# Factor Analysis

## Overview

*Factor analysis*
is used to uncover the latent structure (dimensions) of a set of variables. It reduces attribute space from a larger number of variables to a smaller number of factors and as such is a "non-dependent" procedure (that is, it does not assume a dependent variable is specified). Factor analysis could be used for any of the following purposes:

- To reduce a large number of variables to a smaller number of factors for modeling purposes, where the large number of variables precludes modeling all the measures individually. As such, factor analysis is integrated in structural equation modeling (SEM), helping confirm the latent variables modeled by SEM. However, factor analysis can be and is often used on a stand-alone basis for similar purposes.

- To select a subset of variables from a larger set, based on which original variables have the highest correlations with the principal component factors.

- To create a set of factors to be treated as uncorrelated variables as one approach to handling multicollinearity in such procedures as multiple regression

- To validate a scale or index by demonstrating that its constituent items load on the same factor, and to drop proposed scale items which cross-load on more than one factor.

- To establish that multiple tests measure the same factor, thereby giving justification for administering fewer tests.

- To identify clusters of cases and/or outliers.

- To determine network groups by determining which sets of people cluster together (using Q-mode factor analysis, discussed below)

A non-technical analogy: A mother sees various bumps and shapes under a blanket at the bottom of a bed. When one shape moves toward the top of the bed, all the other bumps and shapes move toward the top also, so the mother concludes that what is under the blanket is a single thing, most likely her child. Similarly, factor analysis takes as input a number of measures and tests, analogous to the bumps and shapes. Those that move together are considered a single thing, which it labels a factor. That is, in factor analysis the researcher is assuming that there is a "child" out there in the form of an underlying factor, and he or she takes simultaneous movement (correlation) as evidence of its existence. If correlation is spurious for some reason, this inference will be mistaken, of course, so it is important when conducting factor analysis that possible variables which might introduce spuriousness, such as anteceding causes, be included in the analysis and taken into account.

Factor analysis is part of the multiple general linear hypothesis (MLGH) family of procedures and makes many of the same assumptions as multiple regression: linear relationships, interval or near-interval data, untruncated variables, proper specification (relevant variables included, extraneous ones excluded), lack of high multicollinearity, and multivariate normality for purposes of significance testing. Factor analysis generates a table in which the rows are the observed raw indicator variables and the columns are the factors or latent variables which explain as much of the variance in these variables as possible. The cells in this table are factor loadings, and the meaning of the factors must be induced from seeing which variables are most heavily loaded on which factors. This inferential labeling process

can be fraught with subjectivity as diverse researchers impute different labels.

There are several different types of factor analysis, with the most common being principal components analysis (PCA). However, principal axis factoring (PAF), also called common factor analysis, is preferred for purposes of confirmatory factory analysis in structural equation modeling.

# Key Concepts and Terms

- **Exploratory factor analysis**
  (EFA) seeks to uncover the underlying structure of a relatively large set of variables. The researcher's *à priori* assumption is that any indicator may be associated with any factor. This is the most common form of factor analysis. There is no prior theory and one uses factor loadings to intuit the factor structure of the data.

- **Confirmatory factor analysis**
  (CFA) seeks to determine if the number of factors and the loadings of measured (indicator) variables on them conform to what is expected on the basis of pre-established theory. Indicator variables are selected on the basis of prior theory and factor analysis is used to see if they load as predicted on the expected number of factors. The researcher's *à priori* assumption is that each factor (the number and labels of which may be specified *à priori*) is associated with a specified subset of indicator variables. A minimum requirement of confirmatory factor analysis is that one hypothesize beforehand the number of factors in the model, but usually also the researcher will posit expectations about which variables will load on which factors (Kim and Mueller, 1978b: 55). The researcher seeks to determine, for instance, if measures created to represent a latent variable really belong together.

  There are two approaches to confirmatory factor analysis:

  > *The Traditional Method*. Confirmatory factor analysis can be accomplished through any general-purpose statistical package which supports factor analysis. Note that for SEM CFA one uses principle axis factoring (PAF) rather than principle components analysis (PCA) as the type of factoring. This method allows the researcher to examine factor loadings of indicator variables to determine if they load on latent variables (factors) as predicted by the researcher's model. This can provide a more detailed insight into the measurement model than can the use of single-coefficient goodness of fit measures used in the SEM approach. As such the traditional method is a useful analytic supplement to the SEM CFA approach when the measurement model merits closer examination.

  > *The SEM Approach*. Confirmatory factor analysis can mean the analysis of alternative measurement (factor) models using a structural equation modeling package such as AMOS or LISREL. While SEM is typically used to model causal relationships among latent variables (factors), it is equally possible to use SEM to explore CFA measurement models. This is done by removing from the model all straight arrows connecting latent variables, adding curved arrows representing covariance between every pair of latent variables, and leaving in the straight arrows from each latent variable to its indicator variables as well as leaving in the straight arrows from error and disturbance terms to their respective variables. Such a measurement model is run like any other model and is evaluated like other models, using goodness of fit measures generated by the SEM package.

  >> *Testing error in the measurement model*. Using SEM, the researcher can explore CFA models with or without the assumption of certain correlations among the error terms of the indicator variables. Such measurement error terms represent causes of variance due to unmeasured variables as well as random measurement error. Depending on theory, it may well be that the researcher should assume unmeasured causal variables will be shared by indicators or will correlate, and thus SEM testing may well be merited. That is, including correlated measurement error in the model tests the possibility that indicator variables correlate not just because of being caused by a common factor, but also due to common or correlated unmeasured variables. This possibility would be ruled out if

the fit of the model specifying uncorrelated error terms was as good as the model with correlated error specified. In this way, testing of the confirmatory factor model may well be a desirable validation stage preliminary to the main use of SEM to model the causal relations among latent variables.

- *Redundancy test of one-factor vs. multi-factor models*. Using SEM, the *redundancy test* is to use chi-square difference (discussed in the section on <u>structural equation modeling)</u> to compare an original multifactor model with one which is constrained by forcing all correlations among the factors to be 1.0. If the constrained model is not significantly worse than the unconstrained one, the researcher concludes that a one-factor model would fit the data as well as a multi-factor one and, on the principle of parsimony, the one-factor model is to be preferred.

- *Measurement invariance test comparing a model across groups*. Using SEM, the *measurement invariance test* is to use chi-square difference to assess whether a set of indicators reflects a latent variable equally well across groups in the sample. The constrained model is one in which factor loadings are specified to be equal for each class of the grouping variable. If the constrained model is not significantly worse, then the researcher concludes the indicators are valid across groups. This procedure is also called *multiple group CFA*. If the model fails this test, then it is necessary to examine each indicator for group invariance, since some indicators may still be invariant. This procedure, called the *partial measurement invariance test* is discussed by Kline (1998: 225 ff.). Note that because standard errors of factor loadings cannot be computed, there are <u>indirect methods</u> but no direct method for comparing models across groups and hence the need for the SEM approach.

- *Orthogonality tests*. Using SEM, the *orthogonality test* is similar to the redundancy test, but factor correlations are set to 0. If the constrained model is not significantly worse than the unconstrained one, the factors in the model can be considered orthogonal (uncorrelated, independent). This test requires at least three indicators per factor.

- **Factors and components:**
  Both are the dimensions (or latent variables) identified with clusters of variables, as computed using factor analysis. Technically speaking, *factors* (as from PFA -- principal factor analysis, a.k.a. principal axis factoring, a.k.a. common factor analysis) represent the <u>common</u> variance of variables, excluding unique variance, and is thus a correlation-focused approach seeking to reproduce the intercorrelation among the variables. By comparison, *components* (from PCA - principal components analysis) reflect <u>both</u> common and unique variance of the variables and may be seen as a variance-focused approach seeking to reproduce both the total variable variance with all components and to reproduce the correlations. PCA is far more common than PFA, however, and it is common to use "factors" interchangeably with "components."

  PCA is generally used when the research purpose is data reduction (to reduce the information in many measured variables into a smaller set of components). PFA is generally used when the research purpose is to identify latent variables which contribute to the common variance of the set of measured variables, excluding variable-specific (unique) variance.

  *Warning:* Simulations comparing factor analysis with <u>structural equation modeling</u> (SEM) using simulated data indicate that at least in some circumstances, factor analysis may not correctly identify the correct number of latent variables, or sometimes even come close. While factor analysis may demonstrate that a particular model with a given predicted number of latent variables is not inconsistent with the data by this technique, researchers should understand that other models with different numbers of latent variables may also have good fit by SEM techniques.

- **Types of Factoring**

  There are different methods of extracting the factors from a set of data. The method chosen will matter more to the extent that the sample is small, the variables are few, and/or the communality estimates of the variables

differ.

- **Principal components analysis (PCA):**
  By far the most common form of factor analysis, PCA seeks a linear combination of variables such that the maximum variance is extracted from the variables. It then removes this variance and seeks a second linear combination which explains the maximum proportion of the remaining variance, and so on. This is called the principal axis method and results in orthogonal (uncorrelated) factors. PCA analyzes total (common and unique) variance.

  - **SPSS procedure:**
    Select Analyze - Data Reduction - Factor - Variables (input variables) - Descriptives - Under Correlation Matrix, check KMO and Anti-image to get overall and individual KMO statistics - Extraction - Method (principal components) and Analyze (correlation matrix) and Display (Scree Plot) and Extract (eigenvalues over 1.0) - Continue - Rotation - under Method, choose Varimax - Continue - Scores - Save as variables - Continue - OK.

  - **Canonical factor analysis**, also called *Rao's canonical factoring*, is a different method of computing the same model as PCA, which uses the principal axis method. CFA seeks factors which have the highest canonical correlation with the observed variables. CFA is unaffected by arbitrary rescaling of the data.

- **Principal factor analysis (PFA):**
  Also called principal axis factoring, PAF, and common factor analysis, PFA is a form of factor analysis which seeks the least number of factors which can account for the common variance (correlation) of a set of variables, whereas the more common principal components analysis (PCA) in its full form seeks the set of factors which can account for all the common <u>and</u> unique (specific plus error) variance in a set of variables. PFA uses a PCA strategy but applies it to a correlation matrix in which the diagonal elements are not 1's, as in PCA, but iteratively-derived estimates of the communalities ($R^2$ of a variable using all factors as predictors; see below).

  - **PFA and SEM:**
    PFA is preferred for purposes of modeling, as in structural equation modeling (SEM). PFA accounts for the covariation among variables, whereas PCA accounts for the total variance of variables. Because of this difference, in theory it is possible under PFA but not under PCA to add variables to a model without affecting the factor loadings of the original variables in the model. Widaman (1993) notes, "principal component analysis should not be used if a researcher wishes to obtain parameters reflecting latent constructs or factors." However, when commonalities are similar under PFA and PCA, then similar results will follow.

  - **PCA vs. PFA.** For most datasets, PCA and PFA will lead to similar substantive conclusions (Wilkinson, Blank, and Gruber, 1996). PCA is generally preferred for purposes of data reduction (translating variable space into optimal factor space), while PFA is generally preferred when the research purpose is detecting data structure or causal modeling.

- **Other Extraction Methods**. In addition to PCA and PFA, there are other less-used extraction methods:

  1. *Image factoring*: based on the correlation matrix of predicted variables rather than actual variables, where each variable is predicted from the others using multiple regression.

  2. *Maximum likelihood factoring*: based on a linear combination of variables to form factors, where the parameter estimates are those most likely to have resulted in the observed correlation matrix, using MLE methods and assuming multivariate normality. Correlations are weighted by each variable's uniqueness. (As discussed below, uniqueness is the variability of a variable minus its communality.) MLF generates a chi-square goodness-of-fit test. The researcher can increase the number of factors one at a time until a satisfactory goodness of fit is obtained. Warning: for large samples, even very small improvements in explaining variance can be significant by the

goodness-of-fit test and thus lead the researcher to select too many factors.

3. *Alpha factoring*: based on maximizing the reliability of factors, assuming variables are randomly sampled from a universe of variables. All other methods assume cases to be sampled and variables fixed.

4. *Unweighted least squares (ULS) factoring*: based on minimizing the sum of squared differences between observed and estimated correlation matrices, not counting the diagonal.

5. *Generalized least squares (GLS) factoring*: based on adjusting ULS by weighting the correlations inversely according to their uniqueness (more unique variables are weighted less). Like MLF, GLS also generates a chi-square goodness-of-fit test. The researcher can increase the number of factors one at a time until a satisfactory goodness of fit is obtained.

- **Factor Analytic Data Modes**

  ○ **R-mode factor analysis**. R-mode is by far the most common, so much so that it is normally assumed and not labeled as such. In R-mode, rows are cases, columns are variables, and cell entries are scores of the cases on the variables. In R-mode, the factors are clusters of variables on a set of people or other entities, at a given point of time.

  ○ **Q-mode factor analysis**, also called *inverse factor analysis*, is factor analysis which seeks to cluster the cases rather than the variables. That is, in Q-mode the rows are variables and the columns are cases (ex., people), and the cell entries are scores of the cases on the variables. In Q-mode the factors are clusters of people for a set of variables. Q-mode is used to establish the factional composition of a group on a set of issues at a given point in time.

    A Q-mode issue has to do with negative factor loadings. In conventional factor analysis of variables, loadings are loadings of variables on factors and a negative loading indicates a negative relation of the variable to the factor. In Q-mode factor analysis, loadings are loadings of cases (often individuals) on factors and a negative loading indicates that the case/individual displays responses opposite to those who load positively on the factor. In conventional factor analysis, loading approaching zero indicates the given variable is unrelated to the factor. In Q-mode factor analysis, a loading approaching zero indicates the given case is near the mean for the factor. Cluster analysis is now more common than Q-mode factor analysis.

    *The following modes are rare.*

  ○ **O-mode factor analysis**
    is an older form of time series analysis in which data are collected on a single entity (ex., one U. S. Senator), the columns are years, and the rows are measures (variables). In this mode, factors show which years cluster together on a set of measures for a single entity. Based on this, one can compare entities or, in a history of the entity, one can differentiate periods for purposes of explanation of behavior.

  ○ **T-mode factor analysis**
    is similar to O-mode in that the columns are years. However, the rows are entities (ex., cases are people) and data are gathered for a single variable. In T-mode, the factors show which years cluster together on that variable for a set of people or other entities. One might investigate, for instance, if Senators' positions on military spending are differentiated between war years and peacetime years.

  ○ **S-mode factor analysis**
    uses entities for columns (ex., Senators), years for rows (cases), and cell entries measure a single variable. In S-mode, factors show which Senators or other entities cluster together over a period of years on a single variable. S-mode would be used, for instance, to establish the underlying factional composition of a group on an issue over a long period of time.

- **Factor loadings:**

The factor loadings, also called component loadings in PCA, are the correlation coefficients between the variables (rows) and factors (columns). Analogous to Pearson's r, the squared factor loading is the percent of variance in that variable explained by the factor. To get the percent of variance in all the variables accounted for by each factor, add the sum of the squared factor loadings for that factor (column) and divide by the number of variables. (Note the number of variables equals the sum of their variances as the variance of a standardized variable is 1.) This is the same as dividing the factor's eigenvalue by the number of variables.

- *Factor, component, pattern, and structure matrices*. In SPSS, the factor loadings are found in a matrix labeled Factor Matrix if PFA is requested, or in one labeled Component Matrix if PCA is requested. (Note SPSS output gives both a factor or component matrix and a rotated factor or component matrix. The rotated version is used to induce factor meanings).

- In *oblique rotation*, one gets both a pattern matrix and a structure matrix. The *structure matrix* is simply the factor loading matrix as in orthogonal rotation, representing the variance in a measured variable explained by a factor on both a unique and common contributions basis. The *pattern matrix*, in contrast, contains coefficients which just represent unique contributions. The more factors, the lower the pattern coefficients as a rule since there will be more common contributions to variance explained. For oblique rotation, the researcher looks at <u>both</u>
  the structure and pattern coefficients when attributing a label to a factor.

- *The sum of the squared factor loadings*
  for all factors for a given variable (row) is the variance in that variable accounted for by all the factors, and this is called the *communality*. In a complete PCA, with no factors dropped, this will be 1.0, or 100% of the variance. The ratio of the squared factor loadings for a given variable (row in the factor matrix) shows the relative importance of the different factors in explaining the variance of the given variable. <u>Factor loadings are the basis for imputing a label to the different factors.</u>

- **Communality, $h^2$,** is the *squared multiple correlation* for the variable as dependent using the factors as predictors. The communality measures the percent of variance in a given variable explained by all the factors jointly and may be interpreted as the *reliability of the indicator*.

  - *Low communality*. When an indicator variable has a low communality, the factor model is not working well for that indicator and possibly it should be removed from the model. Low communalities across the set of variables indicates the variables are little related to each other. However, communalities must be interpreted in relation to the interpretability of the factors. A communality of .75 seems high but is meaningless unless the factor on which the variable is loaded is interpretable, though it usually will be. A communality of .25 seems low but may be meaningful if the item is contributing to a well-defined factor. That is, what is critical is not the communality coefficient per se, but rather the extent to which the item plays a role in the interpretation of the factor, though often this role is greater when communality is high.

  - *Spurious solutions*. If the communality exceeds 1.0, there is a spurious solution, which may reflect too small a sample or the researcher has too many or too few factors.

  - *Computation*. Communality for a variable is computed as the sum of squared factor loadings for that variable (row). Recall r-squared is the percent of variance explained, and since factors are uncorrelated, the squared loadings may be added to get the total percent explained, which is what communality is. For full orthogonal PCA, the initial communality will be 1.0 for all variables and all of the variance in the variables will be explained by all of the factors, which will be as many as there are variables. The "extracted" communality is the percent of variance in a given variable explained by the factors which are extracted, which will usually be fewer than all the possible factors, resulting in coefficients less than 1.0. For PFA and other extraction methods, however, the communalities for the various factors will be less than 1 even initially. Communality does not change when rotation is carried out, hence in SPSS there is only one communalities table.

- **Uniqueness** of a variable is $1 - h^2$. That is, uniqueness is the variability of a variable minus its communality.

- **Eigenvalues:** Also called *characteristic roots*. The eigenvalue for a given factor measures the variance in all the variables which is accounted for by that factor. The ratio of eigenvalues is the ratio of explanatory importance of the factors with respect to the variables. If a factor has a low eigenvalue, then it is contributing little to the explanation of variances in the variables and may be ignored as redundant with more important factors.

    - *Interpretation*. Eigenvalues measure the amount of variation in the total sample accounted for by each factor. Note that the eigenvalue is <u>not</u>
    the percent of variance explained but rather a measure of amount of variance in relation to total variance (since variables are standardized to have means of 0 and variances of 1, total variance is equal to the number of variables). SPSS will output a corresponding column titled '% of variance'. A factor's eigenvalue may be computed as the sum of its squared factor loadings for all the variables.

    - *Extraction sums of squared loadings*. Initial eigenvalues and eigenvalues after extraction (listed by SPSS as "Extraction Sums of Squared Loadings") are the same for PCA extraction, but for other extraction methods, eigenvalues after extraction will be lower than their initial counterparts. SPSS also prints "Rotation Sums of Squared Loadings" and even for PCA, these eigenvalues will differ from initial and extraction eigenvalues, though their total will be the same.

- **Trace**
  is the sum of variances for all factors, which is equal to the number of variables since the variance of a standardized variable is 1.0. A factor's eigenvalue divided by the trace is the percent of variance it explains in all the variables, usually labeled *percent of trace*
  in computer output. Computer output usually lists the factors in descending order of eigenvalue, along with a cumulative percent of trace for as many factors as are extracted.

- **Factor scores:** Also called *component scores*
  in PCA, factor scores are the scores of each case (row) on each factor (column). To compute the factor score for a given case for a given factor, one takes the case's standardized score on each variable, multiplies by the corresponding factor loading of the variable for the given factor, and sums these products. Computing factor scores allows one to look for factor outliers. Also, factor scores may be used as variables in subsequent modeling.

    - *SPSS*. The SPSS FACTOR procedure saves standardized factor scores as variables in your working data file. In SPSS, click Scores; select 'Save as Variables' and 'Display factor score coefficient matrix'. The factor (or in PCA, component) score coefficient matrix contains the regression coefficients used down the columns to compute scores for cases, were one to want to do this manually. By default SPSS will name them FAC1_1,FAC2_1, FAC3_1, etc., for the corresponding factors (factor 1, 2 and 3) of analysis 1; and FAC1_2, FAC2_2, FAC3_2 for a second set of factor scores, if any, within the same procedure, and so on. Although SPSS adds these variables to the right of your working data set automatically, they will be lost when you close the dataset unless you re-save your data.

- **Criteria for determining the number of factors**

    - *Comprehensibility*. Though not a strictly mathematical criterion, there is much to be said for limiting the number of factors to those whose dimension of meaning is readily comprehensible. Often this is the first two or three. Using one or more of the methods below, the researcher determines an appropriate range of solutions to investigate. For instance, the Kaiser criterion may suggest three factors and the scree test may suggest 5, so the researcher may request 3-, 4-, and 5-factor solutions and select the solution which generates the most comprehensible factor structure.

    - *Kaiser criterion:*
    A common rule of thumb for dropping the least important factors from the analysis is the K1 rule. Though originated earlier by Guttman in 1954, the criterion is usually referenced in relation to Kaiser's 1960 work which relied upon it. The Kaiser rule is to drop all components with eigenvalues under 1.0. It may overestimate or underestimate the true number of factors; the preponderance of simulation study

evidence suggests it usually overestimates the true number of factors, sometimes severely so (Lance, Butts, and Michels, 2006). The Kaiser criterion is the default in SPSS and most computer programs but is not recommended when used as the sole cut-off criterion for estimated the number of factors.

- ○ *Scree plot:*
  The Cattell scree test plots the components as the X axis and the corresponding eigenvalues as the Y axis. As one moves to the right, toward later components, the eigenvalues drop. When the drop ceases and the curve makes an elbow toward less steep decline, Cattell's scree test says to drop all further components after the one starting the elbow. This rule is sometimes criticised for being amenable to researcher-controlled "fudging." That is, as picking the "elbow" can be subjective because the curve has multiple elbows or is a smooth curve, the researcher may be tempted to set the cut-off at the number of factors desired by his or her research agenda.Researcher bias may be introduced due to the subjectivity involved in selecting the elbow. The scree criterion may result in fewer or more factors than the Kaiser criterion. [Scree plot example](#)

- ○ *Parallel analysis (PA)*, also known as *Humphrey-Ilgen parallel analysis*. PA is now often recommended as the best method to assess the true number of factors (Velicer, Eaton, and Fava, 2000: 67; Lance, Butts, and Michels, 2006). PA selects the factors which are greater than random. The actual data are factor analyzed, and separately one does a factor analysis of a matrix of random numbers representing the same number of cases and variables. For both actual and random solutions, the number of factors on the x axis and cumulative eigenvalues on the y axis is plotted. Where the two lines intersect determines the number of factors to be extracted. Though not available directly in SPSS or SAS, O'Connor (2000) presents programs to implement PA in SPSS, SAS, and MATLAB. These programs are located at http://flash.lakeheadu.ca/~boconno2/nfactors.html.

- ○ *Minimum average partial (MAP) criterion*. Developed by Velicer, this criterion is similar to PA in good resuls, but more complex to implement. O'Connor (2000), linked above, also presents programs for MAP.

- ○ *Variance explained criteria:*
  Some researchers simply use the rule of keeping enough factors to account for 90% (sometimes 80%) of the variation. Where the researcher's goal emphasizes parsimony (explaining variance with as few factors as possible), the criterion could be as low as 50%.

- ○ *Joliffe criterion:*
  A less used, more liberal rule of thumb which may result in twice as many factors as the Kaiser criterion. The Joliffe rule is to crop all components with eigenvalues under .7.

- ○ *Mean eigenvalue*. This rule uses only the factors whose eigenvalues are at or above the mean eigenvalue. This strict rule may result in too few factors.

  Before dropping a factor below one's cut-off, however, the researcher should check its correlation with the dependent variable. A very small factor can have a large correlation with the dependent variable, in which case it should not be dropped. Also, as a rule of thumb, factors should have at least three high, interpretable loadings -- fewer may suggest that the reasearcher has asked for too many factors.

- **Using reproduced correlation residuals to validate the choice of number of factors**

  - ○ *Reproduced correlations*
    is the correlation matrix of original items which would result on the assumption that the computed factors were the true and only factors. For any given pair of variables, the reproduced correlation is the product of their factor loadings on the first factor plus the product on the second factor, etc., for all factors. The diagonal values are the extracted communalities.

  - ○ *Reproduced correlation residuals*
    or "residual correlation matrix" is the matrix of differences between the reproduced and actual correlations. The closer the residuals are to zero (i.e., low or non-significant), the more confidence the

researcher has in his or her selection of the number of factors in the solution. In SPSS, footnotes to the table of residual correlations reports the percentage of non-redundant residual correlations greater than .05. In a good factor analysis, this percentage is low. (This is not a test used to reject a model.)

The reproduced correlation residuals matrix may help the researcher to identify particular correlations which are ill reproduced by the factor model with the current number of factors. By experimenting with different models with different numbers of factors, the researcher may assess which model best reproduces the correlations which are most critical to his or her research purpose.

- *In SPSS*, click the Descriptives button in the "Factor Analysis" dialog, then check "Reproduced" in the Correlations area of the "Factor Analysis: Descriptives" dialog. This option prints a table containing two subtables, the reproduced correlations on top and the reproduced correlation residuals on the bottom.

- **Rotation methods**. Rotation serves to make the output more understandable and is usually necessary to facilitate the interpretation of factors. The sum of eigenvalues is not affected by rotation, but rotation will alter the eigenvalues (and percent of variance explained) of particular factors and will change the factor loadings. Since alternative rotations may explain the same variance (have the same total eigenvalue) but have different factor loadings, and since factor loadings are used to intuit the meaning of factors, this means that different meanings may be ascribed to the factors depending on the rotation - a problem often cited as a drawback to factor analysis. If factor analysis is used, the researcher may wish to experiment with alternative rotation methods to see which leads to the most interpretable factor structure.

  Oblique rotations, discussed below, allow the factors to be correlated, and so a factor correlation matrix is generated when oblique is requested. Normally, however, an orthogonal method such as varimax is selected and no factor correlation matrix is produced as the correlation of any factor with another is zero.

  - *No rotation*
    is the default in SPSS, but it is a good idea to select a rotation method, usually varimax. The original, unrotated principal components solution maximizes the sum of squared factor loadings, efficiently creating a set of factors which explain as much of the variance in the original variables as possible. The amount explained is reflected in the sum of the eigenvalues of all factors. However, unrotated solutions are hard to interpret because variables tend to load on multiple factors.

  - *Varimax rotation*
    is an orthogonal rotation of the factor axes to maximize the variance of the squared loadings of a factor (column) on all the variables (rows) in a factor matrix, which has the effect of differentiating the original variables by extracted factor. Each factor will tend to have either large or small loadings of any particular variable. A varimax solution yields results which make it as easy as possible to identify each variable with a single factor. This is the most common rotation option.

  - *Quartimax rotation*
    is an orthogonal alternative which minimizes the number of factors needed to explain each variable. This type of rotation often generates a general factor on which most variables are loaded to a high or medium degree. Such a factor structure is usually not helpful to the research purpose.

  - *Equimax rotation* is a compromise between Varimax and Quartimax criteria.

  - *Direct oblimin rotation*
    is the standard method when one wishes a non-orthogonal (oblique) solution -- that is, one in which the factors are allowed to be correlated. This will result in higher eigenvalues but diminished interpretability of the factors. See below. See also hierarchical factor analysis.

  - *Promax rotation*
    is an alternative non-orthogonal (oblique) rotation method which is computationally faster than the direct oblimin method and therefore is sometimes used for very large datasets.

- **PRINCALS:**

A computer program which adapts PCA for non-metric and non-linear relationships. Its use is still rare.

- The **Component Transformation Matrix**
  in SPSS output shows the correlation of the factors before and after rotation.

# Assumptions

- **Valid imputation of factor labels**. Factor analysis is notorious for the subjectivity involved in imputing factor labels from factor loadings. For the same set of factor loadings, one researcher may label a factor "work satisfaction" and another may label the same factor "personal efficacy," for instance. The researcher may wish to involve a panel of neutral experts in the imputation process, though ultimately there is no "correct" solution to this problem.

- **No selection bias/proper specification**. The exclusion of relevant variables and the inclusion of irrelevant variables in the correlation matrix being factored will affect, often substantially, the factors which are uncovered. Although social scientists may be attracted to factor analysis as a way of exploring data whose structure is unknown, knowing the factorial structure in advance helps select the variables to be included and yields the best analysis of factors. This dilemma creates a chicken-and-egg problem. Note this is not just a matter of including all relevant variables. Also, if one deletes variables arbitrarily in order to have a "cleaner" factorial solution, erroneous conclusions about the factor structure will result. See Kim and Mueller, 1978a: 67-8.

- **No outliers**. Outliers can impact correlations heavily and thus distort factor analysis. One may use Mahalanobis distance to identify cases which are multivariate outliers, then remove them from the analysis prior to factor analysis. One can also create a dummy variable set to 1 for cases with high Mahalanobis distance, then regress this dummy on all other variables. If this regression is non-significant (or simply has a low R-squared for large samples) then the outliers are judged to be at random and there is less danger in retaining them. The ratio of the beta weights in this regression indicates which variables are most associated with the outlier cases.

- **Interval data**
  are assumed. However, Kim and Mueller (1978b 74-5) note that ordinal data may be used if it is thought that the assignment of ordinal categories to the data do not seriously distort the underlying metric scaling. Likewise, these authors allow use of dichotomous data if the underlying metric correlations between the variables are thought to be moderate (.7) or lower. The result of using ordinal data is that the factors may be that much harder to interpret.

  Note that categorical variables with similar splits will necessarily tend to correlate with each other, regardless of their content (see Gorsuch, 1983). This is particularly apt to occur when dichotomies are used. The correlation will reflect similarity of "difficulty" for items in a testing context, hence such correlated variables are called *difficulty factors*. The researcher should examine the factor loadings of categorical variables with care to assess whether common loading reflects a difficulty factor or substantive correlation. Improper use of dichotomies can result in too many factors. See the discussion of <u>levels of data</u>.

- **Linearity**. Factor analysis is a linear procedure. Of course, as with multiple linear regression, nonlinear transformation of selected variables may be a pre-processing step. The smaller the sample size, the more important it is to screen data for linearity.

- **Multivariate normality**
  of data is required for related significance tests. PCA and PFA, significance testing apart, have no distributional assumptions. Note, however, that a less-used variant of factor analysis, maximum likelihood factor analysis, does assume multivariate normality. The smaller the sample size, the more important it is to screen data for normality. Moreover, as factor analysis is based on correlation (or sometimes covariance), both correlation and covariance will be attenuated when variables come from different underlying distributions (ex., a normal vs. a

bimodal variable will correlate less than 1.0 even when both series are perfectly co-ordered). Nonetheless, normality is not considered one of the critical assumptions of factor analysis. See further discussion in the FAQ section.

- **Homoscedasticity**. Since factors are linear functions of measured variables, homoscedasticity of the relationship is assumed. However, homoscedasticity is not considered a critical assumption of factor analysis.

- **Orthogonality (for PFA but not PCA)**: the unique factors should be uncorrelated with each other or with the common factors. Recall that PFA factors only the common variance, ignoring the unique variance. This is not an issue for PCA, which factors the total variance.

- **Underlying dimensions**
  shared by clusters of variables are assumed. If this assumption is not met, the "garbage in, garbage out" (GIGO) principle applies. Factor analysis cannot create valid dimensions (factors) if none exist in the input data. In such cases, factors generated by the factor analysis algorithm will not be comprehensible. Likewise, the inclusion of multiple definitionally-similar variables representing essentially the same data will lead to tautological results.

- **Moderate to moderate-high intercorrelations without multicollinearity** are not mathematically required, but applying factor analysis to a correlation matrix with only low intercorrelations will require for solution nearly as many principal components as there are original variables, thereby defeating the data reduction purposes of factor analysis. On the other hand, too high intercorrelations may indicate a multicollinearity problem and colinear terms should be combined or otherwise eliminated prior to factor analysis. KMO statistics may be used to address multicollinearity in a factor analysis, or data may first be screened using VIF or tolerance in regression. Some researchers require correlations > 3.0 to conduct factor analysis.

- **No perfect multicollinearity**. Singularity in the input matrix, also called an ill-conditioned matrix, arises when two or more variables are perfectly redundant. Singularity prevents the matrix from being inverted and prevents a solution.

- **Factor interpretations and labels** must have face validity and/or be rooted in theory. It is notoriously difficult to assign valid meanings to factors. A recommended practice is to have a panel not otherwise part of the research project assign one's items to one's factor labels. A rule of thumb is that at least 80% of the assignments should be correct.

- **Adequate sample size**. At a minimum, there must be more cases than factors.

## SPSS Output Example

- **Annotated SPSS Factor Analysis Output**

## Frequently Asked Questions

- **How many cases do I need to do factor analysis?**
- **How do I input my data as a correlation matrix rather than raw data?**
- **How many variables do I need to do factor analysis? The more, the better?**
- **What is "sampling adequacy" and what is it used for?**
- **Why is normality not required for factor analysis when it is an assumption of correlation, on which factor analysis rests?**
- **Is it necessary to standardize one's variables before applying factor analysis?**
- **Can you pool data from two samples together in factor analysis?**
- **How does *factor comparison* of the factor structure of two samples work?**

- **[Why is rotation of axes necessary?](#)**
- **[Why are the factor scores I get the same when I request rotation and when I do not?](#)**
- **[Why is oblique (non-orthogonal) rotation rare in social science?](#)**
- **[When should oblique rotation be used?](#)**
- **[What is hierarchical factor analysis and how does it relate to oblique rotation?](#)**
- **[How high does a factor loading have to be to consider that variable as a defining part of that factor?](#)**
- **[What is *simple factor structure*, and is the simpler, the better?](#)**
- **[How is factor analysis related to validity?](#)**
- **[What is the matrix of standardized component scores, and for what might it be used in research?](#)**
- **[What are the pros and cons of PFA compared to PCA?](#)**
- **[Why are my PCA results different in SAS compared to SPSS?](#)**
- **[How do I do Q-mode factor analysis of cases rather than variables?](#)**
- **[How else may I use factor analysis to identify clusters of cases and/or outliers?](#)**
- **[What do I do if I want to factor categorical variables?](#)**

- **How many cases do I need to do factor analysis?**
     There is no scientific answer to this question, and methodologists differ. Alternative arbitrary "rules of thumb," in descending order of popularity, include those below. These are not mutually exclusive: Bryant and Yarnold, for instance, endorse both STV and the Rule of 200. There is near universal agreement that factor analysis is inappropriate when sample size is below 50.
     1. Rule of 10. There should be at least 10 cases for each item in the instrument being used.
     2. STV ratio. The subjects-to-variables ratio should be no lower than 5 (Bryant and Yarnold, 1995)
     3. Rule of 100: The number of subjects should be the larger of 5 times the number of variables, or 100. Even more subjects are needed when communalities are low and/or few variables load on each factor. (Hatcher, 1994)
     4. Rule of 150: Hutcheson and Sofroniou (1999) recommends at least 150 - 300 cases, more toward the 150 end when there are a few highly correlated variables, as would be the case when collapsing highly multicollinear variables.
     5. Rule of 200. There should be at least 200 cases, regardless of STV (Gorsuch, 1983)
     6. Rule of 300. There should be at least 300 cases (Norušis, 2005: 400).
     7. Significance rule. There should be 51 more cases than the number of variables, to support chi-square testing (Lawley and Maxwell, 1971)

- **How do I input my data as a correlation matrix rather than raw data?**
     In SPSS, one first creates a "matrix data file" using the MATRIX DATA command, as explained in the *SPSS Syntax Reference Guide*. The format is:

     ```
     MATRIX DATA VARIABLES=varlist.
     BEGIN DATA
     MEAN    meanslist
     STDDEV  stddevlist
     CORR 1
     CORR .22  1
     CORR -.58  .88  1
     CORR  .33  .02   -.17   1
     END DATA.
     EXECUTE.
     ```

     where
     varlist is a list of variable names separated by commas
     meanslist is a list of the means of variables, in the same order as varlist
     stddevlist is a list of standard deviations of variables, in the same order
     CORR statements define a correlation matrix, with variables in the same order (data above are for illustration; one may have more or fewer CORR statements as needed according to the number of variables).
     Note the period at the end of the MATRIX DATA and END DATA commands.

Then if the MATRIX DATA command is part of the same control syntax working file, add the FACTOR command as usual but add the subcommand "/MATRIX=(IN(*)" (but without the quote marks). If the MATRIX DATA is not part of the same syntax set but has been run earlier, the matrix data file name is substituted for the asterisk.

- **How many variables do I need in factor analysis? The more, the better?**

    For confirmatory factor analysis, there is no specific limit on the number of variables to input. For exploratory factory analysis, Thurstone recommended at least three variables per factor (Kim and Mueller, 1978b: 77).

    Using confirmatory factor analysis in [structural equation modeling](#), having several or even a score of indicator variables for each factor will tend to yield a model with more reliability, greater validity, higher generalizability, and stronger tests of competing models, than will CFA with two or three indicators per factor, all other things equal. However, the researcher must take account of the statistical artifact that models with fewer variables will yield apparent better fit as measured by SEM goodness of fit coefficients, all other things equal.

    However, "the more, the better" may <u>not</u> be true when there is a possibility of *suboptimal factor solutions* ("bloated factors"). Too many too similar items will mask true underlying factors, leading to suboptimal solutions. For instance, items like "I like my office," "My office is nice," "I like working in my office," etc., may create an "office" factor when the researcher is trying to investigate the broader factor of "job satisfaction." To avoid suboptimization, the researcher should start with a small set of the most defensible (highest face validity) items which represent the range of the factor (ex., ones dealing with work environment, coworkers, and remuneration in a study of job satisfaction). Assuming these load on the same job satisfaction factor, the researcher then should add one additional variable at a time, adding only items which continue to load on the job satisfaction factor, and noting when the factor begins to break down. This stepwise strategy results in the most defensible final factors.

- **What is "sampling adequacy" and what is it used for?**

    Measured by the Kaiser-Meyer-Olkin (KMO) statistics, sampling adequacy predicts if data are likely to factor well, based on correlation and partial correlation. In the old days of manual factor analysis, this was extremely useful. KMO can still be used, however, to assess which variables to drop from the model because they are too multicollinear.

    There is a KMO statistic for each individual variable, and their sum is the KMO overall statistic. KMO varies from 0 to 1.0 and KMO overall should be .60 or higher to proceed with factor analysis. If it is not, drop the indicator variables with the lowest individual KMO statistic values, until KMO overall rises above .60. (Some researchers use a more lenient .50 cut-off).

    the To compute KMO overall, the numerator is the sum of squared correlations of all variables in the analysis (except the 1.0 self-correlations of variables with themselves, of course). The denominator is this same sum plus the sum of squared partial correlations of each variable i with each variable j, controlling for others in the analysis. The concept is that the partial correlations should not be very large if one is to expect distinct factors to emerge from factor analysis. See Hutcheson and Sofroniou, 1999: 224.

    In SPSS, KMO is found under Analyze - Statistics - Data Reduction - Factor - Variables (input variables) - Descriptives - Correlation Matrix - check KMO and Bartlett's test of sphericity and also check Anti-image - Continue - OK. The KMO output is KMO overall. The diagonal elements on the Anti-image correlation matrix are the KMO individual statistics for each variable.

- **Why is normality not required for factor analysis when it is an assumption of correlation, on which factor analysis rests?**

    Factor analysis is a correlative technique, seeking to cluster variables along dimensions, or it may be used to provide an estimate (factor scores) of a latent construct which is a linear combination of variables. The normality assumption pertains to significance testing of coefficients. If one is just interested in clustering (correlating) factors or developing factor scores, significance testing is beside the point. In correlation,

significance testing is used with random samples to determine which coefficients cannot be assumed to be different from zero. However, in factor analysis the issue of which variables to drop is assessed by identifying those with low communalities since these are the ones for which the factor model is "not working." Communalities are a form of effect size coefficient, whereas significance also depends on sample size. As mentioned in the assumptions section, it is still true that factor analysis also requires adequate sample size, in the absence of which the factor scores and communalities may be unreliable, and if variables come from markedly different underlying distributions, correlation and factor loadings will be attenuated as they will be for other causes of attenuation in correlation.

- **Is it necessary to standardize one's variables before applying factor analysis?**
  No. Results of factor analysis are not affected by standardization, which is built into the procedure. Note, however, that standardization (subtracting the mean, dividing by the standard deviation) scales data in a sample-specific way. If the research purpose is to compare factor structures between two or more samples, then one should use the covariance matrix rather than the correlation matrix as input to factor analysis (Kim and Mueller, 1978b: 76). However, the covariance method has problems when the variables are measured on widely different scales (ex., income measure to $100,000 and education measured to 22 years). Kim and Mueller recommend multisample standardization for this case (subtracting the grand mean of all samples, and dividing by the standard deviation of all cases) prior to computing the sample covariance matrix (p. 76). However, in practice factor comparison between samples usually is based on ordinary factor analysis of correlation matrices.

- **Can you pool data from two samples together in factor analysis?**
  Yes, but only after you have shown both samples have the same factor structure in through factor comparison.

- **How does *factor comparison* of the factor structure of two samples work?**

  The *pooled data method*
  has the researcher pool the data for two samples, adding a dummy variable whose coding represents group membership. The factor loadings of this dummy variable indicate the factors for which the groups' mean factor scores would be most different.

  The *factor invariance test*, discussed above, is a structural equation modeling technique (available in AMOS, for ex.) which tests for deterioration in model fit when factor loadings are constrained to be equal across sample groups.

  The *comparison measures method*
  requires computation of various measures which compare factor attributes of the two samples. Factor comparison is discussed by Levine (1977: 37-54), who describes these factor comparison measures:

  - *RMS, root mean square*. RMS is the root mean square of the average squared difference of the loadings of the variables on each of two factors. RMS varies from 0 to 2, reaching 0 in the case of a perfect match between samples of <u>both</u> the pattern and the magnitude of factors in the two samples. An RMS of 2 indicates all loadings are at unity but differ in sign between the two samples. Intermediate values are hard to interpret.

  - *CC, coefficient of congruence*. The coefficient of congruence is the sum of the products of the paired loadings divided by the square root of the product of the two sums of squared loadings. Like RMS, CC measures both pattern and magnitude similarities between samples. There is a tendency to get a high CC whenever two factors have many variables with the same sign.

  - *S, salient variable similarity index*. The salient variable similarity index is based on classifying factor loadings into positive salient ones (over +.1), hyperplane ones (from -.1 to +.1), and negative salient ones ((below -.1). Hyperplane loadings, which approach 0, indicate variables having only a near-chance relationship to the factor. The S index will be 0 when there are no salient loadings, indicating no factor congruence between the two samples. An S of 1 indicates perfect congruence,

and -1 indicates perfect negative (reflected) congruence. Note that calculating S for all possible pairs of factors between two samples risks coming to conclusions on the basis of chance.

- **Why is rotation of axes necessary?**

    For solutions with two or more factors, prior to rotation the first axis will lie in between the clusters of variables and in general the variables will not sort well on the factors. Rotation of the axes causes the factor loadings of each variable to be more clearly differentiated by factor.

- **Why are the factor scores I get the same when I request rotation and when I do not?**

    If the solution has only one factor, rotation will not be done so the factor scores will be the same whether you request rotation or not.

- **Why is oblique (non-orthogonal) rotation rare in social science?**

    Oblique rotation is often defended on the grounds that (1) it reflects the real world more accurately, since real-world constructs are often correlated; and (2) the resulting factor structure is simpler and more interpretable. However, oblique rotation is relatively rare precisely because the correlation of factors makes much more problematic the greatest difficulty of factor analysis, which is imputing the meaning of the factors from the factor loadings. In general, oblique solutions will result in more cross-loaded variables than will standard orthogonal rotations.

    However, occasionally an oblique rotation will still result in a set of factors whose intercorrelations approach zero. This, indeed, is the test of whether the underlying factor structure of a set of variables is orthogonal. Orthogonal rotation mathematically assures resulting factors w

    Also, oblique rotation is necessary as part of [hierarchical factor analysis](), which seeks to identify higher-order factors on the basis of correlated lower-level ones..

- **When <u>should</u> oblique rotation be used?**

    In confirmatory factor analysis (CFA), if theory suggests two factors are correlated, then this measurement model calls for oblique rotation. In exploratory factor analysis (EFA), the researcher does not have a theoretical basis for knowing how many factors there are or what they are, much less whether they are correlated. Researchers conducting EFA usually assume the measured variables are indicators of two or more *different*
    factors, a measurement model which implies orthogonal rotation. That EFA is far more common than CFA in social science is another reason why orthogonal rotation is far more common than oblique rotation.

    When modeling, oblique rotation may be used as a filter. Data are first analyzed by oblique rotation and the factor correlation matrix is examined. If the factor correlations are small (ex., < .32, corresponding to 10% explained), then the researcher may feel warranted in assuming orthogonality in the model. If the correlations are larger, then covariance between factors should be assumed (ex., in structural equation modeling, one adds double-headed arrows between latents).

    For purposes other than modeling, such as seeing if test items sort themselves out on factors as predicted, orthogonal rotation is almost universal.

- **What is hierarchical factor analysis and how does it relate to oblique rotation?**

    *Hierarchical factor analysis*
    (HFA) seeks to differentiate higher-order factors from a set of correlated lower-order factors. For instance, HFA has been used to support Spearman's "g factor" theory in the study of intelligence, where g is the highest-order common factor emerging from a hierarchical factor analysis of a large number of diverse cognitive tests of ability (Jensen, 1998). In this theory, the g factor is the higher-order factor underlying such correlated lower-order factors as arithmetic ability, verbal ability, and reasoning ability. Likewise, HFA has been used prominently in trait theory, in the debate over whether personality is best characterized in terms of three, five, or some other number of higher-order traits (see Matthews, Deary, & Whiteman, 2003).

HFA is a two-stage process. First an oblique (oblimin) factor analysis is conducted on the raw dataset. As it is critical in HFA to obtain the simplest factor structure possible, it is recommended to run oblimin for several different values of delta, not just the default delta=0. A delta of 0 gives the most oblique solutions, but the more the researcher specifies (in the SPSS "Factor Analysis" Rotation" dialog, invoked by clicking the Rotation button) a more and more negative delta, the factors become less and less oblique. To override the default delta of 0, the researcher enters a value less than or equal to 0.8.

When the researcher feels the simplest factor structure has been obtained, one has a correlated set of lower-order factors. Factor scores or a correlation matrix of factors from the first stage can be input to a second-stage orthogonal factor analysis (ex., varimax) to generate one or more higher-order factors.

- **How high does a factor loading have to be to consider that variable as a defining part of that factor?**
  This is purely arbitrary, but common social science practice uses a minimum cut-off of .3 or .35. Norman and Streiner (1994: 139) give this alternative formula for minimum loadings when the sample size, N, is 100 or more: Min FL = 5.152/[SQRT(N-2)]. Another arbitrary rule-of-thumb terms loadings as "weak" if less than .4, "strong" if more than .6, and otherwise as "moderate." These rules are arbitrary. The meaning of the factor loading magnitudes varies by research context. For instance, loadings of .45 might be considered "high" for dichotomous items but for Likert scales a .6 might be required to be considered "high."

- **What is *simple factor structure*, and is the simpler, the better?**
  A factor structure is simple to the extent that each variable loads heavily on one and only one factor. Usually rotation is necessary to achieve simple structure, if it can be achieved at all. Oblique rotation does lead to simpler structures in most cases, but it is more important to note that oblique rotations result in correlated factors, which are difficult to interpret. Simple structure is only one of several sometimes conflicting goals in factor analysis.

- **How is factor analysis related to validity?**
  In confirmatory factor analysis (CFA), a finding that indicators have high loadings on the predicted factors indicates *convergent validity*. In an oblique rotation, *discriminant validity* is demonstrated if the correlation between factors is not so high (ex., > ,85) as to lead one to think the two factors overlap conceptually.

- **What is the matrix of standardized component scores, and for what might it be used in research?**
  These are the scores of all the cases on all the factors, where cases are the rows and the factors are the columns. They can be used for orthogonalization of predictors in multiple regression. In a case where there is multicollinearity, one may use the component scores in place of the X scores, thereby assuring there is no multicollinearity of predictors.

  Note, however, that this orthogonalization comes at a price. Now, instead of explicit variables, one is modeling in terms of factors, the labels for which are difficult to impute. Statistically, multicollinearity is eliminated by this procedure, but in reality it is hidden in the fact that all variables have some loading on all factors, muddying the purity of meaning of the factors.

  A second research use for component scores is simply to be able to use fewer variables in, say, a correlation matrix, in order to simplify presentation of the associations.

  Note also that factor scores are quite different from factor loadings. Factor scores are coefficients of cases on the factors, whereas factor loadings are coefficients of variables on the factors.

- **What are the pros and cons of PFA compared to PCA?**
  PCA determines the factors which can account for the total (unique and common) variance in a set of variables. This is appropriate for creating a typology of variables or reducing attribute space. PCA is appropriate for most social science research purposes and is the most often used form of factor analysis.

  PFA determines the least number of factors which can account for the common variance in a set of

variables. This is appropriate for determining the dimensionality of a set of variables such as a set of items in a scale, specifically to test whether one factor can account for the bulk of the common variance in the set, though PCA can also be used to test dimensionality. PFA has the disadvantage that it can generate negative eigenvalues, which are meaningless.

- **Why are my PCA results different in SAS compared to SPSS?**
  There are different algorithms for computing PCA, yielding quite similar results. SPSS by used "iterated principal factors" by default, whereas SAS did not. In SAS, specify METHOD=PRINIT to get the iterated solution within PROC FACTOR. This assumes, of course, that the same rotation is also specified in both programs. Also, SPSS by default does 25 iterations and this could make a minor difference if SAS differs, though SPSS allows the user to change this to another number.

- **How do I do Q-mode factor analysis of cases rather than variables?**
  Simply transpose the data matrix, reversing rows and columns. Note that there must be more cases than variables.

- **How else may I use factor analysis to identify clusters of cases and/or outliers?**
  If there are only two or at most three principal component factors which explain most of the total variation in the original variables, then one can calculate the factor scores of all cases on these factors, and then a plot of the factor scores will visually reveal both clusters of cases and also outliers. See Dunteman, 1989: 75-79.

- **What do I do if I want to factor categorical variables?**
  A nominal and ordinal analog to factor analysis is *latent class analysis*. Also, SPSS offers offers other procedures for factoring categorical data.

# Bibliography

- Bryant and Yarnold (1995). Principal components analysis and exploratory and confirmatory factor analysis. In Grimm and Yarnold, *Reading and understanding multivariate analysis*. American Psychological Association Books.
- Dunteman, George H. (1989). *Principal components analysis*. Thousand Oaks, CA: Sage Publications, Quantitative Applications in the Social Sciences Series, No. 69.
- Fabrigar, L. R., Wegener, D. T., MacCallum, R. C., & Strahan, E. J. (1999). Evaluating the use of exploratory factor analysis in psychological research. *Psychological Methods*, 4: 272-299.
- Gorsuch, R. L. (1983). *Factor Analysis*. Hillsdale, NJ: Lawrence Erlbaum. Orig. ed. 1974.
- Hatcher, Larry (1994). *A step-by-step approach to using the SAS system for factor analysis and structural equation modeling*. Cary, NC: SAS Institute. Focus on the CALIS procedure.
- Hutcheson, Graeme and Nick Sofroniou (1999). *The multivariate social scientist: Introductory statistics using generalized linear models*. Thousand Oaks, CA: Sage Publications.
- Jensen, A. R. (1998). *The g factor: The science of mental ability.* Westport, CT: Praeger
- Kim, Jae-On and Charles W. Mueller (1978a). *Introduction to factor analysis: What it is and how to do it*. Thousand Oaks, CA: Sage Publications, Quantitative Applications in the Social Sciences Series, No. 13.
- Kim, Jae-On and Charles W. Mueller (1978b). *Factor Analysis: Statistical methods and practical issues*. Thousand Oaks, CA: Sage Publications, Quantitative Applications in the Social Sciences Series, No. 14.
- Kline, Rex B. (1998). *Principles and practice of structural equation modeling*. NY: Guilford Press. Covers confirmatory factor analysis using SEM techniques. See esp. Ch. 7.
- Lance, Charles E, Marcus M. Butts, and Lawrence C. Michels (2006). The sources of four commonly reported cutoff criteria: What did they really say? Organizational Research Methods 9(2): 202-220. Discusses Kaiser and other criteria for selecting number of factors.
- Lawley, D. N. and A. E. Maxwell (1971). *Factor analysis as a statistical method*. London: Butterworth and Co.
- Levine, Mark S. (1977). *Canonical analysis and factor comparison*. Thousand Oaks, CA: Sage Publications, Quantitative Applications in the Social Sciences Series, No. 6.

- Matthews, G., Deary, I. J., & Whiteman, M. C. (2003). *Personality traits, Second edition*. Cambridge: Cambridge University Press.
- Norman, G. R., and D. L. Streiner (1994). *Biostatistics: The bare essentials*. St. Louis, MO: Mosby.
- O'Connor, B. P. (2000). SPSS and SAS programs for determining the number of components using parallel analysis and Velicer's MAP test. *Behavior Research Methods, Instrumentation, and Computers* 32: 396-402. .
- Norušis. Marija J. (2005). SPSS 13.0 Statistical Procedures Companion. Chicago: SPSS, Inc.
- Pett, Marjorie A., Nancy R. Lackey, and John J. Sullivan (2003). *Making sense of factor analysis: The use of factor analysis for instrument development in health care research*. Thousand Oaks, CA: Sage Publications.
- Velicer, W. F., Eaton, C. A., and Fava, J. L. (2000). Construct explication through factor or component analysis: A review and evaluation of alternative procedures for determining the number of factors or components. Pp. 41-71 in R. D. Goffin and E. Helmes, eds., *Problems and solutions in human assessment*. Boston: Kluwer. Upholds PA over K1 as a number of factors cutoff criterion.
- Widaman, K. F. (1993). Common factor analysis versus principal components analysis: Differential bias in representing model parameters?" *Multivariate Behavioral Research* 28: 263-311. Cited with regard to preference for PFA over PCA in confirmatory factor analysis in SEM.
- Wilkinson, L., G. Blank, and C. Gruber (1996). *Desktop Data Analysis with SYSTAT*. Upper Saddle River, NJ: Prentice-Hall.

---

Copyright 1998, 2007 by G. David Garson.

---

Back

---