then conclude

the different cultures achieve partial equivalence where they differ on one specific indicator. As mentioned by “partial equivalence enables a researcher to control for a limited number of violations of the equivalence requirement and to proceed with substantive analysis of cross-cultural data” .

The aforementioned approaches to assessing measurement equivalence have been widely used in 3MC survey analysis. However, it has recently been criticized for being overly strict. As mentioned by , it is difficult to achieve scalar equivalence, or even metric equivalence, in surveys with many countries or cultural groups. A Bayesian approximate equivalence testing approach has been recently proposed by . This approach allows for “small variations in parameters across different cultural groups; thus, when approximate scalar measurement equivalence is reached, one can compare across cultures meaningfully, even though the traditional method may indicate scalar inequivalence. It also points out that approximate measurement inequivalence or invariance are favorable for cases with large-scale data containing many groups or repeated measurements, small differences in intercepts and factor loadings, and small differences which are cancelled out both within and between groups. In contrast, approximate measurement invariance testing seems to lead to bias in latent mean estimates if large differences in intercepts or factor loadings or systematic differences between groups exist . For introductions and references of Bayesian methods, see [Guideline 9](#).

6.3 Latent class analysis (LCA): unlike the previously mentioned approach, such as CFA and SEM, where the variables are continuous, LCA treats the latent variables as categorical (nominal or ordinal) (see [Figure 3](#)). The categories of the latent variable in LCA are referred to as classes, which represent “a mixture of subpopulation membership is not known but is inferred from the data” . That is to say, LCA can classify respondents into different groups based on their attitudes or behaviors, such as classifying respondents by their drinking behavior. Respondents in the same group are similar to each other, regarding the behavior/attitudes, and they differ from those in other groups —i.e., heavy drinkers vs. non-drinkers. One can also add covariates to the model if those measures can influence class membership. In a second step, the class membership from the model can then be used for followup analysis. For example, to better understand the differences between respondents, a logistic (or multinomial logistic, if more than two groups) regression model can be run in which selected covariates are used to predict class membership. Alternatively, to evaluate the influence of class membership on other variables, LCA can also be used in regression models as a covariate to predict other outcomes. For more information on LCA, please see:

- [Resource 1](#)
-

As mentioned by , when testing for measurement invariance with latent class analysis, “the model selection procedure usually starts by determining the required number of latent classes or discrete latent factors for each group. ... If the number of classes is the same across groups, then the heterogeneous model is fitted to the data; followed by a series of nested, restricted models which are evaluated in terms of model fit.” That is to say, unlike multigroup CFA, the multigroup LCA will need to identify whether the number of classes are the same across groups before testing different models at different invariance levels. See , , and for more information.

6.4 Structural equation modeling (SEM) is a multivariate analysis technique used in many disciplines which test the causal relationship hypothesis between variables . It usually includes two components: 1) the measurement model, which summarizes several observed variables using their latent construct (e.g., CFA as discussed in [Guideline 6.1](#)), and 2) the structural model, which describes the relationship between multiple constructs (e.g., relationships among both latent and observed variables).

6.4.1 Similar to previously discussed latent variable models, SEM can have both observed and latent variables, where observed variables are the data collected from respondents and latent variables represent unobserved constructs and factors . The observed variables which are used as measures of a construct are indicators of the latent variable. In other words, the latent variable is indicated by these observed variables. Besides observed variables, SEM models also include error terms, similar to the error terms in a regression analysis. As mentioned by , “a residual term represents variance unexplained by the factor that the corresponding indicator

supposed to measure. Part of this unexplained variance is due to random measurement error, or score unreliability.”

6.4.2 In SEM analysis, parameter estimation is done by comparing the model-based covariance matrix with based covariance matrix. The goal of this approach is to evaluate whether the model with best fit is supported by the data—that is, whether the two covariance matrices are consistent with each other and whether the model explains as much of the variance of the data.

6.4.3 Structural equation models can also estimate the means of latent variables. It also allows researchers to estimate between- and within-group mean differences. In 3MC analysis, one can estimate the group mean differences of latent variables, such as those between two cultures.

6.4.4 Suggested reading:

-
-
-

6.5 **Item response theory (IRT)** is commonly used for psychometric and educational testing, and is becoming increasingly popular in 3MC analysis. It begins with the idea that when answering a specific question, the response provided by an individual depends on the ability/qualities of the individual and the qualities of the question item. As mentioned earlier, “the mathematical foundation of IRT is a function that relates the probability of a person responding to an item in a specific manner to the standing of that person on the trait that the item is measuring. In other words, the function describes, in probabilistic terms, how a person with a higher standing on a trait (i.e., more of the trait) is likely to provide a response in a different response category to a person with a low standing on the trait.” Therefore, IRT allows researchers to model the probability of a specific response to a question item, given the item and the individual’s latent level.

The simplest IRT model is the Rasch model, also called the one-parameter IRT model, which assumes equal item discrimination (“the extent to which the item is able to distinguish between individuals on the latent construct being measured”). The Rasch model starts from the premise that the probability of giving a 'positive' answer to a yes/no question is “a logistic function of the distance between the item’s location, also referred to as item difficulty, and the person’s location on the latent construct being measured,” also known as the person’s latent trait level. There are other types of IRT models available as well, which can be categorized by the number of parameters and the question response option for (e.g., binary or multiple response options and whether ordinal or nominal). Table 1 below summarizes the different types of IRT models.

In a two-parameter (2PL) IRT model, an item discrimination parameter is also included in the model. The two parameters are “analogous” to the factor loadings in CFA and EFA, since they all represent “the relationship between the latent trait and item responses”. A three-parameter (3PL) IRT model also includes a 'guessing' parameter, which describes the situation wherein when a question can be answered by guessing, the probability of giving a correct answer is higher than zero even for those with low latent trait level.

For items with multiple response options (ordinal or nominal variables), polytomous IRT models can be used. See Table 1 for details. In this chapter, we will not discuss these models in detail. See suggested readings on IRT for more information.

Table 1: Variable type and IRT model choices.

Type of observed variable	Model
Binary	1 parameter-logistic (1PL) model/Rasch model

	2PL model
	3PL model
Multiple response options (ordinal)	Graded response model/Thurstone/Samejima polytomous models
	Partial credit model (PCM)/Graded PCM
Multiple response options (nominal)	Rating scale model
	Nominal response model/Bock's model

6.5.1 Suggested reading:

-
-
-
-
-
-

6.6 Besides what is discussed above, other types of latent variable models are available. Some examples are listed below:

- The **latent transition model** is “a special kind of latent class factor model that represents the shift from one different states, such as from nonmastery to mastery of a skill, is a latent transition model” . See [these slides](#) for more information.
- In **latent profile models**, the latent variable is categorical and the indicators are continuous. It is commonly used in cluster analysis. See for more information.
- The **mixed Rasch model** is “a combination of the polytomous Rasch model with latent class analysis” . See for more information.
- When we have data where the population of individuals are divided into different groups, such as in a 3MC **multilevel structural equation modeling** (MLSEM) can be used. This model decomposes individual data into within-group and between-group components, and can simultaneously estimate within- and between-group parameters. For more information on MLSEM, see .

6.7 Potential uses in 3MC research:

- As mentioned by , the observed mean does not equal the latent mean, where the observed mean is a function of intercepts, factor loadings, and the latent mean. Similarly, “observed mean differences between two or more groups (e.g., cultures) do not necessarily indicate latent mean differences as unequal intercepts and/or factor loadings also lead to observed differences” . To conduct more valid comparisons across different groups (e.g., cultures), measurement invariance testing is a widely used method which aims to evaluate whether the latent means of different groups are comparable. In other words, it evaluates whether the different groups differ in factor loadings and intercepts of the measures. See for more information.
- Measurement invariance testing is usually conducted within the multigroup analysis (MGA) framework. The most commonly used method is multigroup confirmatory factor analysis (MG-CFA) (e.g.,). Other types of MGA include multigroup structural equation modeling (MG-SEM) analysis (e.g.,), multigroup latent class analysis (e.g.,), multigroup IRT model (e.g.,) and multigroup mixed Rasch model (e.g.,). See for more information.
- A recent paper by discusses misconceptions in measurement equivalence analysis. Using data from the World Values Survey, they show that “constructs can entirely lack convergence at the individual level and nevertheless

- exhibit powerful and important linkages at the aggregate level”.
- examine new approaches to measurement invariance evaluation in the context of international large-scale assessments (ILSAs), which aim to produce cross-national comparisons of student outcome measures (e.g., achievement, behaviors, and attitudes). begin by reviewing the extent to which exact measurement invariance can be confirmed; then, using data from the International Civic Citizenship Education Study 2009, they extend the model to include partial invariance, indirect/covariate effects, and a Bayesian approach to evaluating the approximate invariance assumption. They find that the Bayesian approach, assuming a parsimonious variance-covariance procedure, can play a key role in evaluating measurement invariance across national samples used in ILSAs. They find that the integration of a strong model with less stringent assumptions is sufficient for estimating comparable estimates of the latent constructs and allows the reporting of a table charting mean comparisons across a large number of groups. Yet, a high degree of consistency between their results and models assuming exact invariance indicates that previously published results are not invalidated in any substantial way.

[↑ Back](#)

7. Differential item functioning.

Differential item functioning (DIF) is a statistical concept developed to identify to what extent a question item might be measuring different properties for individuals of separate groups (i.e., ethnicity, culture, region, language, sex, etc.). Items that function differently across groups are used as indicators for 'item bias' if the items function in a systematically different way across cultures. To detect DIF, several different methods can be used, as listed below:

- Mantel-Haenszel (MH) statistic is regarded as a “reference” technique of detecting DIF due to its ease of use and the fact that it can be applied to small samples. The disadvantage of MH statistic is that it does not allow for statistical significance testing.
- Logistic regression can be used as an alternative method to detect DIF. For more information, see [Larsen and Bracken \(2002\)](#).
- DIF can be detected using an IRT framework. Item characteristic curves (ICCs) of the same item but from different groups can be compared to evaluate whether there is DIF. For more information, see [Larsen and Bracken \(2002\)](#) and [Larsen and Bracken \(2002\)](#).

[↑ Back](#)

8. Machine learning.

Machine learning is “a general term for a diverse number of classification and prediction algorithms” which have applications in many different fields. Unlike statistical modeling approaches, machine learning evaluates the relationship between outcome variables and predictors using a “learning algorithm without an a priori model”. Below, we introduce several machine learning methods.

- 8.1 The **classification tree** is a data-driven method which allows researchers to evaluate the complex interactions between variables when there are many predictor variables present. In binary trees, the nodes of the tree are divided into two branches. To reasonably construct and prune a given tree, deviance measure is used to choose the split. The 'rpart' package is used for classification tree analysis (see [Appendix A](#) for more information). The classification tree result can be evaluated through apparent error rate and true error rate. The former is the error rate when the tree is applied to a training data set, and the latter is when it is applied to a new data set or to test data. In evaluating error rate, researchers usually divide the data into two parts—training data and test data—and validate the tree on the test data set.
- 8.2 Random forest is an algorithm for classification which uses an “ensemble” of classification trees. Through averaging over a large ensemble of “low-bias, high-variance but low correlation trees,” the algorithm yields an ensemble that can achieve both “low bias and low variance”.
- 8.3 Suggested reading:

- [Resource 1](#)
- [Resource 2](#)
-
-
-
-

8.4 Potential uses in 3MC research:

- Classification tree analysis in cross-cultural research allows researchers to evaluate 1) the important factors culture and 2) how the factor interactions differ across cultures. One study used classification tree to evaluate college student alcohol consumptions across American and Greek students, and found that “student attitude drinking were important in the classification of American and Greek drinkers” .

[↑ Back](#)

9. Incorporate complex survey data features.

It is usually difficult to draw a simple random sample from the population due to cost and practical considerations such as a comprehensive sampling frame being available. As discussed in [Sample Design](#), complex samples, such as surveys involving stratified/cluster sample design, are commonly used in surveys. In a simple random sample, one can assume observations are independent from one another. However, in a complex sample design such as multi-stage samples (e.g., schools, classes, and students), students from one classroom are likely to be more correlated than those from another classroom. Therefore, as described in [Sample Design](#), in the analysis phase, we need to compensate for complex survey designs with features including, but not limited to, unequal likelihoods of selection, differences in response rates across subgroups, and deviations from distributions on critical variables found in the target population from external sources (e.g., as a national census), most commonly through the development of survey weights for statistical adjustment. If complex sample designs are implemented in data collection but the analysis assumes simple random sampling, the variances of survey estimates can be underestimated, and the confidence interval and test statistics are likely to be biased .

In a recent meta-analysis of 150 sampled research papers analyzing several surveys with complex sampling designs, researchers found that analytic errors caused by ignorance or incorrect use of the complex sample design features were frequent. These analytic errors define an important component of the larger total survey error framework, produce misleading descriptions of populations, and ultimately yield misleading inferences . It is thus of critical importance to incorporate the complex design features in statistical analysis.

For many of the aforementioned statistical models, various statistical software programs, such as the 'svy' statement and SURVEY procedures in SAS, have enabled the analysis of complex survey data features. See [Appendix A](#) for more information.

9.1 Suggested reading:

-
-
-
-
-

[↑ Back](#)

10. Introduction to Bayesian inference.

This section presents an overview of Bayesian theory, which follows closely the overview of , , and . In surveys, respondents' answers, denoted as $\{y\}$, reflect our measure of the true population's $\{Y\}$ —a random variable that takes a realized value $\{y\}$. In other words, $\{Y\}$ is unobserved, and the probability distribution $\{Y\}$ is of researchers' interest. We use $\{\theta\}$ to denote a parameter which reflects the characteristics of the distribution of $\{Y\}$. For example, $\{\mu\}$ can be the mean of the distribution. The goal is to estimate the unknown parameter $\{\theta\}$ based on the data, which is $\{p(\theta|y)\}$. Based on Bayes' theorem ($\{p(\theta|y)\} = \frac{\{p(\theta,y)\}}{\{p(y)\}} = \frac{\{p(y|\theta)p(\theta)\}}{\{p(y)\}}$), $\{p(y)\}$ is the probability distribution of the data, which is known for researchers; $\{p(y|\theta)\}$ refers to the probability of the data given the unknown parameter $\{\theta\}$; and $\{p(\theta)\}$ is the *prior distribution* of the parameters. $\{p(\theta|y)\}$ is thus referred to as the *posterior distribution* of the parameter $\{\theta\}$ given the data, which is also the result of the

In summary, Bayesian methods use both the prior information (which indicates the distribution of parameters) and the distribution of data to estimate the model results—the posterior distributions of the parameters. The key difference between Bayesian and frequentist approaches relates to the unknown parameter $\{\theta\}$. In the frequentist approach, $\{\theta\}$ is viewed as unknown but fixed. On the other hand, in the Bayesian approach, $\{\theta\}$ is random, which has a posterior distribution taking into account the uncertainty of $\{\theta\}$.

10.1 There are generally two types of priors: **noninformative** and **informative priors**. The choice between the two depends on our confidence about how much information we have about the priors and how accurate they are. Noninformative priors are also referred to as 'vague' or 'diffuse' priors. They are used when there is little information about the priors, and thus their influence on the posterior distribution of $\{\theta\}$ is minimal. An example of a noninformative prior can be a density with a huge variance, so that the Bayesian estimation is mainly affected by the data. Informative priors are used when we have sufficient information about the priors, such as from expert knowledge and similar datasets.

10.2 There are multiple Bayesian model comparison statistics. Two of the most commonly used are the **Bayes factor** and the **deviance information criterion (DIC)**. The Bayes factor quantifies the odds that the data favor one hypothesis over another. As discussed in [Guideline 5](#), Bayes factors are not well defined when using noninformative priors, and the evaluations can be computationally difficult. DIC compromises both goodness of fit and model complexity. In practical applications, the model with the smaller DIC value is preferred.

10.3 When we have estimated the posterior distributions of the parameters, we would like summaries of the distributions (such as mean and variance) for hypothesis testing. One important way to evaluate the distribution is based on the **credible interval**, which is often considered to be a similar measure to the 'confidence interval' in the frequentist approach. A credible interval is based on the quantiles of the posterior distributions. Based on the quantiles, we can evaluate the probability that the parameter lies in a particular interval. When this probability is 0.95, it is referred to as a 95% credible interval. If the credible intervals from two models do not overlap, we say that the two posterior distributions of this parameter differ.

10.4 **Markov chain Monte Carlo (MCMC)** is the most common computational algorithm for Bayesian methods. It generates Markov chains, which simulate the posterior distribution. The basic idea is that by simulating a sufficiently large number of observations from the posterior distribution, $\{p(\theta|y)\}$, we can approximate the mean and other summary statistics of the distribution. The use of MCMC for posterior simulation in latent variable models is particularly useful for the latent variables as missing data, which enables the augmentation of the observed variables. The most common MCMC algorithm is the Gibbs sampler, which performs on alternating conditional sampling at each of its iterations. More specifically, it draws each component conditional on the values of all the other components. In a Markov chain, a small proportion of the chain which may not converge to target distribution is called burn-in.

10.5 Multiple **convergence diagnostics** exist. In practice, it is common to inspect several different diagnostics, since there is no single adequate assessment. One of the most common statistics in a multiple chain condition is the **trace plot** diagnostic, which compares the within-chain and between-chain variance. A value above 1.1 is an indication of lack of convergence. The common diagnostics for single chain condition include the convergence diagnostic and the

convergence diagnostic, which can help to decide how many iterations are needed and how many can be treated as burn-in in a long enough chain.

10.6 Suggested reading:

-
-
-
-
-

10.7 Potential uses in 3MC research:

- As previously mentioned in [Guideline 6](#), approximate Bayesian measurement equivalence approach can be used in cross-cultural comparison research (e.g., or). See also the suggested readings in [Guideline 9.1](#) above for more information.

[↑ Back](#)

References

[↑ Back](#)