

Hypothesis testing for coefficient alpha: An SEM approach

ALBERTO MAYDEU-OLIVARES

University of Barcelona, Barcelona, Spain

DONNA L. COFFMAN

Pennsylvania State University, State College, Pennsylvania

AND

CARLOS GARCÍA-FORERO AND DAVID GALLARDO-PUJOL

University of Barcelona, Barcelona, Spain

We show how to test hypotheses for coefficient alpha in three different situations: (1) hypothesis tests of whether coefficient alpha equals a prespecified value, (2) hypothesis tests involving two statistically independent sample alphas as may arise when testing the equality of coefficient alpha across groups, and (3) hypothesis tests involving two statistically dependent sample alphas as may arise when testing the equality of alpha across time or when testing the equality of alpha for two test scores within the same sample. We illustrate how these hypotheses may be tested in a structural equation-modeling framework under the assumption of normally distributed responses and also under asymptotically distribution free assumptions. The formulas for the hypothesis tests and computer code are given for four different applied examples. Supplemental materials for this article may be downloaded from <http://brm.psychonomic-journals.org/content/supplemental>.

Assessing the reliability of a questionnaire or test score is one of the most frequent tasks in psychological research. Often, researchers wish to go beyond providing a point estimate of the reliability of their test score and are interested in testing hypotheses concerning the reliability of their test score. A typical situation is one in which a researcher wishes to determine whether the reliability of his or her test score is larger than some predetermined cutoff value (say .80). Another commonly encountered situation is one in which a researcher wishes to determine whether the reliability of his or her test score is equal across two or more populations. For instance, the researcher may wish to determine whether the reliability of a test score is equal across genders, or he or she may wish to determine whether the reliability of a test score is the same across several countries. Finally, sometimes researchers are interested in determining whether, within a population, the reliabilities of two test scores are equal. For instance, the researcher may wish to test whether the reliability of a test score based on p items equals the reliability of a test score based on a subset of those p items (such as when a full form and a short form of a questionnaire are available). A special case of this instance is when a researcher wishes to test whether the reliability changes when a single item is removed from the scale. As another example, a researcher may wish to test whether the scores based on two subsets of items drawn from the same item domain are equally reliable.

Most often, reliability assessment is performed by means of coefficient alpha (Hogan, Benjamin, & Brezinski, 2000). Coefficient alpha (α) was first proposed by Guttman (1945), with important contributions by Cronbach (1951). For some discussions on coefficient alpha, see Cortina (1993); Miller (1995); Schmitt (1996); Shevlin, Miles, Davies, and Walker (2000); and ten Berge (2000). Coefficient α is a population parameter and thus an unknown quantity. In applications, it is typically estimated using the sample coefficient alpha, a point estimator of the population coefficient alpha. As with any point estimator, sample coefficient alpha is subject to variability around the true parameter, particularly in small samples. Methods for performing hypothesis testing based on coefficient alpha rely on the estimation of the variability of sample coefficient alpha (i.e., its standard error). The initial proposals for estimating the standard error of coefficient alpha were based on model and distributional assumptions (see Duhachek & Iacobucci, 2004, for an overview). Thus, if a particular model held for the covariance matrix among the test items, and the test items followed a particular distribution, a confidence interval for coefficient alpha could be obtained. The sampling distribution for coefficient alpha was first derived (independently) by Kristof (1963) and Feldt (1965), who assumed that the test items were strictly parallel (see Lord & Novick, 1968) and normally distributed. This model implies that all of the item variances are equal and that all of the item covariances are

A. Maydeu-Olivares, amaydeu@ub.edu



equal. However, Barchard and Hakstian (1997) found that standard errors for coefficient alpha obtained using these results were not sufficiently accurate when model assumptions were violated (i.e., the items were not strictly parallel). The lack of robustness of the standard errors for coefficient alpha to violations of model assumptions may have hindered the widespread use of hypothesis tests for alpha in applications.

A major breakthrough occurred when van Zyl, Neu-decker, and Nel (2000) derived the asymptotic (i.e., large sample) distribution of sample coefficient alpha without model assumptions. In particular, van Zyl et al. assumed only that the items composing the test were normally distributed and that their covariance matrix is positively definitive. Hence, their approach is model free. All previous derivations, which assumed particular models (e.g., tau equivalence) for the covariance matrix, can be treated as special cases of van Zyl et al.'s result. Duhachek and Iacobucci (2004) compared the performance of these model-free standard errors for coefficient alpha with those of the model-based procedures proposed by Feldt (1965) and Hakstian and Whalen (1976) under violations of the model assumptions underlying coefficient alpha. They found that the model-free, normal theory (NT) interval estimator proposed by van Zyl et al. uniformly outperformed competing procedures across all conditions.

However, van Zyl et al. (2000) assumed that the items composing the test can be well approximated by a normal distribution. In practice, tests are most often composed of binary or Likert-type items for which the normal distribution can be a poor approximation. Yuan, Guarnaccia, and Hayslip (2003) have proposed a model-free and asymptotically distribution-free (ADF) standard error for sample coefficient alpha that overcomes this limitation. Maydeu-Olivares, Coffman, and Hartmann (2007) showed that for sample sizes over 100 observations, ADF standard errors are preferable to NT standard errors, because the latter are not sufficiently accurate when the skewness or excess kurtosis of the items is larger than 1.

The aim of this study is to show how hypotheses for coefficient alpha can be tested using the NT results of van Zyl et al. (2000) and also using the ADF results of Yuan et al. (2003), using a structural equations modeling (SEM) framework. In particular, we show how to perform hypothesis testing in three cases. Case 1 involves a single sample alpha. Case 2 involves two statistically independent sample alphas. Case 3 involves two statistically dependent sample alphas. Case 1 arises when testing whether the population alpha exceeds some predetermined cutoff value. Case 2 arises when comparing the population alpha across two independent samples, such as male and female subjects or across countries. Case 3 arises when comparing the population alpha for two sets of items in a single sample. Typical examples of Case 3 are testing the equality of population alpha when an item is dropped, testing the equality of alpha for a full-scale score and a reduced-scale score, and testing the equality of alpha for the same score measured at two time points. An SEM framework is not needed for testing the Case 1 and 2 hypotheses. Indeed, the formulae involved are straightforward. All that is needed are the stan-

dard errors of sample alpha, which can be computed using the code provided by Duhachek and Iacobucci (2004) and Maydeu-Olivares et al. (2007). However, adopting an SEM framework for Cases 1 and 2 is convenient, because it provides a link to Case 3 hypothesis testing, which cannot be easily performed without using an SEM framework. Also, by adopting an SEM framework, we can integrate the results of van Zyl et al. and Yuan et al. with the large literature on reliability assessment using SEM.

The three cases considered are illustrated using four examples. The test statistics discussed in the present article are based on large-sample theory and may not be accurate in small samples. Since it is questionable to present results using arbitrary parameter values and to draw generalizable conclusions from them, we show how a Monte Carlo investigation can be performed using the simulation capabilities of SEM packages to determine the accuracy of the obtained p values, and we do so for each of the examples presented.

Coefficient Alpha

Consider a test or questionnaire composed of p items, Y_1, \dots, Y_p , intended to measure a single attribute. The reliability of the test score, $X = Y_1 + \dots + Y_p$, is defined as the percentage of variance of X that is due to the attribute of which the items are indicators. The most widely used procedure to assess the reliability of X is coefficient alpha (Cronbach, 1951; Guttman, 1945). In the population of respondents, coefficient alpha is

$$\alpha = \frac{p}{p-1} \left(1 - \frac{\sum_i \sigma_{ii}}{\sum_i \sigma_{ii} + 2 \sum_{i<j} \sigma_{ij}} \right), \tag{1}$$

where $\sum_i \sigma_{ii}$ denotes the sum of the p item variances in the population, and $\sum_{i<j} \sigma_{ij}$ denotes the sum of the $p(p-1)/2$ distinct item covariances. In applications, a sample of N respondents from the population is procured, and a point estimator of the population alpha given in Equation 1 can be obtained using the sample coefficient alpha:

$$\hat{\alpha} = \frac{p}{p-1} \left(1 - \frac{\sum_i s_{ii}}{\sum_i s_{ii} + 2 \sum_{i<j} s_{ij}} \right), \tag{2}$$

where s_{ij} denotes the sample covariance between items i and j , and s_{ii} denotes the sample variance of item i .

Coefficient alpha and the reliability of a test score.

If the items satisfy a true-score equivalent model (a.k.a. an essentially tau-equivalent model), coefficient alpha equals the reliability of the test score (Lord & Novick, 1968, p. 50; McDonald, 1999, chap. 6). A true-score equivalent model is a one-factor model in which the factor loadings are equal for all items. The model implies that the population covariances are all equal, but the population variances need not be equal for all items. Coefficient alpha also equals the reliability of the test score when the items are strictly parallel, because these are special cases of the true-score equivalent model. In the parallel items model, in addition to the assumptions of the true-score equivalent model, the unique variances of the error terms

in the factor model are assumed to be equal for all items. A more constrained version of the parallel items model is the strictly parallel items model, in which the item means are additionally assumed to be equal across items.

When the items do not conform to a true-score equivalent model, coefficient alpha does not equal the reliability of the test score. For instance, if the items conform to a one-factor model with distinct factor loadings (i.e., congeneric items), the reliability of the test score is given by coefficient omega (see McDonald, 1999). Under a congeneric measurement model, coefficient alpha underestimates the true reliability. However, the difference between coefficient alpha and coefficient omega is generally in the third decimal, except in the rare cases in which one of the factor loadings is very large (e.g., .9) and all of the other factor loadings are very small (e.g., .2) (Raykov, 1997).

The large-sample distribution of sample alpha.

Equation 2 shows that sample coefficient alpha is a function of the sample variances and covariances. These variances and covariances are normally distributed in large samples, not only when the item responses are normally distributed, but also under the ADF assumptions set forth by Browne (1982, 1984). As a result, and without any model assumptions, in large samples, $\hat{\alpha}$ is normally distributed with mean α and variance φ^2 . The standard error of sample alpha, $\hat{\varphi}$, can be estimated using the delta method (e.g., Agresti, 2002) from the large-sample covariance matrix of the sample variances and covariances. This matrix is different under NT and ADF assumptions. As a result, when the normality assumption for the items is replaced by the milder ADF assumption, the standard error for sample alpha will differ, and we will use $\hat{\varphi}_{NT}$ and $\hat{\varphi}_{ADF}$ to distinguish them. However, the point estimate of sample coefficient alpha, $\hat{\alpha}$, remains unchanged when NT or ADF assumptions are invoked. Note that ADF assumptions replace the normality assumption by the milder assumption that eighth-order moments of the distribution of the data are finite. This assumption is satisfied in the case of Likert-type items, in which the distribution of each item is multinomial. The assumption ensures that the fourth-order sample moments are consistent estimators of their population counterparts (Browne, 1984).

The accuracy of statistical inferences for coefficient alpha rests on the accuracy of the standard errors for sample efficient alpha. Both the NT and ADF model-free standard errors for sample coefficient alpha, proposed by van Zyl et al. (2000) and Yuan et al. (2003), respectively, are based on large-sample theory. Fortunately, Duhachek and Iacobucci (2004) showed that the NT standard errors might be well estimated with sample sizes as small as 30, provided that the item responses are approximately normally distributed. Also, sample sizes as small as 100 observations (and in some cases 50 observations) might suffice to adequately estimate ADF standard errors (Maydeu-Olivares et al., 2007).

Hypothesis Testing for Coefficient Alpha

In this section, we describe the statistical theory underlying hypothesis testing for coefficient alpha based on the results of van Zyl et al. (2000) and Yuan et al. (2003).

Case 1: Hypothesis testing involving a single-sample coefficient alpha. Consider testing whether coefficient alpha equals some a priori value, α_0 (e.g., .8 or .7). The null and alternative hypotheses are $H_0: \alpha_{dif} = 0$ and $H_1: \alpha_{dif} > 0$, where $\alpha_{dif} = \alpha - \alpha_0$. Since in large samples, $\hat{\alpha}$ is normally distributed, a suitable test statistic is

$$z = \frac{\hat{\alpha}_{dif}}{\hat{\varphi}_{dif}} = \frac{\hat{\alpha} - \alpha_0}{\hat{\varphi}}, \quad (3)$$

where $\hat{\varphi}_{dif}$ is the standard error for $\hat{\alpha}_{dif}$, and $\hat{\varphi}$ is either the NT standard error, $\hat{\varphi}_{NT}$, or the ADF standard error, $\hat{\varphi}_{ADF}$, depending on the distributional assumptions made. Then, the observed significance level (p value) for the test is the area under the standard normal curve to the left of the observed z value.

Case 2: Hypothesis testing involving two statistically independent sample coefficient alphas. This case arises when a researcher is interested in comparing coefficient alpha in two populations (e.g., male vs. female subjects) or in two disjoint samples from the same population. For testing the equality of alpha across two populations, the null and alternative hypotheses are $H_0: \alpha_{dif} = 0$ and $H_1: \alpha_{dif} \neq 0$, where $\alpha_{dif} = \alpha_1 - \alpha_2$, and α_1 and α_2 are the alpha coefficients for a test score in Populations 1 and 2, respectively. An appropriate test statistic is

$$z = \frac{\hat{\alpha}_{dif}}{\hat{\varphi}_{dif}} = \frac{\hat{\alpha}_1 - \hat{\alpha}_2}{\sqrt{\hat{\varphi}_1^2 + \hat{\varphi}_2^2}}, \quad (4)$$

where $\hat{\varphi}_1$ and $\hat{\varphi}_2$ are the (NT or ADF) standard errors for the estimates $\hat{\alpha}_1$ and $\hat{\alpha}_2$. For this two-tailed alternative, the p value of the test is obtained as twice the area under the standard normal curve to the left of $|z|$.

Case 3: Hypothesis testing involving two statistically dependent sample coefficient alphas. Armed with the code for NT standard errors provided by Duhachek and Iacobucci (2004), and with the code for ADF standard errors provided by Maydeu-Olivares et al. (2007), Case 1 and 2 hypothesis testing can readily be performed. In particular, for testing the equality of alpha across two populations, we used the fact that the variance of the difference between $\hat{\alpha}_1$ and $\hat{\alpha}_2$ equals the sum of the variances of each sample alpha. However, a researcher may be interested in testing whether the alpha coefficients for two test scores obtained from the same sample are equal. In this case, the null and alternative hypotheses are, as in Case 2, $H_0: \alpha_{dif} = 0$ and $H_1: \alpha_{dif} \neq 0$, where $\alpha_{dif} = \alpha_1 - \alpha_2$, and α_1 and α_2 are the alpha coefficients for two different test scores in the same population. An appropriate test statistic is

$$z = \frac{\hat{\alpha}_{dif}}{\hat{\varphi}_{dif}} = \frac{\hat{\alpha}_1 - \hat{\alpha}_2}{\sqrt{\hat{\varphi}_1^2 + \hat{\varphi}_2^2 + 2\text{cov}(\hat{\alpha}_1, \hat{\alpha}_2)}}, \quad (5)$$

where $\hat{\varphi}_1$ and $\hat{\varphi}_2$ are the (NT or ADF) standard errors of $\hat{\alpha}_1$ and $\hat{\alpha}_2$. The p values are obtained as in Case 2. Notice, however, that in this case, the variance for the difference between $\hat{\alpha}_1$ and $\hat{\alpha}_2$ also depends on the covariance between the two sample alphas, because they are obtained from the same sample. There are a variety of situations

in which Case 3 hypothesis testing can arise. These situations are easily handled by adopting an SEM framework. The simpler Cases 1 and 2 can also be tested using an SEM framework that directly yields the z test statistics and associated p values.

Hypothesis Testing for Coefficient Alpha Using an SEM Framework

In this section, we describe how to test hypotheses concerning coefficient alpha using SEM and the model-free approach of van Zyl et al. (2000) and Yuan et al. (2003). All that is needed is an SEM package that has the capabilities for defining additional parameters that are functions of the parameters of the model. In this article, we used Mplus Version 5 (Muthén & Muthén, 2008). We provide the annotated Mplus input files as supplementary material that can be downloaded from <http://brm.psychonomic-journals.org>. We also provide the data used in the examples as supplementary materials so the examples below can be reproduced.

Case 1: Hypothesis testing involving a single-sample coefficient alpha. Within an SEM framework, a model-free standard error for coefficient alpha can be obtained as follows:

1. Specify the model to be a $p \times p$ symmetric matrix.
2. Following Equation 1, define three additional parameters: $\gamma_1 = \sum_i \sigma_{ii}$, $\gamma_2 = \sum_{i < j} \sigma_{ij}$, and $\alpha = [p/(p-1)]\{1 - [\gamma_1/(\gamma_1 + 2\gamma_2)]\}$.
3. Define $\alpha_{\text{dif}} = \alpha - \alpha_0$.

The z statistic given by Equation 3 appears in the computer output as the ratio of the estimated α_{dif} divided by its standard error, along with its associated two-tailed p value. Because in this case the alternative is one-tailed, the two-tailed p value shown in the computer output must be divided by 2 to obtain the desired one-tailed p value.

Note that since a fully saturated model is used in Step 1, there are zero degrees of freedom, and the model fits perfectly. Also, the additional parameters in Step 2 do not introduce additional constraints on the model.

Different estimation methods can be used to estimate the parameters. Some popular choices are generalized least squares (GLS) estimation, maximum likelihood (ML) estimation, and weighted least squares (WLS) estimation. GLS and ML estimation can be performed under normality assumptions or with standard errors that are robust to normality (e.g., ADF) assumptions. WLS estimation assumes ADF assumptions. Because a saturated model is being fitted, all estimators (GLS, ML, or WLS) lead to the same point estimate for coefficient alpha, as was given in Equation 2. Also, when estimating the model under normality assumptions, GLS and ML lead to the same standard error for sample coefficient alpha, as was given by van Zyl et al. (2000). Similarly, when estimating the model without normality assumptions, robust GLS, robust ML, and WLS lead to the same standard error for sample alpha, as was given by Yuan et al. (2003). Mplus 5 implements NT GLS estimation, NT ML estimation, robust ML estimation, and WLS estimation. Also, Mplus yields as optional output confidence intervals for any parameter in

the model (including additional parameters, such as α or α_{dif}). In Mplus 5, GLS and ML denote GLS and ML estimation, respectively, under normality assumptions. ML estimation with robust standard errors is performed by using MLM or MLMV. MLM and MLMV yield the same parameter estimates and standard errors and differ only in the goodness-of-fit statistics provided. For the models considered here, MLM and MLMV yield the same fit—a perfect fit—because the models are saturated.

Case 2: Hypothesis testing involving two statistically independent sample coefficient alphas. For two populations, we need to extend the previous SEM setup to two populations as follows:

- (1) For each population, specify the model to be a $p \times p$ symmetric matrix.
- (2) Define three additional parameters as above for each population. Thus, for the first population, define γ_{11} , γ_{21} , and α_1 . For the second population, define γ_{12} , γ_{22} , and α_2 .
- (3) Define $\alpha_{\text{dif}} = \alpha_1 - \alpha_2$.

Again, the model fits perfectly, and the z statistic given by Equation 4 appears in the Mplus output as the ratio of the estimated α_{dif} divided by its standard error, along with the desired two-tailed p value. Also, a confidence interval for α_{dif} may be requested.

Case 3: Hypothesis testing involving two statistically dependent sample coefficient alphas. Consider two test scores computed on the same sample of respondents. This may occur when the two test scores being compared are alternate forms of the same test (possibly with no items in common) or when the two test scores correspond to pretest and posttest administrations of the same test. The first test score is based on p_1 items, and the second is based on p_2 items. Some items may appear on both test scores, so that the overall number of items is $p \leq p_1 + p_2$. When no item appears on both test scores, the overall number of items is $p = p_1 + p_2$. On the other hand, $p < p_1 + p_2$ when one or more items appear on both test scores. This may occur when one test score corresponds to the full form of a test and the other test score corresponds to a reduced form of the test.

The procedure involved for testing a hypothesis involving the difference between the alphas is very similar to the previous ones:

- (1) Specify the model to be a $p \times p$ symmetric matrix.
- (2) Define three additional parameters for each test score: γ_{11} , γ_{21} , and α_1 for the first test score and γ_{12} , γ_{22} , and α_2 for the second test score.
- (3) Define $\alpha_{\text{dif}} = \alpha_1 - \alpha_2$.

Again, we do not impose any constraints among the p items. The model fits perfectly. The z statistic given by Equation 5 appears in the Mplus output as the ratio of the estimated α_{dif} divided by its standard error, along with the desired two-tailed p value.

In the next section, we provide numerical examples to illustrate hypothesis testing for coefficient alpha under both

normality assumptions and the less stringent ADF assumptions. As an example of Case 1, we test whether the population coefficient alpha of a scale score equals .9. As an example of Case 2, we test whether the population coefficient alphas across genders are equal. We provide two examples of Case 3. In the first example, we test whether the population alpha of the scale score equals the scale score when only half of the items are used. These scale scores correspond to the full and short forms of a questionnaire. In the second example, we test whether the population alpha of a scale score changes when the questionnaire is administered to the same respondents at two time points.

A Numerical Example: Testing Reliability Hypotheses Based on Coefficient Alpha for the Negative Problem Orientation Scale Scores

The negative problem orientation (NPO) scale is one of the five scales of the Social Problem Solving Inventory (SPSI-R; D’Zurilla, Nezu, & Maydeu-Olivares, 2002; see also Maydeu-Olivares, Rodríguez-Fornells, Gómez-Benito, & D’Zurilla, 2000). Two forms of this inventory are available: the full form and the short form. In its full form, the NPO scale consists of 10 items. Each item is to be answered using a 5-point response scale. The short scale consists of a subset of 5 items. Here, we shall use two random samples, 100 male and 100 female respondents, from the normative U.S. sample. The correlations and standard deviations among the 10 NPO items in these samples are provided in Table 1. The first 5 items shown in Table 1 correspond to the items composing the short form. Note that the correlations and standard deviations provided suffice for NT hypothesis testing involving coefficient alpha. For ADF hypothesis testing, the raw data are needed. The raw data are provided as supplementary materials.

Example 1: Testing the hypothesis that $\alpha = .9$ for the full NPO scale in the female population (Case 1). Using ML and assuming that the items are approximately normally distributed, $\hat{\alpha} = .88$ and $\hat{\phi}_{NT} = .02$. The z statistic of Equation 3 is $z = -1.04$, yielding a two-tailed p value of .30. Hence, the one-tailed p value is .15, and we cannot reject the hypothesis that $\alpha = .9$ in the female population. Because in this sample the NPO items do not markedly depart from a normal distribution, we obtain

almost identical results when using the milder ADF assumptions. In that case, the z statistic is -1.08 and $p = .14$ (one-tailed). Mplus also yields (upon request) confidence intervals for α and for α_{dif} . A 95% confidence interval for α under normality assumptions is (.85; .92) and the ADF interval is the same (to two significant digits).

Example 2: Testing the hypothesis that the population alpha for the full NPO scale is equal across genders (Case 2). For the male sample, under normality assumptions, $\hat{\alpha} = .84$ and $\hat{\phi}_{NT} = .02$. The Mplus output yields an estimated alpha difference (males – females) of $-.05$ with a standard error of .03 under normality assumptions. The z statistic from Equation 4 is -1.52 , yielding a p value of .13. We cannot reject the hypothesis that population alpha is equal across genders. In the male sample, the NPO items do not markedly depart from a normal distribution either. As a result, when we replace the NT assumptions by ADF assumptions, a similar result is obtained: $z = -1.48$ and $p = .14$.

Example 3: Testing the equality of alpha between the full and short forms of the NPO scale in the male population (Case 3). For the short form in the male sample, $\hat{\alpha} = .72$ and $\hat{\phi}_{NT} = .04$. Also, the estimated alpha difference (full – short) is .12 with a standard error of .03. The z statistic from Equation 5 is 4.21 and $p < .01$. We reject the hypothesis of equality of alpha for the full and short NPO scale scores in the male population. Again, similar results are obtained under ADF assumptions: $z = 3.60$ and $p < .01$.

As a final example, we provide another example of Case 3. This example involves testing the hypothesis of equality of alpha for two repeated administrations of the short forms of the NPO scale. The two administrations are 3 weeks apart. The sample includes both male and female respondents ($N = 138$). Table 2 provides the correlations and standard deviations among the five items at each administration. The first five items shown in Table 2 correspond to the first administration, and the last five items correspond to the second administration.

Example 4: Testing the equality of alpha in two repeated administrations of the short form of the NPO scale (Case 3). Under ADF assumptions, a 95% confidence interval for alpha for the first administration is (.69; .82), whereas a 95% confidence interval for the sec-

Table 1
Correlations and Standard Deviations Among the Items of the Negative Problem Orientation Scale

	Male Respondents ($n = 100$)										Female Respondents ($n = 100$)									
	i1	i2	i3	i4	i5	i6	i7	i8	i9	i10	i1	i2	i3	i4	i5	i6	i7	i8	i9	i10
i1	1.00										1.00									
i2	.32	1.00									.39	1.00								
i3	.23	.24	1.00								.39	.42	1.00							
i4	.28	.29	.31	1.00							.37	.39	.46	1.00						
i5	.41	.44	.22	.28	1.00						.37	.46	.18	.43	1.00					
i6	.44	.35	.47	.38	.39	1.00					.26	.41	.37	.51	.45	1.00				
i7	.21	.48	.37	.25	.33	.47	1.00				.43	.54	.47	.62	.50	.50	1.00			
i8	.29	.32	.50	.26	.33	.32	.28	1.00				.20	.24	.25	.41	.33	.29	.50	1.00	
i9	.24	.35	.29	.29	.35	.37	.46	.33	1.00											
i10	.20	.23	.46	.26	.25	.47	.50	.28	.55	1.00	.46	.50	.60	.55	.38	.39	.64	.44	.63	1.00
SD	1.09	1.15	1.19	1.02	.97	1.19	1.02	1.36	1.19	1.08	1.13	1.14	1.19	1.30	1.17	1.21	1.23	1.20	1.34	1.31

Table 2
Correlations and Standard Deviations Among the Items of the Short Scale of the Negative Problem Orientation Scale Measured at Two Time Points (N = 138)

	i1	i2	i3	i4	i5	r1	r2	r3	r4	r5
i1	1.00									
i2	.53	1.00								
i3	.41	.33	1.00							
i4	.37	.37	.35	1.00						
i5	.47	.31	.42	.24	1.00					
r1	.52	.39	.35	.34	.45	1.00				
r2	.45	.58	.37	.36	.45	.57	1.00			
r3	.45	.41	.69	.35	.59	.51	.51	1.00		
r4	.34	.35	.28	.49	.34	.45	.36	.50	1.00	
r5	.48	.48	.37	.36	.64	.61	.59	.56	.44	1.00
<i>SD</i>	1.16	1.16	1.22	1.10	1.25	1.12	1.19	1.20	1.02	1.29

Note—Items measured at Time 1 are denoted as i1–i5, and items measured at Time 2 are denoted as r1–r5.

ond administration is (.79; .89). Given these intervals, it is difficult to determine whether coefficient alpha is equal across administrations. In contrast, the *z* statistic from Equation 5 is $-3.06, p < .01$. We clearly reject the hypothesis of invariance of alpha across administrations. A higher alpha was obtained for the second administration. For these data, a similar result is obtained under normality assumptions: $z = -2.98, p < .01$.

Accuracy of the *p* values. As we have pointed out, the accuracy of the test statistics rely on the accuracy of the standard errors. Duhachek and Iacobucci (2004) and Maydeu-Olivares et al. (2007) investigated the accuracy of the NT and ADF standard errors in a variety of situations and reported that they are accurate with sample sizes of 100 (and, in some cases, even fewer) observations. However, when applying these test statistics, the applied researcher may be in doubt as to whether the conditions confronted in his or her study match those investigated in previous studies. In other words, the *p* values obtained may be in doubt.

To verify the accuracy of the *p* values for a particular study, a simulation study can be performed using the estimated parameters of the model as the true parameter values. Using the capabilities of Mplus for Monte Carlo simulation, we investigated the accuracy of the *p* values in each of our four examples, using the actual sample size from each of the studies as the sample size in our simulations. We provide the annotated Mplus files for the simulation as supplementary materials that can be downloaded from <http://brm.psychonomic-journals.org>. Because in our examples the items were approximately normally distributed, multivariate normal data were generated in each case using the estimated mean and covariance matrix from

each example as the true mean and covariance matrix. In each case, 1,000 random samples were drawn. Table 3 provides the empirical rejection rates for each of the examples. As can be seen in this table, given the small sample sizes considered, the *p* values obtained are reliable. Also notice that the *p* values for Example 2 are somewhat more accurate than those for the remaining three examples. This may be related to the sample sizes involved. In Example 2, the sample size is larger (100 male and 100 female respondents) than in the other examples (100 female respondents in Example 1, 138 individuals in Example 3, and 100 male respondents in Example 4).

Readers familiar with SEM methods may wonder why ADF methods work so well in this situation when extant research reveals that much larger sample sizes are needed for ADF methods to work well. We believe that the reason is that ADF methods are used here only in the estimation of standard errors and not in the point estimation of coefficient alpha. Indeed, in this setup, coefficient alpha is estimated by sample coefficient alpha as given in Equation 2 (see Maydeu-Olivares, Coffman, & Hartmann, 2007, for further details).

Discussion

The two most widely used strategies for drawing statistical inferences about the reliability of a test score are the use of coefficient alpha and the use of a model-based reliability coefficient.

The model-based approach begins by fitting a measurement model to the items composing the test. When a model cannot be rejected at the usual significance level, a reliability coefficient based on the fitted model can be employed. For instance, suppose that interest lies in per-

Table 3
Empirical Rejection Rates of the Test Statistic at the Exact Settings for Each of Our Examples

Example	Sample Size	Empirical Rejection Rates (%)										
		1	5	10	20	30	40	50	60	70	80	90
1	100	2.3	7.3	11.6	19.9	29.4	39.0	48.1	60.2	71.1	81.3	90.7
2	200	1.5	5.4	10.5	20.8	30.4	39.3	50.2	60.2	69.6	79.9	90.5
3	100	0.1	3.1	7.7	19.8	30.3	40.0	48.8	57.8	67.0	76.0	86.7
4	138	2.2	6.8	12.1	21.5	30.0	39.1	49.8	61.2	72.1	83.1	92.5

Note—In each case, 1,000 replications were generated using the actual example’s sample size.

forming hypothesis testing on a single reliability coefficient (as in Case 1). If a one-factor model fits the items, statistical inferences about the reliability of the test score can be performed using coefficient omega, because this coefficient equals the reliability of the test. There is an extensive literature on using structural equation modeling to perform statistical inferences using model-based reliability estimates for a variety of measurement models (see, e.g., Kano & Azuma, 2003; Raykov, 2004; and Raykov & ShROUT, 2002). However, implementing the model-based approach may prove difficult in applications. Often, all measurement models under consideration will be rejected. In this case, model-based reliability inferences can be based on the best-fitting model found, which will fit the data only approximately. However, because the model does not fit the data exactly, the model-based reliability estimate will be biased, and the direction and magnitude of the bias will be unknown. Also, in some cases, in which a model can be found that is not rejected by the exact goodness of fit test, the measurement model may be too complex for applied researchers to easily compute the appropriate reliability estimate. Consider for instance the random intercept factor model of Maydeu-Olivares and Coffman (2006). Computing a model-based reliability coefficient for this model need not be immediately obvious to an applied researcher. Finally, implementing a model-based approach becomes more difficult when interest lies in drawing reliability inferences across different populations or for test-retest situations (Cases 2 and 3, discussed above), because different measurement models may be needed across populations or time points.

The alternative strategy—drawing inferences using population coefficient alpha—is easy to implement, because it is model-free in the sense that only a positive definite covariance matrix is assumed. However, researchers using this strategy should bear in mind that they are performing statistical inferences on population alpha and not on population reliability. In general, these are two different population parameters that are only equal when a tau-equivalent model fits the items. To claim inferences about population reliability using coefficient alpha, researchers need to test the adequacy of the tau-equivalent model. Furthermore, coefficient alpha need not always be a lower bound to population reliability. Coefficient alpha is a lower bound to the reliability of a test score whenever the item scores have the decomposition

$$Y_i = T_i + E_i, \quad (6)$$

where $i = 1, \dots, p$; where the true scores, T_i , and errors, E_i , are uncorrelated; and where the E_i s are uncorrelated with each other (e.g., Bentler, 2007; Novick & Lewis, 1967). This condition is quite general. For instance, if a k -factor model fits the data, Equation 6 is satisfied, and coefficient alpha is a lower bound to the reliability of the test score. To see this, consider the k -factor model $Y_i = \mu_i + \lambda'_i \eta + \varepsilon_i$, where μ_i , λ'_i , and ε_i denote the mean, the $1 \times k$ vector of factor loadings, and the unique factor for the i th item, respectively, and η denotes the $k \times 1$ vector of factors. Letting $\mu_i + \lambda'_i \eta = T_i$ and $\varepsilon_i = E_i$, the k -factor model is a

special case of Equation 6. Thus, in many instances, such as when a k -factor model fits the data, coefficient alpha is a lower bound to population reliability. Nevertheless, it may be best to simply claim inferences about population alpha rather than population reliability. Claiming lower bound properties for coefficient alpha without fitting a measurement model should be avoided, because when Equation 6 is not satisfied, such as when some of the errors, E_i , are correlated, population alpha may be larger than population reliability (Green & Hershberger, 2000; Komaroff, 1997; Raykov, 2001). In our experience, however, the situations in which alpha is larger than reliability are rather rare.

Concluding Remarks

In this article, we have shown that drawing statistical inferences for population alpha is quite straightforward. Statistical inferences for coefficient alpha are model free and do not require assuming that the items composing the test score are normally distributed. Because of this computational ease, researchers interested in drawing statistical inferences for population reliability may want to consider drawing inferences for population alpha instead. We do believe that researchers should attempt to draw inferences for population reliability whenever possible. However, this requires that a well-fitting measurement model can be found and the model-based reliability estimate is easy to compute. If a well-fitting model can be found but the model-based reliability estimate is cumbersome to compute, researchers may consider drawing inferences for coefficient alpha instead. If the fitted model satisfies the conditions for coefficient alpha to be a lower bound to reliability, drawing inferences for coefficient alpha becomes an attractive option to drawing inferences for population reliability from a computational viewpoint. When no well-fitting measurement model can be found, researchers may still draw inferences for coefficient alpha, because this is a meaningful parameter per se. However, in this case, researchers drawing inferences about coefficient alpha should carefully avoid extrapolating their conclusions to population reliability or claiming that coefficient alpha is a lower bound of population reliability. These claims need to be supported by model fitting.

AUTHOR NOTE

A.M.-O., C.G.-F., and D.G.-P. were supported by Grant SEJ2006-08204 from the Spanish Ministry of Education awarded to the first author. D.L.C. was supported by NIDA Center Grant P50 DA100075 and NIDA Training Grant T32 DA017629-01A1. Correspondence concerning this article should be addressed to A. Maydeu-Olivares, Faculty of Psychology, University of Barcelona, P. Valle de Hebrón, 171, 08035 Barcelona, Spain (e-mail: amaydeu@ub.edu).

REFERENCES

- AGRESTI, A. (2002). *Categorical data analysis* (2nd ed.). Hoboken, NJ: Wiley.
- BARCHARD, K. A., & HAKSTIAN, A. R. (1997). The robustness of confidence intervals for coefficient alpha under violation of the assumption of essential parallelism. *Multivariate Behavioral Research*, *32*, 169-191.
- BENTLER, P. M. (2007). Covariance structure models for maximal reli-

- ability of unit-weighted composites. In S.-Y. Lee (Ed.), *Handbook of latent variable and related models* (pp. 1-17). Amsterdam: Elsevier.
- BROWNE, M. W. (1982). Covariance structures. In D. M. Hawkins (Ed.), *Topics in applied multivariate analysis* (pp. 72-141). Cambridge: Cambridge University Press.
- BROWNE, M. W. (1984). Asymptotically distribution-free methods for the analysis of covariance structures. *British Journal of Mathematical & Statistical Psychology*, **37**, 62-83.
- CORTINA, J. M. (1993). What is coefficient alpha? An examination of theory and applications. *Journal of Applied Psychology*, **78**, 98-104.
- CRONBACH, L. J. (1951). Coefficient alpha and the internal structure of tests. *Psychometrika*, **16**, 297-334.
- DUHACHEK, A., & IACOBUCCI, D. (2004). Alpha's standard error (ASE): An accurate and precise confidence interval estimate. *Journal of Applied Psychology*, **89**, 792-808.
- D'ZURILLA, T. J., NEZU, A. M., & MAYDEU-OLIVARES, A. (2002). *Social problem-solving inventory—revised (SPSI-R)*. North Tonawanda, NY: MultiHealth Systems.
- FELDT, L. S. (1965). The approximate sampling distribution of Kuder–Richardson reliability coefficient twenty. *Psychometrika*, **30**, 357-370.
- GREEN, S. B., & HERSHBERGER, S. L. (2000). Correlated errors in true score models and their effect on coefficient alpha. *Structural Equation Modeling*, **7**, 251-270.
- GUTTMAN, L. (1945). A basis for analyzing test–retest reliability. *Psychometrika*, **10**, 255-282.
- HAKSTIAN, A. R., & WHALEN, T. E. (1976). A *k*-sample significance test for independent alpha coefficients. *Psychometrika*, **41**, 219-231.
- HOGAN, T. P., BENJAMIN, A., & BREZINSKI, K. L. (2000). Reliability methods: A note on the frequency of use of various types. *Educational & Psychological Measurement*, **60**, 523-531.
- KANO, Y., & AZUMA, Y. (2003). Use of SEM programs to precisely measure scale reliability. In H. Yanai, A. Okada, K. Shigemasa, Y. Kano, & J. J. Meulman (Eds.), *New developments in psychometrics* (pp. 141-148). Tokyo: Springer.
- KOMAROFF, E. (1997). Effect of simultaneous violations of essential tau-equivalence and uncorrelated error on coefficient alpha. *Applied Psychological Measurement*, **21**, 337-348.
- KRISTOF, W. (1963). The statistical theory of stepped-up reliability when a test has been divided into several equivalent parts. *Psychometrika*, **28**, 221-238.
- LORD, F. M., & NOVICK, M. R. (1968). *Statistical theories of mental test scores*. Reading, MA: Addison-Wesley.
- MAYDEU-OLIVARES, A., & COFFMAN, D. L. (2006). Random intercept item factor analysis. *Psychological Methods*, **11**, 344-362.
- MAYDEU-OLIVARES, A., COFFMAN, D. L., & HARTMANN, W. M. (2007). Asymptotically distribution free (ADF) interval estimation of coefficient alpha. *Psychological Methods*, **12**, 157-176.
- MAYDEU-OLIVARES, A., RODRÍGUEZ-FORNELLS, A., GÓMEZ-BENITO, J., & D'ZURILLA, T. J. (2000). Psychometric properties of the Spanish adaptation of the Social Problem-Solving Inventory—Revised (SPSI-R). *Personality & Individual Differences*, **29**, 699-708.
- MCDONALD, R. P. (1999). *Test theory: A unified treatment*. Mahwah, NJ: Erlbaum.
- MILLER, M. B. (1995). Coefficient alpha: A basic introduction from the perspectives of classical test theory and structural equation modeling. *Structural Equation Modeling*, **2**, 255-273.
- MUTHÉN, L., & MUTHÉN, B. (2008). *Mplus 5*. Los Angeles, CA: Muthén & Muthén.
- NOVICK, M. R., & LEWIS, C. (1967). Coefficient alpha and the reliability of composite measurements. *Psychometrika*, **32**, 1-13.
- RAYKOV, T. (1997). Scale reliability, Cronbach's coefficient alpha, and violations of essential tau-equivalence with fixed congeneric components. *Multivariate Behavioral Research*, **32**, 329-353.
- RAYKOV, T. (2001). Bias of coefficient α for fixed congeneric measures with correlated errors. *Applied Psychological Measurement*, **25**, 69-76.
- RAYKOV, T. (2004). Point and interval estimation of reliability for multiple-component measuring instruments via linear constraint covariance structure modeling. *Structural Equation Modeling*, **11**, 342-356.
- RAYKOV, T., & SHROUT, P. E. (2002). Reliability of scales with general structure: Point and interval estimation using a structural equation modeling approach. *Structural Equation Modeling*, **9**, 195-212.
- SCHMITT, N. (1996). Uses and abuses of coefficient alpha. *Psychological Assessment*, **8**, 350-353.
- SHEVLIN, M., MILES, J. N. V., DAVIES, M. N. O., & WALKER, S. (2000). Coefficient alpha: A useful indicator of reliability? *Personality & Individual Differences*, **28**, 229-237.
- TEN BERGE, J. M. F. (2000). Linking reliability and factor analysis: Recent developments in some classical psychometric problems. In S. E. Hampson (Ed.), *Advances in personality psychology* (Vol. 1, pp. 138-156). Philadelphia: Psychology Press.
- VAN ZYL, J. M., NEUDECKER, H., & NEL, D. G. (2000). On the distribution of the maximum likelihood estimator of Cronbach's alpha. *Psychometrika*, **65**, 271-280.
- YUAN, K.-H., GUARNACCIA, C. A., & HAYSLIP, B., JR. (2003). A study of the distribution of sample coefficient alpha with the Hopkins symptom checklist: Bootstrap versus asymptotics. *Educational & Psychological Measurement*, **63**, 5-23.

SUPPLEMENTAL MATERIALS

Mplus and data files for the hypothesis testing procedure discussed in this article may be downloaded from <http://brm.psychonomic-journals.org/content/supplemental>.

(Manuscript received September 2, 2009;
revision accepted for publication October 24, 2009.)