

Sex Differences in Sum Scores May Be Hard to Interpret

The Importance of Measurement Invariance

M. C. T. Slof-Op 't Landt

Center for Eating Disorders Ursula, Leidschendam, VU University, Amsterdam, Leiden University Medical Centre, Leiden, The Netherlands

E. F. van Furth

Center for Eating Disorders Ursula, Leidschendam, The Netherlands

I. Rebollo-Mesa

M. Bartels

C. E. M. van Beijsterveldt

P. E. Slagboom

D. I. Boomsma

VU University, Amsterdam, The Netherlands

I. Meulenbelt

Leiden University Medical Centre, Leiden, The Netherlands

C. V. Dolan

University of Amsterdam, Amsterdam, The Netherlands

In most assessment instruments, distinct items are designed to measure a trait, and the sum score of these items serves as an approximation of an individual's trait score. In interpreting group differences with respect to sum scores, the instrument should measure the same underlying trait across groups (e.g., male/female, young/old). Differences with respect to the sum score should accurately reflect differences in the latent trait of interest. A necessary condition for this is that the instrument is measurement invariant. In the current study, the authors illustrate a stepwise approach for testing measurement invariance with respect to sex in a four-item instrument designed to assess disordered eating behavior in a large epidemiological sample (1,195 men and 1,507 women). This approach can be applied to other phenotypes for which group differences are expected. Any analysis of such variables may be subject to measurement bias if a lack of measurement invariance between grouping variables goes undetected.

Keywords: *measurement invariance, confirmatory factor analysis, sex differences, eating disorders, sex*

Questionnaires are often used to assess psychological and behavioral traits on a quantitative scale. Well-known examples are the Beck Depression Inventory (Beck, Ward, Mendelson, Mock, & Erbaugh, 1961), Eysenck Personality Questionnaire (Eysenck & Eysenck, 1975), and the Temperament and Character

Inventory (Cloninger, Svrakic, & Przybeck, 1993). In these assessment instruments, items are designed to measure an underlying trait or latent (i.e., unobserved) variable and scores on the items are summed to derive a total score on the trait of interest. The *Diagnostic and Statistical Manual of Mental Disorders*, fourth edition (*DSM-IV*; American Psychiatric Association, 1994) also employs a weighted sum score in diagnosing psychiatric disorders.

When comparing groups, it is vital that an instrument measures the same underlying trait across groups (e.g., male/female, young/old). Observed group differences in the sum scores should accurately reflect group

Authors' Note: The Netherlands Organization of Scientific Research (NWO; grant numbers 575-25-012 and NWO/SPI 56-464-14192) NWO: VENI:451-04-034 (M. Bartels). Please address correspondence to Margarita C.T. Slof-Op 't Landt, Center for Eating Disorders Ursula, PO Box 422, 2260 AK Leidschendam, The Netherlands; email: r.optlandt@centrummeetstoornissen.nl

differences with respect to the latent variable. A necessary condition for this is that the instrument displays measurement invariance with respect to the groups under consideration (Mellenbergh, 1989; Meredith, 1993). If there is a sex difference with respect to the latent trait, men should, for example, score lower on all the items of the instrument measuring this trait. If however, men score lower on all the items but one, this one item displays differential item functioning, and the scale is not measurement invariant with respect to sex (Dolan, 2000; Mellenbergh, 1989; Meredith, 1993; Millsap & Yun-Tein, 2004). In that case, group differences in sum scores reflect, at least in part, measurement bias. The interpretation of differences between groups with respect to the sum scores thus hinges on the establishment of measurement invariance or at least on the understanding of the violations, if any, of measurement invariance. Ideally, differences in sum scores should reflect true differences in the latent variable that the psychometric instrument purports to measure.

Measurement invariance can be investigated by fitting a measurement model that relates item scores to the underlying trait(s) across groups. Several methods have been suggested for both continuous and categorical variables (Dolan, 2000; Mellenbergh, 1989; Meredith, 1993; Millsap & Yun-Tein, 2004; Muthén & Asparouhov, 2002; Muthén, & Muthén, 2005). In the current study, we described a stepwise approach that was derived from previous studies to investigate measurement invariance for ordered categorical items. Our goal was to provide a comprehensive overview of the different steps accumulating into a model of complete measurement invariance. To illustrate this approach, we investigated whether a four-item instrument, designed to measure disordered eating behavior is measurement invariant with respect to sex. As eating disorders mainly affect young women (90% to 95% of cases; Fairburn, & Harrison, 2003; Hoek, 1993; Van Hoeken, Lucas, & Hoek, 1998), one might expect sex differences in the endorsement of the four eating disorder items. Multigroup discrete factor analyses were applied to test whether the disordered eating behavior instrument is measurement invariant with respect to sex.

Method

Participants

All participants were registered with the Netherlands Twin Registry, which is maintained at the Department of Biological Psychology at the VU

University in Amsterdam (Bartels, Van Beijsterveldt, Stroet, Hudziak, & Boomsma, 2007; Boomsma et al., 2006). In this study, we used data from the 1986 to 1992 birth cohorts. In January 2005, questionnaires were sent to adolescent twins (mean age = 15.2 years, $SD = 1.3$) and their nontwin siblings (mean age = 16.7 years, $SD = 2.8$). The twins and siblings were asked to complete a survey containing items relevant for eating disorders. Questionnaires were sent to 2,000 families. A total of 2,175 twins (twin response rate 54.4%) and 527 siblings from 1,144 families returned the questionnaire (family response rate 57.2%). The total sample consisted of 1,195 men and 1,507 women (956 male twins, 1,219 female twins, 239 brothers and 288 sisters, respectively), mean age was 15.5 years ($SD = 1.8$).

Measures

Participants filled out a self-report questionnaire containing measures of health and behavior (Bartels et al., 2007; Boomsma et al., 2006). The eating disorder section included four items: (a) dieting (Question: Have you ever gone on a diet to lose weight or to stop gaining weight?), (b) fear of weight gain (Question: How afraid are you to gain weight or become fat?), (c) importance of body weight or shape on self-evaluation (Question: How important are body weight and/or shape in how you feel about yourself?), and (d) binge eating (Question: Have you ever had episodes of binge eating?). Responses were given on 5-point Likert-type scales, ranging from *never* to *always* for dieting (DIET), from *not afraid* to *extremely afraid* for fear of weight gain (FEAR), from *not important* to *most important* for importance of body weight and shape on self-evaluation (ISE), and from *never* to *more than once a week* for binge eating (BE). For the multigroup confirmatory factor analyses it was essential that, for every item, each category was endorsed by both groups. Because none of the men reported that they were always on a diet, the fourth and fifth categories of the dieting item were merged. As a consequence, three items with five categories and one item with four categories were used in the analyses.

Data Analysis

We performed multigroup confirmatory factor analyses to establish whether the four eating disorder items formed a unidimensional scale and whether the scale was measurement invariant with respect to sex. To conduct a confirmatory factor analysis, a minimum of three items is required. Measurement invariance with respect to sex held if the probability of a

certain response on a given item was the same for all participants with the same value on the underlying trait (disordered eating behavior [DEB]) regardless of the sex of the participant. This definition gave rise to a highly constrained multigroup factor model (Chen, Sousa, & West, 2005; Meredith, 1993; Millsap & Yun-Tein, 2004). To establish measurement invariance, we fitted several increasingly restrictive models derived from approaches described in previous studies (Dolan, 2000; Mellenbergh, 1989; Meredith, 1993; Millsap & Yun-Tein, 2004; Muthén & Asparouhov, 2002; Muthén & Muthén, 2005), cumulating in this highly constrained model.

In the first step, a saturated model was fitted to the data simply to obtain estimates of the item thresholds and the polychoric correlation among items. To this end, we assumed that a latent continuous variable, called the liability, was underlying the responses to each discrete item. Assuming the liability underlying each item was standard normally distributed, the discrete responses were modeled to items by estimating thresholds on the standard normal distributions of the liability (3 thresholds for the DIET item and 4 thresholds for the other three items). The positions of these thresholds determined the marginal response probabilities of each item. In addition, the (polychoric) correlations among the liability underlying the four items were estimated. Thresholds and correlations were estimated separately in men and women.

In the second model it was tested whether the four items were unidimensional in men and women. The four continuous latent liabilities were regressed on a single common factor, without imposing any equality constraints over sex. Thresholds in men were constrained to equal those in the women. By imposing this constraint, the thresholds were estimated on a common metric. The distribution of the liability for each item was standard normal in the women as in Model 1. In the men the means and variances of the liability underlying the four items were estimated freely. Thus, in this step, we fitted a single factor model to the correlation matrix of the liabilities in the women and a single factor model to the covariance matrix of the liabilities in the men. In both sexes, the common factor was scaled to have a mean of zero and a variance of one (i.e., standard scaling constraints in the common factor model). By estimating all the factor loadings freely, the item reliability in the women and the men were obtained separately. Note that these reliability estimates need not be equal over sex.

In Model 3, the factor loadings were constrained to be equal over sex. This constraint allowed estimation of the variance of the common factor in one group (men), while retaining the scaling constraint (variance of factor equal to one) in the other group (women). We thus allowed for a difference in common factor variance between men and women. This model included sex differences in the residual variances of the items, in the liability means, and in the common factor variance.

In Model 4, mean liabilities (intercepts) in the male sample were constrained at zero and the common factor mean was estimated. As before, the mean liabilities and common factor mean were fixed to zero in women. In the preceding model, the estimated mean in liabilities in men gave an indication of the sex differences per item. By fixing these intercepts at zero in men, while freely estimating the mean of the common factor, any sex difference in means of the liabilities was explained by a difference in the mean of the common factor, that is, a difference with respect to the latent variable of interest.

In Model 5, we added the final constraint of "invariance of residual variances over sex." As a consequence, the amount of the variance in the separate items that was not explained by the common factor was constrained to be equal in the women and men. This model represented full measurement invariance. Note that in this model any observed sex difference in the observed test scores was attributable to a difference with respect to the latent variable that we purported to measure. With respect to the interpretation of sex differences in test scores, Model 5 represented the ideal. Model 4 represented a weaker form of invariance in which sex differences in the residuals were permitted. Model 4 was still useful as it allowed us to interpret sex differences in the mean scale score as a manifestation of a mean difference with respect to the latent variable. Weaker forms of measurement invariance are entertained in the literature (e.g., Model 3: equality of factor loadings), but we did not consider these to be sufficient for the interpretation of sex differences with respect to the test scores (Meredith, 1993).

All analyses were performed in Mplus 4.0 (Muthén & Asparouhov, 2002; Muthén & Muthén, 2005). Because our sample consisted of families, the individual cases were not independent. To correct for the effect of this dependence on the standard errors and overall goodness-of-fit indices, we used the weighted least squares with mean adjusted chi-square test statistics in combination with the

“Complex” option in Mplus. The latter corrects the statistical effect of clustering on the results. Rebollo, de Moor, Dolan, and Boomsma (2006) found this method to be satisfactory to correct for dependency due to family grouping.

As suggested by Schermelleh-Engel, Moosbrugger, and Müller (2003), several fit statistics were used to evaluate the fit of the models; hierarchical chi-square tests, the comparative fit index (CFI), and the root mean square error of approximation (RMSEA). For the hierarchical chi-square test, the difference between the chi-square test statistics obtained for each model yielded a new chi-square value with degrees of freedom equal to the difference in the number of parameters in the two models. In the weighted least square with mean adjusted chi-square test statistics approach in Mplus, the reported chi-squares were mean adjusted and a scaling correction factor was applied for each model. As a consequence, in calculating the chi-square difference test, scaling correction factors had to be entered into the equation (Asparouhov & Muthen, 2006). According to the principle of parsimony, models with fewer parameters are preferred, if they do not give a significant deterioration of the fit. Significance can be determined on statistical grounds, but in structural equation modeling, rules of thumb are usually used (Schermelleh-Engel et al., 2003). The CFI ranges from zero to one with higher values indicating better fit; for a good model fit the CFI should be >0.97 , and values >0.95 indicate an acceptable fit (Schermelleh-Engel et al., 2003). The RMSEA is a measure of closeness of fit, and provides a measure of discrepancy per degree of freedom. A value of 0.05 or smaller indicates a close fit and values between 0.05 and 0.08 indicate an acceptable fit (Jöreskog, 1993; Schermelleh-Engel et al., 2003).

There were 257 persons ($n = 127$ men and $n = 130$ women) who completed the survey twice with an interval of 6 months. Retest data obtained in this group will serve to estimate stability of the test scores. The reliability of the eating disorder items was estimated separately in men and women. Polychoric correlations between the two occasions of measurement were calculated for each item using Mplus.

Results

To evaluate how often the different eating disorder attitudes and behaviors were endorsed, we calculated

the frequencies of the item scores greater than three in the adolescent twins and their nontwin siblings for the four items. These frequencies showed significant sex differences for three features ($p < .001$). For the DIET item, 0.4% of the men compared with 3.4% of the women had been on a diet often or always. Few men (1.3%) reported being very or extremely afraid to gain weight or become fat (FEAR). In women this item was endorsed more often with 8.7%. A large proportion of both men and women reported that “their body weight and or shape played an important role in how they felt about themselves” (ISE). The frequency of this feature was 40.9% in the women compared with 26.8% in the men. No sex differences were found for the BE item, 5.1% of the women and 5.5% of the men reported having binge eating episodes at least once a week.

In Model 1 polychoric correlations among items and the thresholds for each item were estimated per sex. These are reported in Table 1. Small to moderate correlations between the items were found in both sexes. Although the magnitude of the correlations differed between groups, similar patterns were observed with the highest correlation between DIET and FEAR and the lowest between ISE and BE. The thresholds of the liabilities represent the cut-points of the response categories in the corresponding ordinal items on a sex-specific z scale. The mainly positive thresholds indicate that the majority of women and men did not engage in eating disordered behaviors and/or attitudes.

In Table 2, fit statistics of the nested models are given. Model 2, which tested whether one factor could account for the correlations among the four eating disorder variables, fitted significantly worse compared with Model 1 according to the chi-square. However, both the RMSEA and the CFI indicated a good fit of this model. The parameter estimates of Model 2 are presented in Table 3. The factor loadings of DIET and BE were comparable between men and women. On the other hand, the factor loading in the men for FEAR was higher and for ISE was lower compared with the women. The least reliable item was BE, whereas the FEAR item had the highest reliability.

The estimates of the mean liability in men were all significantly lower than zero. As these means were fixed to zero in the women, we established, as expected, that the men scored lower than the women on all eating disorder items. The estimated variances of the liability of FEAR, ISE, and BE were significantly smaller than one in the men. The variances were fixed at one in the women.

Table 1
Correlations and Thresholds for Women and Men (Saturated Model)

	DIET	FEAR	ISE	BE
Correlations ^a				
DIET	1.00	0.59 (0.52, 0.66)	0.39 (0.30, 0.48)	0.41 (0.31, 0.51)
FEAR	0.53 (0.38, 0.67)	1.00	0.59 (0.54, 0.64)	0.33 (0.24, 0.41)
ISE	0.27 (0.13, 0.40)	0.39 (0.29, 0.48)	1.00	0.27 (0.19, 0.36)
BE	0.22 (0.03, 0.41)	0.20 (0.06, 0.34)	0.16 (0.06, 0.27)	1.00
Women				
Threshold 1	0.68 (0.57, 0.78)	-0.43 (-0.52, -0.33)	-1.54 (-1.68, -1.40)	0.64 (0.54, 0.74)
Threshold 2	1.36 (1.24, 1.49)	0.67 (0.57, 0.77)	-0.56 (-0.66, -0.46)	1.12 (1.00, 1.23)
Threshold 3	1.83 (1.66, 1.99)	1.36 (1.24, 1.49)	0.23 (0.14, 0.32)	1.63 (1.48, 1.78)
Threshold 4	—	2.13 (1.93, 2.33)	1.88 (1.70, 2.06)	2.08 (1.87, 2.29)
Men				
Threshold 1	1.62 (1.44, 1.80)	0.72 (0.60, 0.83)	-0.94 (-1.06, -0.83)	0.95 (0.83, 1.07)
Threshold 2	2.29 (2.02, 2.56)	1.71 (1.52, 1.89)	-0.10 (-0.20, 0.001)	1.27 (1.13, 1.40)
Threshold 3	2.64 (2.25, 3.02)	2.24 (1.97, 2.51)	0.58 (0.48, 0.69)	1.60 (1.44, 1.76)
Threshold 4	—	2.71 (2.28, 3.14)	1.94 (1.73, 2.15)	1.87 (1.68, 2.06)

Note: DIET = dieting; FEAR = fear of weight gain; ISE = importance of body weight or shape in self-evaluation; BE = binge eating.
 a. The correlations in the women are listed above the diagonal, the correlations in the men are listed below the diagonal. Numbers in parentheses represent 95% confidence intervals. The thresholds are estimated on a sex-specific z scale.

Table 2
Model Fit Statistics

Model	χ^2	df	CFI	RMSEA	CM	$\Delta\chi^2$	Δdf	p
Model 1 (saturated)	0.00	0	1.00	0.00	—	—	—	—
Model 2 (one-factor model)	37.98	11	0.99	0.04	1	37.98	11	.0001
Model 3	35.62	14	0.99	0.03	2	2.32	3	.5100
Model 4	101.07	17	0.96	0.06	3	50.99	3	.0001
Model 5 (full measurement invariance)	246.53	21	0.90	0.09	4	99.57	4	.0001

Note: CFI = comparative fit index; RMSEA = root mean square error of approximation; CM = compared to model; $\Delta\chi^2$ = chi-square test statistic between two models adjusted for scaling correction factor; Δdf = degrees of freedom for the chi-square difference test.

Table 3
Parameter Estimates for Model 2 in the Female Reference Group and the Male Group

	DIET	FEAR	ISE	BE
Women				
Factor loading	0.68 (0.60, 0.75)	0.88 (0.81, 0.94)	0.66 (0.59, 0.72)	0.44 (0.35, 0.52)
Mean	0	0	0	0
Variance	1	1	1	1
Reliability	.46	.77	.43	.19
Men				
Factor loading	0.69 (0.35, 1.03)	0.97 (0.71, 1.24)	0.55 (0.41, 0.70)	0.45 (0.19, 0.71)
Mean	-1.11 (-1.83, -0.39)	-1.30 (-1.56, -1.05)	-0.44 (-0.57, -0.31)	-0.84 (-1.28, -0.40)
Variance	0.91 (0.59, 1.24)	0.84 (0.70, 0.98)	0.85 (0.77, 0.93)	0.65 (0.50, 0.79)
Reliability	.48	.94	0.30	0.20

Note: Numbers in parentheses represent 95% confidence intervals for the factor loadings and residual variances. DIET = dieting; FEAR = fear of weight gain; ISE = importance of body weight or shape in self-evaluation; BE = binge eating.

The chi-square test statistic suggested some violation of unidimensionality (Model 2). But because both the RMSEA and the CFI indicated a good fit, the invariance of factor loadings across sexes was tested next. For this model, all three fit statistics indicated a good fit. The estimate of variance of the common factor (DEB) in the male group was 0.96. Given the 95% confidence interval (CI) of 0.62 and 1.30, we concluded that the variance was not significantly different between the men and women in Model 3.

In Model 4, the mean of the liabilities were constrained to be zero in men (as they were in women). The mean of the common factor was fixed to zero in the women, as before, and estimated freely in the men. This model did not fit very well in comparison with Model 3. The chi-square test statistic indicated a significantly worse fit for this model. However, the fit was acceptable according to the RMSEA and the CFI. The estimated common factor mean in the men was -0.99 , which differed significantly from zero (95% CI = -1.18 to -0.80). In other words, the mean of DEB was lower in men than in women (factor mean fixed at zero).

Because the fit of Model 4 was acceptable based on the RMSEA and the CFI, the final model of complete measurement invariance was tested. In this fifth model, the residual variances were also constrained to be equal across the groups. The chi-square statistic indicated deterioration in fit compared with Model 4. In addition, the CFI and the RMSEA indicated a bad fit. This implied that the eating disorder items were not fully measurement invariant with respect to sex. The variances presented in Table 3, give an indication of which item might be underlying this bad fit. The variance of BE showed the largest deviation from 1, suggesting that the greatest difference between both groups in residual variance was observed for this item.

Finally the stability of the item responses and the DEB total score were considered. The four eating disorder items were moderately to highly correlated over a period of 6 months. The polychoric correlation was 0.59 (95% CI = 0.28-0.89) for DIET, 0.75 (95% CI = 0.59-0.90) for FEAR, 0.56 (95% CI = 0.41-0.71) for ISE, and 0.74 (95% CI = 0.55-0.93) for BE in men. In women, the polychoric correlation was 0.75 (95% CI = 0.60-0.89) for DIET, 0.67 (95% CI = 0.55-0.79) for FEAR, 0.43 (95% CI = 0.27-0.59) for ISE, and 0.58 (95% CI = 0.42-0.74) for BE.

Discussion

In most assessment instruments, distinct items are designed to measure a trait, and the sum score of these items serves as an approximation of an individual's trait score. The interpretation of differences between groups with respect to these sum scores hinges on the establishment of measurement invariance. Ideally, differences in sum scores should reflect true differences in the latent variable that the psychometric instrument purports to measure. If there is a lack of measurement invariance, group differences in sum scores reflect, at least in part, measurement bias.

We described a stepwise multigroup confirmatory factor analysis to investigate measurement invariance for categorical items with respect to a grouping variable. Previously, several methods have been reported to test for measurement invariance both for continuous and categorical items (Dolan, 2000; Mellenbergh, 1989; Meredith, 1993; Millsap & Yun-Tein, 2004; Muthen & Asparouhov, 2002; Muthén & Muthén, 2005). All these methods cumulated in an identical highly constrained model in which strict factorial invariance, or complete measurement invariance, was tested. However, the number and order of the constraints in the intermediate models differed between the reported methods. In contrast to previous studies, our analysis began by fitting a saturated model to the data, to obtain estimates of the polychoric correlation among items and the thresholds for each item. The second model, which tested for unidimensionality of the items, was more comparable with the baseline models described by other groups (Millsap & Yun-Tein, 2004; Muthen & Asparouhov, 2002; Muthén & Muthén, 2005), although there was a difference in the constraints. In our model, thresholds were constrained across groups, whereas factor loadings were estimated freely. This enabled us to calculate the reliability of the separate item scores. Means and variances of the liabilities provided insight in the between-group differences. In the third model, both item thresholds and factor loadings were constrained to be equal across groups. The between-group differences in this model were represented by the residual variances of the items, the liability means, and by the common factor variance. In addition to the previous constraints, the liability means were constrained at zero in all groups in Model 4. Within this model, any

group difference in the means of the latent indicators would be explained by a difference in the mean of the common factor. This model represented a weaker form of invariance in which group differences were permitted in the residuals, and was similar to the third model described by Millsap and Yun-Tein (2004). The final model of strict factorial invariance, added the constraint of invariance of residual variances over groups; that is, the amount of the variance in each item that was not explained by the common factor was constrained to be equal in the groups.

The method was illustrated by investigating whether a scale comprised of four eating disorder items was measurement invariant with respect to sex. The model of full measurement invariance with respect to sex (Model 5) did not fit the data well. If this model had fitted, the probability of a certain response on a given item would have been the same for all participants with the same value on the underlying trait (DEB) regardless of the sex of the participant. However, this was not the case. The underlying common factor might not be the only source of difference between the sexes with respect to the four items. The sum score based on the four eating disorder items therefore cannot be taken to represent exactly the same underlying trait in men and women. This means that sex differences in this sum score might be due to measurement bias instead of a true difference in the underlying trait.

What implication does this finding have for existing eating disorder measurement instruments? We acknowledge that a scale consisting of four items might not be ideal to measure the underlying latent trait in eating disorders. However, in large epidemiological studies such as becoming common for gene finding, short scales might be a requirement to obtain phenotyping in sufficiently large samples. With the selection of the items we have tried to capture a variety of eating disorder symptoms. Three of the items (FEAR, ISE, and BE) used in this study are based on eating disorder criteria from the *DSM-IV* (American Psychiatric Association, 1994). The fourth item (DIET) has been identified as a potent risk factor (Jacobi, Hayward, de Zwaan, Kraemer, & Agras, 2004). However, one eating disorder symptom, compensatory behavior, is missing in our assessment instrument.

There has been a lot of debate about whether eating disorders are dimensional like proposed in the "continuum of eating disorders" (Fairburn & Harrison, 2003; Hay, & Fairburn, 1998) or whether they are

discrete syndromes (Williamson, Gleaves, & Stewart, 2005). Some studies suggest that eating disorders can be conceptualized as having at least two latent features (Williamson et al., 2002; Williamson et al., 2005); binge eating and general psychopathology. Accordingly, the FEAR, DIET, and ISE items would load on one factor and the BE item would load on a second factor. The correlations presented in Table 1, however, show substantial correlations between DIET and BE, especially in women (0.41). Bulimic behavior has been correlated with dieting and body concerns in several other studies (Williamson et al., 2005), although this correlation appears to exist exclusively in nonclinical samples. Because our sample is also nonclinical, this may be the cause of the high correlation between DIET and BE. Hence, the factor structure discussed above might not be suitable in nonclinical groups. On the other hand, the low reliability of the BE item and the fact that the variance of this item showed the largest deviation from one in Model 2, might be supportive of the two factor structure underlying eating disorders. However, investigating partial measurement invariance by omitting the final constraints on the BE item did not lead to a model of strict factorial invariance for the remaining three items.

The finding of a lack of strict factorial invariance in the 4-item DEB scale might not generalize to existing eating disorder scales. However, this form of measurement invariance has never been tested in the eating disorder field. Many studies have used both exploratory and confirmatory factor analysis to test whether existing measurement instruments have the same factor structure across, for example, different types of patients and different ethnic groups and to establish different factors within eating disorders (Calugi, Grave, Ghisi, & Sanavio, 2006; Fernandez, Malacarne, Wifley, & McQuaid, 2006; Hrabosky et al., 2008; Lee et al., 2007; Peterson et al., 2007; Varnado, Williamson, & Netemeyer, 1995; Wade, Byrne, & Bryant-Waugh, 2008; Williamson et al., 2002; Williamson et al., 2005). Until now, only one study has investigated measurement equivalence (Warren et al., 2008). Warren et al. tested for the equivalence of factor loadings for the Body Shape Questionnaire in American and Spanish women with and without an eating disorder diagnosis. For a subscale of 10 items, the constraint of invariant factor loadings fitted the data well. However, because the intercepts were not constrained to equivalence in this study, the scores in the different groups may not have the same origin

(Chen et al., 2005). Thus, differences on factor means between groups could still be caused by measurement bias.

The responses to the four eating disorder items were fairly stable over a 6 month period, with correlations ranging from 0.43 for ISE in the women to 0.75 for FEAR in the men and DIET in the women. The prevalence for the DIET, FEAR, and BE item were low to moderate. The prevalence of ISE was substantially higher. Comparable rates were found in other population-based studies in adolescents with the exception of the DIET item, which had a lower prevalence (Kjelsas, Bjornstrom, & Gotestam, 2004; Neumark-Sztainer, Wall, Haines, Story, & Eisenberg, 2007; Rowe, Pickles, Simonoff, Bulik, & Silberg, 2002; Silberg, & Bulik, 2005). Because of the low endorsement rates of dieting in the men, we had to merge the fourth and fifth category for the DIET item. As a consequence, the number of response frequencies differed between the four items. This difference in response categories does not appear to affect the results. When all items are merged into four or even three categories, the same results were found throughout the different steps of the confirmatory factor analyses. Comparable correlations, thresholds, and factor loadings for the four items were found. In addition, the model of weak measurement invariance (Model 4) remained the best-fitting model.

The framework we presented in this article can serve as a valuable tool for examining the psychometric qualities of other interviews and questionnaires with respect to sex. In addition, other kinds of grouping variables (e.g., age, level of education) can also be studied using this method.¹ An advantage of our approach is that it provides a better understanding of the consequences of the different constraints per model. As a consequence, it gives a better insight into the violations of measurement invariance and the underlying causes of this measurement bias. It is essential to test for measurement invariance before sum scores or scale scores are used to compare groups. This is not only the case in the eating disorder field, but applies to other fields of research as well.

Acknowledgments

Financial support by The Netherlands Organization of Scientific Research is gratefully acknowledged.

Note

1. The Mplus scripts can be obtained from the first author on request.

References

- American Psychiatric Association. (1994). *Diagnostic and statistical manual of mental disorders* (4th ed.). Washington, DC: Author.
- Asparouhov, T., & Muthen, B. (2006). *Robust chi square difference testing with mean and variance adjusted test statistics*. Mplus Web Notes, No. 10.
- Bartels, M., Van Beijsterveldt, C. E. M., Stroet, T. M., Hudziak, J. J., & Boomsma, D. I. (2007). Young-Netherlands Twin Register (Y-NTR): A longitudinal multiple informant study of problem behavior. *Twin Research and Human Genetics, 10*, 3-12.
- Beck, A. T., Ward, C. H., Mendelson, M., Mock, J., & Erbaugh, J. (1961). An inventory for measuring depression. *Archives of General Psychiatry, 4*, 561-571.
- Boomsma, D. I., de Geus, E. J., Vink, J. M., Stubbe, J. H., Distel, M. A., Hottenga, J. J., et al. (2006). Netherlands twin register: From twins to twin families. *Twin Research and Human Genetics, 9*, 849-857.
- Calugi, S., Grave, R. D., Ghisi, M., & Sanavio, E. (2006). Validation of the body checking questionnaire (BCQ) in an eating disorders population. *Behavioural and Cognitive Psychotherapy, 34*, 233-242.
- Chen, F. F., Sousa, K. H., & West, S. G. (2005). Testing measurement invariance of second-order factor models. *Structural Equation Modeling, 12*, 471-492.
- Cloninger, C. R., Svrakic, D. M., & Przybeck, T. R. (1993). A psychobiological model of temperament and character. *Archives of General Psychiatry, 50*, 975-990.
- Dolan, C. V. (2000). Investigating Spearman's hypothesis by means of multi-group confirmatory factor analysis. *Multivariate Behavioral Research, 35*, 21-50.
- Eysenck, H. J., & Eysenck, S. B. G. (1975). *Manual of the Eysenck Personality Questionnaire*. San Diego, CA: Digits.
- Fairburn, C. G., & Harrison, P. J. (2003). Eating disorders. *Lancet, 361*, 407-416.
- Fernandez, S., Malacrne, V. L., Wifley, D. E., & McQuaid, J. (2006). Factor structure of the Bulimia Test-Revised in college women from four ethnic groups. *Cultural Diversity and Ethnic Minority Psychology, 12*, 403-419.
- Hay, P., & Fairburn, C. (1998). The validity of the DSM-IV scheme for classifying bulimic eating disorders. *International Journal of Eating Disorders, 23*, 7-15.
- Hoek, H. W. (1993). Review of the epidemiological studies of eating disorders. *International Review of Psychiatry, 5*, 61-74.
- Hrabosky, J. I., White, M. A., Masheb, R. M., Rothschild, B. S., Burke-Martindale, C. H., & Grilo, C. M. (2008). Psychometric evaluation of the eating disorder examination-questionnaire for bariatric surgery candidates. *Obesity, 16*, 763-769.
- Jacobi, C., Hayward, C., de Zwaan, M., Kraemer, H. C., & Agras, W. S. (2004). Coming to terms with risk factors for eating disorders: Application of risk terminology and suggestions for a general taxonomy. *Psychological Bulletin, 130*, 19-65.

- Jöreskog, K. G. (1993). Testing structural equation models. In K. A. Bollen & J. Scott Long (Eds.), *Testing structural equation models* (pp. 294-317). Newbury Park, CA: Sage.
- Kjelsas, E., Bjornstrom, C., & Gotestam, K. G. (2004). Prevalence of eating disorders in female and male adolescents (14-15 years). *Eating Behaviors*, *5*, 13-25.
- Lee, S. W., Stewart, S. M., Striegel-Moore, R. H., Lee, S., Ho, S. Y., Lee, P. W. H., et al. (2007). Validation of the eating disorder diagnostic scale for use with Hong Kong adolescents. *International Journal of Eating Disorders*, *40*, 569-574.
- Mellenbergh, G. J. (1989). Item bias and item response theory. *International Journal of Educational Research*, *13*, 127-143.
- Meredith, W. (1993). Measurement invariance, factor analysis and factorial invariance. *Psychometrika*, *58*, 525-543.
- Millsap, R. E., & Yun-Tein, J. (2004). Assessing factorial invariance in ordered-categorical measures. *Multivariate Behavioral Research*, *39*, 479-515.
- Muthén, B., & Asparouhov, T. (2002). *Latent variable analysis with categorical outcomes: Multiple-group and growth modeling in Mplus*. Mplus Web Notes, No. 4.
- Muthén, L. K., & Muthén, B. O. (2005). *Mplus user's guide* (3rd ed.) Los Angeles: Muthén & Muthén.
- Neumark-Sztainer, D., Wall, M., Haines, J., Story, M., & Eisenberg, M. E. (2007). Why does dieting predict weight gain in adolescents? Findings from project EAT-II: A 5-year longitudinal study. *Journal of the American Dietetic Association*, *107*, 448-455.
- Peterson, C. B., Crosby, R. D., Wonderlich, S. A., Joiner, T., Crow, S. J., Mitchell, J. E., et al. (2007). Psychometric properties of the Eating Disorder Examination-Questionnaire: Factor structure and internal consistency. *International Journal of Eating Disorders*, *40*, 386-389.
- Rebollo, I., de Moor, M. H., Dolan, C. V., & Boomsma, D. I. (2006). Phenotypic factor analysis of family data: Correction of the bias due to dependency. *Twin Research and Human Genetics*, *9*, 367-376.
- Rowe, R., Pickles, A., Simonoff, E., Bulik, C. M., & Silberg, J. L. (2002). Bulimic symptoms in the Virginia Twin Study of Adolescent Behavioral Development: Correlates, comorbidity, and genetics. *Biologic Psychiatry*, *51*, 172-182.
- Schermelleh-Engel, K., Moosbrugger, H., & Müller, H. (2003). Evaluating the fit of structural equation models: Tests of significance and descriptive goodness-of-fit measures. *Methods of Psychological Research Online*, *8*, 23-74.
- Silberg, J. L., & Bulik, C. M. (2005). The developmental association between eating disorders symptoms and symptoms of depression and anxiety in juvenile twin girls. *Journal of Child Psychology and Psychiatry*, *46*, 1317-1326.
- Van Hoeken, D., Lucas, A. R., & Hoek, H. W. (1998). Epidemiology. In H. W. Hoek, J. Treasure, & M. Katzman (Eds.), *The integration of neurobiology in the treatment of eating disorders* (pp. 97-126). London: Wiley.
- Varnado, P. J., Williamson, D. A., & Netemeyer, R. (1995). Confirmatory factor analysis of eating disorder symptoms in college women. *Journal of Psychopathology and Behavioral Assessment*, *17*, 69-79.
- Wade, T. D., Byrne, S., & Bryant-Waugh, R. (2008). The eating disorder examination: Norms and construct validity with young and middle adolescent girls. *International Journal of Eating Disorders*, *41*, 551-558.
- Warren, C. S., Cepeda-Benito, A., Gleaves, D. H., Moreno, S., Rodriguez, S., Fernandez, M. C., et al. (2008). English and Spanish versions of the Body Shape Questionnaire: Measurement equivalence across ethnicity and clinical status. *International Journal of Eating Disorders*, *41*, 265-272.
- Williamson, D. A., Gleaves, D. H., & Stewart, T. M. (2005). Categorical versus dimensional models of eating disorders: An examination of the evidence. *International Journal of Eating Disorders*, *37*, 1-10.
- Williamson, D. A., Womble, L. G., Smeets, M. A. M., Netemeyer, R. G., Thaw, J. M., Kutlesic, V., et al. (2002). Latent structure of eating disorder symptoms: A factor analytic and taxometric investigation. *American Journal of Psychiatry*, *159*, 412-418.