

A common measurement scale for self-report instruments in mental health care: T scores with a normal distribution

## Introduction

The importance of measurement for clinical management and quality improvement of Mental Health Care (MHC) is widely acknowledged (Kilbourne et al., 2018; Lambert & Harmon, 2018). Inspired by the recovery movement (Slade et al., 2008) and developments as shared decision making (Patel et al., 2008), feedback informed treatment (Miller et al., 2015), and measurement based care (Harding et al., 2011; Kilbourne et al., 2018), patients' needs and preferences are granted a more prominent role in MHC. Increased patient involvement requires being well informed about the severity of one's condition at the onset of treatment and about progress made in the journey towards recovery. Therefore, a good understanding of measurement results is needed, when findings of Routine Outcome Monitoring (ROM; de Beurs et al., 2011) are shared. However, the huge diversity of measurement instruments in use in clinical practice, each with its own measurement scale, complicates a straightforward communication among professionals and between professionals and patients about their measurement results. This might hamper further implementation of measurement-based care. Use of common, measure-independent metrics for results of clinical tests may be a solution to this problem (Authors, 2021).

Two metrics have been proposed and researched in recent years: T scores and percentile rank scores. T scores, first proposed by McCain (1922), are standardized scores (Z-scores) multiplied by 10 and 50 added, resulting in a metric with  $M=50$  ( $SD=10$ ). The T score denotes the commonness of a test result by its

distance from the mean of a reference group in standard units. T scores require the assumption that test results are measured on an interval measurement scale. Z scores (and T scores) based on data from a reference population may have a highly skewed frequency distribution. It is advisable to first transform the raw scores to have a normal frequency distribution before converting them to T scores, which results in *normalized* T scores. When the normalized T score metric is calibrated on the general population, most psychiatric patients will score before treatment around 65-70 on measures of psychopathology and, when treated successfully, their score will decrease to 55-60 over time. Percentile rank scores have been proposed as an alternative way to express how common or how exceptional a test result is (Crawford & Garthwaite, 2009; Ley, 1972). They quantify the rarity of a tested persons score in a percentage. However, percentile scores are not on an interval scale, but ordinal and are, especially at the extremes, easily misunderstood (Bowman, 2002).

There is quite some literature on converting measurement instruments to a common metric (Dorans et al., 2007; Holland et al., 2006) and the subject is closely linked to norming instruments (Mellenbergh, 2011). Most frequently used are distribution based approaches, such as percentile conversion (Kolen & Brennan, 2014), and methods based on Item Response Theory (IRT; Embretson & Reise, 2013). Following the latter approach, several researchers have published reports on linking measures for the same constructs and have proposed population based T scores as common metric for the severity of depression (Choi et al., 2014; Fischer et al., 2011; Schalet et al., 2015; Wahl et al., 2014), anxiety (Schalet et al., 2014), pain (Cook et al., 2015), physical functioning (Schalet et al., 2015), fatigue (Friedrich et al., 2019), psychological distress (Batterham et al., 2018) and quality of life measures (ten Klooster et al., 2013). Usually, general population samples are used for easy

interpretability of the resulting metric (Wahl et al., 2014). If appropriately sampled, such samples reflect the general population. In contrast, clinical samples vary in composition and severity or complexity of the disorder, and as such they are less useful as a reference group for general psychopathology measures. Various clinical measurement instruments are administered to the same group of respondents and an IRT estimate, usually the expected a posteriori (EAP) method (Lord & Wingersky, 1984) is used to estimate  $\theta$  in order to express scores in the common metric. This endeavor is also known as the PROsetta Stone project (Schalet et al., 2015).

A potential drawback of the IRT approach is that it requires a comprehensive dataset and the use of dedicated software to calculate  $\theta$ -scores. In clinical practice this is not always feasible: a clinician may only have a raw scale score from an individual patient and he/she does not have access to the algorithm to obtain an IRT based  $\theta$ -score. To help out, several authors provide crosswalk tables to translate the raw test score into a T score (Batterham et al., 2018; Choi et al., 2014; Schalet et al., 2014). However, reading from crosswalk tables is cumbersome and prone to error. The relation between raw scores and T scores can also be modelled and expressed in a function. Such a function to calculate T scores can easily be used or can be implemented in ROM-software and provides an alternative method to obtain scores on the common metric for individual patients. In line with international developments, T scores may be based on IRT models and stem from data of general population samples (Authors, 2021). However, for everyday clinical practice, we propose an approach in which T scores are calculated with a conversion function. This will be feasible, even if only a test score from a single individual is available.

Insert Figure 1 about here

Figure 1 presents an overview of various approaches to obtain T scores. The first approach (on the left) entails standardizing the sum score of items of a scale into Z scores and converting these to T scores with  $T = 10 \cdot Z + 50$  (standard T scores). The second approach is based on percentile rank scores, which are converted to normalized Z scores and subsequently to normalized T scores and takes the frequency distribution of scores into account (percentile-based T-scores). The third approach is advocated in the present paper, with crosswalk formulas derived from regression of T scores based on the factor scores resulting from an IRT model on sum scores (calculated T scores). The fourth approach is the IRT based approach itself, practiced by -amongst others- the PROMIS group ( $\theta$ -based T scores).

In this article we present cross walk formulas for three frequently used clinical measurement instruments: the Brief Symptom Inventory (BSI; Derogatis, 1975), the Four-Dimensional Symptom Questionnaire (4DSQ; Terluin et al., 2006), and the Outcome Questionnaire (OQ-45; Lambert et al., 2004). A necessary step is to investigate the validity of calculated T scores by comparing them with  $\theta$ -based T scores. If there is high agreement, transformation of raw scale scores with the conversion function is a valid approach for obtaining a proxy for  $\theta$ -based T scores. For each measure, data were used from two samples: the general population and patients. IRT based T scores were founded on both samples, with the general population as reference population. The patient samples allowed us to investigate whether the T scores were normally distributed in clinical samples. Measures for ROM, where patients are repeatedly assessed to determine change over time, have preferably an interval scale of measurement with a normal distribution. Finally, percentile ranks are provided based on the population and clinical samples.

## Methods

### Datasets

BSI and 4DSQ data from the general population were collected on separate occasions in a large sample of the Dutch population, called the LISS-panel ("Long-term Internet Studies for the Social Sciences"). The LISS-panel is maintained by CentERdata, a research institute located on the Tilburg University campus, and includes about 5 300 household respondents who were approached with the help of the municipal population register (van der Laan, 2009). The sample is representative of the Dutch population (Scherpenzeel, 2018; Scherpenzeel & Bethlehem, 2011). In 2007/2008, one third of all households was approached, of which one person each completed the BSI ( $n = 1\,662$ ). This sample was stratified by age (four strata: 18-29, 30-49, 50-64, 65+ years), gender, and ethnic origin. In 2013, normative data on the 4DSQ were collected from the entire LISS-panel ( $n = 5\,273$ ). The sample and the procedure for collecting 4DSQ data were described by Terluin et al. (2016).

For the OQ-45, data were used from  $n = 1\,810$  respondents, which have been described by Timman et al. (2017): 1000 respondents came from a panel of the TNS-NIPO research agency and were stratified by gender, age, socio-economic status, and education level; 810 came from an earlier validation study (de Jong et al., 2007), 448 came from a sample drawn from the telephone directory and 362 were invited via internal mail from 14 companies or non-commercial organizations (Timman et al., 2017).

Patient data for the BSI came from a dataset of patients referred for treatment of depression, anxiety disorders, or somatoform disorders at RijnVeste, the outpatient clinic of GGZ Rivierduinen in the city of Leiden. Data from  $n = 4\,853$  patients, collected between 2002 and 2013, were used. The procedure of data collection has been

described by de Beurs et al. (2011). For the 4DSQ, patient data were used from 199 patients from primary care physicians. All patients were on sick leave, reported elevated levels of stress, and participated in a trial evaluating an intervention for stress-related mental disorders (Bakker et al., 2007). The patient data for the OQ-45 came from 12 436 patients described by Timman et al. (2017). This comprised a mixed sample: patients in day care (n=481) or inpatient care (n=484) from various MHC institutes; patients in outpatient care (basic and specialized MHC, n=1581 and n=9433 respectively), and patients treated by private practitioners (n=457). According to Dutch law, anonymized questionnaire data collected to support treatment, may be used for scientific research and such use is exempt from an informed consent procedure.

## Measures

The Brief Symptom Inventory (BSI; Derogatis, 1975) (Dutch version: de Beurs & Zitman, 2006) consists of 53 items describing symptoms. Respondents can indicate to what extent they were bothered by each symptom on a Likert-type scale from 0 (not at all) to 4 (very much) in the past week. In this study, we limited ourselves to investigating the three most important scales of the BSI: depression (BSI-DEP), anxiety (BSI-ANX), and somatic complaints (BSI-SOM). Scale scores are calculated as the mean score of the comprising items and range from 0 - 4. In addition, the global score for the severity of psychopathology was analyzed: the mean score on all 53 items (Global Severity Index or BSI-GSI, range 0 - 4).

The Four-Dimensional Symptom Questionnaire (4DSQ; Terluin et al., 2006) consists of 50 items, each describing one symptom. Respondents can indicate how often they experienced the symptom in the past week on a scale from 0 (no) to 4

(very often or constantly). The 4DSQ comprises four scales: general distress (4DSQ-DIS, 16 items, range 0-32), depression (4DSQ-DEP, 6 items, range 0-12), anxiety (4DSQ-ANX, 12 items, range 0-24), and somatic complaints (4DSQ-SOM, 16 items, range 0-32). All scores are sum scores (after recoding the two response options for high frequency: 4=2 and 3=2), as recommended in the scoring instruction of this instrument (Terluin et al., 2006).

The Outcome Questionnaire (OQ-45; Lambert et al., 2004) (Dutch version: de Beurs et al., 2005; de Jong et al., 2007) comprises 45 items describing symptoms or problems. The respondent is asked to indicate how often these emerged during the past week on a scale from 0 (never) to 4 (almost always). A total score can be calculated (OQ-TOTAL, 45 items, range 0-180) and four subscale scores: Symptom Distress (OQ-SD, 25 items, range 0-100), Interpersonal Relations (OQ-IR, 11 items, range 0-4), Social Role, (OQ-SR, 9 items, range 0-36) and Anxiety and Social Distress (OQ-ASD, 13 items, range 0-52). All scores are sum scores.

## Statistical analysis

### *Calculations according to IRT*

The "Graded Response Model for polytomous items" and its Expected A Posteriori (EAP) score was used as the estimator for the  $\theta$ -score with the multidimensional IRT (mirt) package (Chalmers, 2012) version in R. There are other estimates available for the latent variable scores in IRT, such as Maximum Likelihood (ML) and Weighted Likelihood Estimates (WLE). We performed a sensitivity analysis, comparing EAP with these alternatives and found sufficiently similar results regarding the mean  $\theta$ 's, with their 95% CI intervals mostly overlapping and highly correlated (results are provided in Table C in the supplementary materials). Item-parameters were

established in combined general population and clinical samples. Thus,  $\theta$  scores were estimated with the multigroup mirt option (Smits, 2016). We fixed the item parameters to be equal across groups. The latent trait ( $\theta$ ) was standardized to a scale with a mean of 0 and a standard deviation of 1 for the general population. Unidimensionality of scales, a requirement for IRT, was investigated with confirmatory factor analysis using the R package lavaan (version 06.5) (Rosseel, 2012). We used the DWLS estimator based on the polychoric correlation matrix for ordinal items and inspected (scaled) fit statistics and set as requirements for unidimensionality: CFI > .95, TLI  $\geq$  0.95, RMSEA < .06 and SRMR < .08. If insufficient fit of a unidimensional model is found, IRT based scores are potentially flawed and alternatives (e.g., percentile conversion or regression-based norming) can be utilized.

We used non-linear least squares modeling of R (nls) to establish the best fitting function (which had the lowest AIC-value and/or the most parsimonious number of coefficients) for the relation between raw scores and  $\theta$ -based T-scores. Linear, polynomial, exponential, logarithmic, power, division, rational, sigmoid, and hyperbolic equations were evaluated. We cross-validated each equation by randomly splitting each dataset in two and using the first dataset to establish the best fitting function and using the second dataset as validation sample (Camstra & Boomsma, 1992). Applying the conversion formula to the raw scale score results in a calculated T score. The distributions of the resulting scores (mean, median, standard deviation, skewness, and kurtosis) were investigated for the population and the patient samples and visually inspected on normality with histograms/density plots and QQ-plots. ICC estimates for absolute agreement and consistency of  $\theta$ -based and calculated T scores and their 95% confidence intervals (CI95) were determined using R and



based on a 2-way mixed-effects model. We established “bias” (mean difference between both T scores (van Stralen et al., 2012), as well as the “percentage error” (the width of the CI95 interval proportional to the population mean (Van Hoeck et al., 2000). If the CI95 interval was within 5 T score points, we would conclude that both approaches yielded sufficiently similar results, since 5 T score points are the proposed limit for statistically reliable change in score over time (de Beurs et al., 2019). We also inspected the agreement between  $\theta$ -based T scores and calculated T scores with Bland-Altman plots for the full range of severity (Bland & Altman, 1986). We did this for the entire sample and for the population and clinical samples, separately. Finally, to establish the effect of normalization, we also compared standard T scores (based on the standard conversion formula “standard T =  $10 \cdot Z + 50$ ”) with  $\theta$ -based T scores and calculated T scores.

## Results

### *Unidimensionality of scales*

Most scales met criteria for unidimensionality, using  $CFI \geq 0.95$ ,  $TLI \geq 0.95$ ,  $RMSEA < 0.06$  to  $0.08$ ;  $SRMR \leq .08$  (Schreiber, Nora, Stage, Barlow, & King, 2006). Table A in the supplementary materials provides CFA results for all scales. The CFI, TLI, and SRMR requirements were met by all scales of the BSI and the 4DSQ, except for the BSI-GSI, and the DIS-SOM. The OQ scales did not meet CFI and TLI requirements. RMSEA was larger than 0.06 for most of the scales, but not substantially larger, except for the BSI-ANX and - again - the OQ-SR ( $RMSEA > 0.12$ ). As a consequence, precision of the  $\theta$ 's for the OQ scales in particular may be compromised, as these scales appears to lack sufficient unidimensionality (Crişan et al., 2017).

*Distribution of raw scores and T scores*

Insert Table 1 and Figure 2 about here

Table 1 presents an overview of raw scores,  $\theta$ -based T scores, and calculated T scores and characteristics of the frequency distribution for scales of the BSI, 4DSQ, and OQ-45. Figure 2 shows for the main score on each instrument, the frequency distribution of the raw score, theta-based T score and calculated T score in the general population sample (upper half) and in the clinical sample (lower half) and QQ-plots. Similarity of the frequency distribution, the density curve, and the normal curve and accordance of the dots in the QQ-plot and the straight diagonal line, indicate how close the distributions approximate normality. (Figure A1 to A3 in the supplementary materials presents plots for all scales). The average BSI-GSI score in the population sample was  $M = 0.38$  ( $sd = 0.34$ ). For the BSI-GSI score skewness was 2.00; kurtosis was 5.97, showing substantial deviation from normality of the scores with a tail to the right. Many respondents had a low score, but this is to be expected when a symptom list is completed by a population sample. The raw scores on the 4DSQ scales in the population were also skewed. For the depression and the anxiety scale, values for skewness and kurtosis were extreme, as 78.7% and 67.6% of the respondents had the lowest possible score. In contrast, raw scores of the general population on the OQ-45 had an almost normal distribution; only the total score and the OQ-SD score showed marginal kurtosis in the population sample (some surplus of scores below the average score). Most T scores based on  $\theta$ 's and calculated T scores had a normal distribution.

Conversion functions for the BSI, 4DSQ, and OQ-45 are shown in Table 2. For BSI scores rational functions provided the best fit; indices of skewness and kurtosis

decreased considerably compared to the raw scores. Raw scale scores in the patient sample had a normal distribution, and this was preserved in the calculated T scores.

For 4DSQ scales, also rational functions best fitted the relation between raw scores and  $\theta$ -based T scores. Again, raw scores were skewed and peaked, whereas  $\theta$ -based and calculated T scores approximated the normal distribution better, although the depression score and the anxiety score of the population sample were still skewed and showed kurtosis, due to a large proportion of respondents with the lowest possible score on these scales. Thus, transformation to a normal distribution was successful for only two of the four subscales of the 4DSQ. Patient scores had a normal distribution. Finally, for the OQ-45 two rational functions, a cubic, a quadratic, and a hyperbolic function provided the best fit.

In the supplementary materials Figure A1 to A3 presents histograms with density lines and QQ-plots for all scales. These graphs reveal a sufficiently normal distribution for most scales, except for the 4DSQ depression and anxiety scales, but also reveal some surplus of extreme high and extreme low scores for all scales.

### *Comparison of $\theta$ -based and calculated T scores*

We cross-validated the equations shown in Table 2 by applying a random split of each dataset into a calibration sample and a cross-validation sample. Table 2 provides the Root Mean Squared Error (RMSE), the coefficient of determination ( $R^2$ ), and the Mean Absolute Error (MAE) for the correspondence of predicted scores (calculated T scores with formulas based on the calibration sample) with observed T scores ( $\theta$ -based T scores in the cross-validation sample). Overall, correspondence was high; the lowest  $R^2 = .869$  for the OQ-SR scale.

Furthermore, we calculated ICCs for absolute agreement between  $\theta$ -based and calculated T scores and consistency (similar ranking of subjects according to both scores) (van Stralen et al., 2012). Table 3 presents the ICCs and additional information on the association between both scores. All ICCs were high, suggesting excellent agreement and consistency. Bias and percentage error were also low, with the exception of a higher percentage error for BSI-ANX, OQ-IR, and OQ-SR and for these scales the CI95 extended beyond 5 T score points. Figure 3 presents Bland-Altman plots for selected scales. Figure B1 to B3 in the supplementary materials presents plots for all scales.

Insert Table 3 and Figure 3 about here

For the BSI,  $\theta$ -based T corresponded well to calculated T scores (mean difference  $M=0.01$  to  $0.02$ , the solid gray line in the BA plot), Only for the BSI-ANX less than 95% of the cases fell within the  $-5$  to  $5$  interval for acceptable difference. For the 4DSQ  $\theta$ -based T scores and calculated T scores corresponded well (mean difference  $M=0.02$ ), 99% of the cases fell within the  $-5$  to  $5$  interval. For the OQ-SD also excellent correspondence was found (mean difference  $M=0.02$ ); however, on average 91.2% of the cases fell within the  $-5$  to  $5$  interval with less correspondence for the OQ-IR and OQ-SR scales. Generally, in the low scoring range ( $<40$ ) and in the high score range ( $>70$ ), calculated T scores were somewhat higher; in the mid-range  $\theta$ -based T scores were higher. We also established ICC's denoting correspondence between theta-based and calculated T scores for the population and the clinical sample, separately. Correspondence was somewhat lower in the clinical sample, especially for the 4DSQ-DEP, 4DSQ\_ANX, OQ-IR, and OQ-SR, but was still very high (see Table B in the supplementary materials).

We also compared the theta-based and the calculated T scores for the BSI-GSI, 4DSQ-DIS, and OQ-SD with standard T-scores (based on the simpler linear equation  $T = 10 \cdot Z + 50$ ). ICCs ranged from  $ICC = .81$  for the BSI-GSI to  $ICC = .99$  for the OQ-SD (see Table 4). Correspondence between standard T scores and calculated T scores was still substantial, especially for the OQ-SD. However, for the BSI-GSI there was a large subgroup with higher standard T scores than theta-based or calculated T scores (for both comparisons  $\Delta < -5$ ; 18.9%), which is understandable, as negative skewness (a low mean score relative to the maximum scale score) results in more respondents with extremely high standard T scores. After all, the non-linear conversion formula yielding calculated T scores corrects for precisely this undesirable effect.

#### *Crosswalk from raw scores to T scores and percentiles*

Table 5a to 5c present crosswalk tables for the conversion of a selection of raw scores to calculated T scores and percentile ranks for the general population and the clinical population.

Insert Table 5a, 5b, and 5c about here

Per measurement instrument the first column gives raw scores (RS), the second T scores calculated according to the functions in Table 2 for all respondents. The conversion functions stretch the scales at the extremes: a change in raw score of one scale point in the low or high score area is larger than a change of one scale point in the mid-range score area and the normalized T score reflects this appropriately. Figure 4 depicts this relationship between raw scores on the instruments and T scores.

Insert Figure 4 about here

## Discussion

We analyzed community data of frequently used generic outcome measures in the Netherlands, to link raw scores to two common metrics: T scores and percentile ranks. In line with previous research (Cook et al., 2015; Fischer & Rose, 2016; Friedrich et al., 2019; Schalet et al., 2015; Wahl et al., 2014), we applied methods based on IRT, which resulted in  $\theta$ -based T scores with a normal frequency distribution. We also determined functions to convert sum scores into T scores and showed that for most scales calculated T scores approximated  $\theta$ -based T scores very well, as the scores were strongly related (all  $ICC > .95$ , see Table 3). Scores were similar across the width of the entire scale, and yielded similar mean values for the groups. Correspondence between calculated and  $\theta$ -based T scores supports the validity of the approach we followed to calculate T scores with a function based on curve fitting. However, at the extreme end of the scales the two approaches diverged somewhat. Thus, caution with extreme scores is in order, especially with  $T < 40$  and  $T > 80$ . Furthermore, the findings of two scales of the 4DSQ revealed that, if raw scale scores are too skewed or too leptokurtic to begin with (due to an excess of respondents with the lowest possible score), conversion to  $\theta$ -based or calculated T scores will not yield scores with a normal frequency distribution.

In comparison to standard T scores ( $T = 10 \cdot Z + 50$ ), the correspondence of the more complex conversion formulas (correcting for a non-normal frequency distribution of raw scores) with  $\theta$ -based T scores was definitely better. For instance, ICC between standard T scores and  $\theta$ -based T scores for the BSI-GSI was  $ICC = .93$  (See Table 4), whereas calculated T scores showed almost perfect correspondence with  $\theta$ -based T scores ( $ICC = .99$ ; see Table 3). Better

correspondence was especially obtained in the higher score range. For each scale, we cross validated the formulas on subsamples (splitting the samples randomly in half and we also compared the population and clinical samples). The results revealed similarity and high correlation between predicted scores from the “learning sample” with obtained scores in the “test sample”. However, thus a substantial number of statistical tests were done per scale in each dataset and the applicability of these formulas still needs further validation with data from other samples.

The results show that some scales evoked the lowest possible score from many respondents (especially the 4DSQ; see Figure 2). This zero-inflation in the data is quite common when measures of psychopathology are administered in the general population. With zero-inflated data, the Graded Response Model may yield biased results of IRT analyses and alternative models to deal with zero-inflation have been proposed (Wall et al., 2015), such as Zero Inflated Mixture GRM (ZIM-GRM). In a simulation study with zero-inflated scores, GRM showed substantial bias (Smits et al., 2020). In future studies the effect of zero-inflation on factor score estimates should be investigated to ascertain the added value of these models.

It should be noted that there are viable alternatives to arrive at normative values and subsequently calculate T scores and percentile rank scores, instead of the IRT-based methods or the approach advocated in the present paper. When scales are not unidimensional and/or the IRT model does not fit, alternatives can be used, e.g., frequency-based percentile rank scores. Especially, regression-based norming (Mellenbergh, 2011) is an interesting option, e.g., with the GAMLSS model (Stasinopoulos et al., 2018), in which the shape on the frequency distribution of raw scores, as well as relevant “norm-predictors”, such as gender, age, and educational level can be taken into account. A useful introduction to the continuous norming

approach, which corrects for all levels on these demographic variables, is offered by Timmerman et al. (2020). However, for the present purposes, we build upon the previous work of the PROsetta stone initiative and the PROMIS group.

Strengths of the study are the use of large data sets with representative samples from the general population and from patients, warranting trust in the findings. The data from the Dutch general population were collected by research institutes with a good reputation, and the representativeness of these population samples has been documented by Scherpenzeel (2018) and Timman et al. (2017). A relatively simple and straightforward approach is outlined, which leads to calculated T scores that approximate  $\theta$ -based T scores well. A potential limitation of the proposed method is that it relies heavily on data from the general population, as this is the basis of the item parameters and  $\theta$ 's used to obtain T scores. The clinical measurement instruments were not developed for the population at large, but rather for patients suffering from psychological complaints or symptoms of psychiatric disorders. Cronbach (1984) already noted that, ideally, instruments should be validated with data from the population for which the instrument was intended. Indeed, the frequency distribution of responses of non-clinical respondents deviated much more from normality than was the case with data obtained in the clinical samples, as the present findings show.

The present results were based on normative data from the Netherlands. Application of the conversion formulas listed in Table 2 or in the cross-walk tables (5a to 5c) is limited to this context, as general population respondents from other countries may score differently on these self-report measures. A further limitation may be the composition of the clinical sample used to evaluate the frequency distribution of the conversion formula in clinical data. All these patients suffered from



mild to moderate common mental disorders and patients with severe mental illness were not included. Future research should investigate other clinical samples.

## Conclusion

The high correlations found in this study between  $\theta$ -based and calculated T scores for assessment with the 4DSQ, the BSI, and the OQ-45, suggests that the proposed approach of using conversion functions provides a good approximation towards a normalized common metric for MHC. Use of such a common metric will make the interpretation of test results easier for therapists and patients (Authors, 2021), will allow for better involvement of patient in shared decision making regarding the treatment (Patel et al., 2008), and will stimulate the future uptake of measurement based mental health care (Kilbourne et al., 2018).

## Acknowledgements

In this paper we gratefully made use of BSI- and 4DSQ data of the LISS (Longitudinal Internet Studies for the Social sciences) panel administered by CentERdata (Tilburg University, The Netherlands) and OQ-45 data from TNS-NIPO, The Netherlands. We would also like to thank MHC providers in the Netherlands for providing patient data on the BSI and the OQ-45, and the VU Medical Center for providing patient data on the 4DSQ.

## References:

- Authors. (2021). From mandating common measures to mandating common metrics: a plea to harmonize measurement results. <https://doi.org/10.31234/osf.io/m4qzb>
- Bakker, I. M., Terluin, B., Van Marwijk, H. W., van der Windt, D. A. M., Rijmen, F., van Mechelen, W., & Stalman, W. A. (2007). A cluster-randomised trial evaluating an intervention for patients with stress-related mental disorders and sick leave in primary care. *PLoS Clinical Trials*, 2(6), e26. <https://doi.org/10.1371/journal.pctr.0020026>
- Batterham, P. J., Sunderland, M., Slade, T., Cleave, A. L., & Carragher, N. (2018). Assessing distress in the community: psychometric properties and crosswalk comparison of eight measures of psychological distress. *Psychological Medicine*, 48(8), 1316-1324. <https://doi.org/10.1017/S0033291717002835>
- Bland, J. M., & Altman, D. G. (1986). Statistical methods for assessing agreement between two methods of clinical measurement. *Lancet*, 84767, 307-310. [https://doi.org/10.1016/S140-6736\(86\)90837-8](https://doi.org/10.1016/S140-6736(86)90837-8)
- Bowman, M. L. (2002). The perfidy of percentiles. *Archives of Clinical Neuropsychology*, 17(3), 295-303. [https://doi.org/10.1016/S0887-6177\(01\)00116-0](https://doi.org/10.1016/S0887-6177(01)00116-0)
- Camstra, A., & Boomsma, A. (1992). Cross-validation in regression and covariance structure analysis: An overview. *Sociological Methods & Research*, 21(1), 89-115. <https://doi.org/10.1177/0049124192021001004>
- Chalmers, R. P. (2012). mirt: A multidimensional item response theory package for the R environment. *Journal of Statistical Software*, 48(6), 1-29. <https://doi.org/10.18637/jss.v048.i06>
- Choi, S. W., Schalet, B., Cook, K. F., & Cella, D. (2014). Establishing a common metric for depressive symptoms: Linking the BDI-II, CES-D, and PHQ-9 to PROMIS Depression. *Psychological Assessment*, 26(2), 513-527. <https://doi.org/10.1037/a0035768>
- Cook, K. F., Schalet, B. D., Kallen, M. A., Rutsohn, J. P., & Cella, D. (2015). Establishing a common metric for self-reported pain: linking BPI Pain Interference and SF-36 Bodily Pain Subscale scores to the PROMIS Pain Interference metric. *Quality of Life Research*, 24(10), 2305-2318. <https://doi.org/10.1007/s11136-014-0790-9>
- Crawford, J. R., & Garthwaite, P. H. (2009). Percentiles please: The case for expressing neuropsychological test scores and accompanying confidence limits as percentile ranks. *The Clinical Neuropsychologist*, 23(2), 193-204. <https://doi.org/10.1080/13854040801968450>
- Crişan, D. R., Tendeiro, J. N., & Meijer, R. R. (2017). Investigating the Practical Consequences of Model Misfit in Unidimensional IRT Models. *Applied Psychological Measurement*, 41(6), 439-455. <https://doi.org/10.1177/0146621617695522>
- Cronbach, L. J. (1984). *Essentials of psychological testing* (4th ed.). Harper & Row.
- de Beurs, E., Carlier, I. V., & van Hemert, A. M. (2019). Approaches to denote treatment outcome: Clinical Significance and Clinical Global Impression compared. *International Journal of Methods in Psychiatric Research*, 28. <https://doi.org/10.1002/mpr.1797>
- de Beurs, E., den Hollander-Gijsman, M., Buwalda, V., Trijsburg, W., & Zitman, F. G. (2005). De Outcome Questionnaire (OQ-45): een meetinstrument voor meer dan alleen psychische klachten [The Outcome Questionnaire (OQ-45): a measure for psychiatric symptoms and more]. *De Psycholoog*, 40(1), 53-63.
- de Beurs, E., den Hollander-Gijsman, M. E., van Rood, Y. R., van der Wee, N. J., Giltay, E. J., van Noorden, M. S., van der Lem, R., E., v. F., & Zitman, F. G. (2011). Routine outcome monitoring in the Netherlands: Practical experiences with a web-based strategy for the assessment of

- treatment outcome in clinical practice. *Clinical Psychology & Psychotherapy*, 18(1), 1-12. <https://doi.org/10.1002/cpp.696>
- de Beurs, E., & Zitman, F. G. (2006). De Brief Symptom Inventory (BSI): De betrouwbaarheid en validiteit van een handzaam alternatief voor de SCL-90 [The Brief Symptom Inventory: Reliability and validity of a handy alternative for the SCL-90]. *Maandblad Geestelijke Volksgezondheid*, 61, 120-141.
- de Jong, K., Nugter, M. A., Polak, M. G., Wagenborg, J. E. A., Spinhoven, P., & Heiser, W. J. (2007). The Outcome Questionnaire (OQ-45) in a Dutch population: A cross-cultural validation. *Clinical Psychology & Psychotherapy*, 14(4), 288-301. <https://doi.org/10.1002/cpp.529>
- Derogatis, L. R. (1975). *The Brief Symptom Inventory*. Clinical Psychometric Research.
- Dorans, N. J., Pommerich, M., & Holland, P. W. (Eds.). (2007). *Linking and aligning scores and scales*. Springer.
- Embretson, S. E., & Reise, S. P. (2013). *Item response theory for psychologists*. Lawrence Erlbaum Associates.
- Fischer, H. F., & Rose, M. (2016). www.common-metrics.org: a web application to estimate scores from different patient-reported outcome measures on a common scale. *BMC Medical Research Methodology*, 16(1), 142. <https://doi.org/10.1186/s12874-016-0241-0>
- Fischer, H. F., Tritt, K., Klapp, B. F., & Fliege, H. (2011). How to compare scores from different depression scales: Equating the Patient Health Questionnaire (PHQ) and the ICD-10-Symptom Rating (ISR) using item response theory [10.1002/mpr.350]. *International Journal of Methods in Psychiatric Research*, 20(4), 203-214. <https://doi.org/10.1002/mpr.350>
- Friedrich, M., Hinz, A., Kuhnt, S., Schulte, T., Rose, M., & Fischer, F. (2019). Measuring fatigue in cancer patients: a common metric for six fatigue instruments. *Quality of Life Research*, 28(6), 1615-1626. <https://doi.org/10.1007/s11136-019-02147-3>
- Harding, K. J., Rush, A. J., Arbuckle, M., Trivedi, M. H., & Pincus, H. A. (2011). Measurement-based care in psychiatric practice: a policy framework for implementation. *Journal of Clinical Psychiatry*, 72(8), 1136-1143. <https://doi.org/10.4088/JCP.10r06282whi>
- Holland, P. W., Dorans, N. J., & Petersen, N. S. (2006). 6 Equating Test Scores. *Handbook of statistics*, 26, 169-203.
- Kilbourne, A. M., Beck, K., Spaeth-Rublee, B., Ramanuj, P., O'Brien, R. W., Tomoyasu, N., & Pincus, H. A. (2018). Measuring and improving the quality of mental health care: a global perspective. *World Psychiatry*, 17(1), 30-38. <https://doi.org/10.1002/wps.20482>
- Kolen, M. J., & Brennan, R. L. (2014). *Test equating, scaling, and linking: Methods and practices* (3rd ed.). Springer Science & Business Media.
- Lambert, M. J., Gregersen, A. T., & Burlingame, G. M. (2004). The Outcome Questionnaire-45. In M. E. Maruish (Ed.), *The use of psychological testing for treatment planning and outcomes assessment: Volume 3: Instruments for adults (3rd ed)* (pp. 191-234). Lawrence Erlbaum Associates Publishers. <http://search.ebscohost.com/login.aspx?direct=true&db=psyh&AN=2004-14941-006&site=ehost-live>
- Lambert, M. J., & Harmon, K. L. (2018). The merits of implementing routine outcome monitoring in clinical practice. *Clinical Psychology: Science and Practice*, 25(4), e12268. <https://doi.org/10.1111/cpsp.12268>
- Ley, P. (1972). *Quantitative aspects of psychological assessment* (Vol. 1). London: Duckworth.
- Lord, F. M., & Wingersky, M. S. (1984). Comparison of IRT true-score and equipercentile observed-score" equatings". *Applied Psychological Measurement*, 8(4), 453-461. <https://doi.org/10.1177/014662168400800409>
- McCall, W. A. (1922). *How to measure in education*. MacMillan.
- Mellenbergh, G. J. (2011). *A Conceptual Introduction to Psychometrics: Development, Analysis and Application of Psychological and Educational Tests*. Eleven International. <https://books.google.es/books?id=jRJYAAACAAJ>

- Miller, S. D., Hubble, M. A., Chow, D., & Seidel, J. (2015). Beyond measures and monitoring: Realizing the potential of feedback-informed treatment. *Psychotherapy*, 52(4), 449-457.  
<https://doi.org/10.1037/pst0000031>
- Patel, S. R., Bakken, S., & Ruland, C. (2008). Recent advances in shared decision making for mental health. *Current Opinion in Psychiatry*, 21(6), 606-6012.  
<https://doi.org/10.1097/YCO.0b013e32830eb6b4>
- Rosseel, Y. (2012). Lavaan: An R package for structural equation modeling and more. Version 0.5–12 (BETA). *Journal of Statistical Software*, 48(2), 1-36. <https://doi.org/10.18637/jss.v048.i02>
- Schalet, B. D., Cook, K. F., Choi, S. W., & Cella, D. (2014). Establishing a common metric for self-reported anxiety: Linking the MASQ, PANAS, and GAD-7 to PROMIS Anxiety. *Journal of Anxiety Disorders*, 28(1), 88-96. <https://doi.org/10.1016/j.janxdis.2013.11.006>
- Schalet, B. D., Revicki, D. A., Cook, K. F., Krishnan, E., Fries, J. F., & Cella, D. (2015). Establishing a common metric for physical function: Linking the HAQ-DI and SF-36 PF subscale to PROMIS® Physical Function. *Journal of General Internal Medicine*, 30(10), 1517-1523.  
<https://doi.org/10.1007/s11606-015-3360-0>
- Scherpenzeel, A. C. (2018). “True” Longitudinal and Probability-Based Internet Panels: Evidence From the Netherlands. In *Social and behavioral research and the Internet* (pp. 77-104). Routledge.
- Scherpenzeel, A. C., & Bethlehem, J. G. (2011). How representative are online panels? Problems of coverage and selection and possible solutions. In M. Das, P. Ester, & L. Kaczmarek (Eds.), *Social and behavioral research and the Internet: Advances in applied methods and research strategies* (pp. 105-132). Taylor & Francis.
- Slade, M., Amering, M., & Oades, L. (2008). Recovery: an international perspective. *Epidemiology and Psychiatric Sciences*, 17(2), 128-137. <https://doi.org/10.1017/S1121189X00002827>
- Smits, N. (2016). On the effect of adding clinical samples to validation studies of patient-reported outcome item banks: a simulation study. *Quality of Life Research*, 25(7), 1635-1644.  
<https://doi.org/10.1007/s11136-015-1199-9>
- Smits, N., Ögreden, O., Garnier-Villareal, M., Terwee, C. B., & Chalmers, R. P. (2020). A study of alternative approaches to non-normal latent trait distributions in item response theory models used for health outcome measurement. *Statistical Methods in Medical Research*, 29(4), 1030-1048. <https://doi.org/10.1177/0962280220907625>
- Stasinopoulos, M. D., Rigby, R. A., & Bastiani, F. D. (2018). GAMLSS: a distributional regression approach. *Statistical Modelling*, 18(3-4), 248-273.  
<https://doi.org/10.1177/1471082X18759144>
- ten Klooster, P. M., Oude Voshaar, M. A. H., Gandek, B., Rose, M., Bjorner, J. B., Taal, E., Glas, C. A. W., van Riel, P. L. C. M., & van de Laar, M. A. F. J. (2013). Development and evaluation of a crosswalk between the SF-36 physical functioning scale and Health Assessment Questionnaire disability index in rheumatoid arthritis. *Health and Quality of Life Outcomes*, 11, 199-199. <https://doi.org/10.1186/1477-7525-11-199>
- Terluin, B., Smits, N., Brouwers, E. P. M., & de Vet, H. C. W. (2016). The Four-Dimensional Symptom Questionnaire (4DSQ) in the general population: scale structure, reliability, measurement invariance and normative data: a cross-sectional survey. *Health and Quality of Life Outcomes*, 14(1), 130. <https://doi.org/10.1186/s12955-016-0533-4>
- Terluin, B., van Marwijk, H. W., Adèr, H. J., de Vet, H. C., Penninx, B. W., Hermens, M. L., van Boeijen, C. A., van Balkom, A. J., van der Klink, J. J., & Stalman, W. A. (2006). The Four-Dimensional Symptom Questionnaire (4DSQ): a validation study of a multidimensional self-report questionnaire to assess distress, depression, anxiety and somatization. *BMC Psychiatry*, 6(1), 1. <https://doi.org/10.1186/1471-244X-6-34>
- Timman, R., de Jong, K., & de Neve-Enthoven, N. (2017). Cut-off scores and clinical change indices for the Dutch Outcome Questionnaire (OQ-45) in a large sample of normal and several psychotherapeutic populations. *Clinical Psychology & Psychotherapy*, 24(1), 72-81.  
<https://doi.org/10.1002/cpp.1979>

- Timmerman, M. E., Voncken, L., & Albers, C. J. (2020). A tutorial on regression-based norming of psychological tests with GAMLSS. *Psychological Methods*.  
<https://doi.org/10.1037/met0000348>
- van der Laan, J. (2009). *Representativity of the LISS panel*. Statistics Netherlands.
- Van Hoeck, K. J. M., Lilien, M. R., Brinkman, D. C., & Schroeder, C. H. (2000). Comparing a urea kinetic monitor with Daugirdas formula and dietary records in children. *Pediatric Nephrology*, 14(4), 280-283. <https://doi.org/10.1007/s004670050759>
- van Stralen, K. J., Dekker, F. W., Zoccali, C., & Jager, K. J. (2012). Measuring Agreement, More Complicated Than It Seems. *Nephron Clinical Practice*, 120(3), c162-c167.  
<https://doi.org/10.1159/000337798>
- Wahl, I., Löwe, B., Bjorner, J. B., Fischer, F., Langs, G., Voderholzer, U., Aita, S. A., Bergemann, N., Brähler, E., & Rose, M. (2014). Standardization of depression measurement: a common metric was developed for 11 self-report depression measures. *Journal of Clinical Epidemiology*, 67(1), 73-86. <https://doi.org/10.1016/j.jclinepi.2013.04.019>
- Wall, M. M., Park, J. Y., & Moustaki, I. (2015). IRT Modeling in the presence of zero-Inflation with application to psychiatric disorder severity. *Applied Psychological Measurement*, 39(8), 583-597. <https://doi.org/10.1177/0146621615588184>

## Open Science

We report on a reanalysis of data about which has been published before. We report and refer to previous publications on how we determined our sample size, all data exclusions, all data inclusion/exclusion criteria, whether inclusion/exclusion criteria were established prior to data analysis, all measures in the study, and all analyses including all tested models. If we use inferential tests, we report exact p values, effect sizes, and 95% confidence or credible intervals.

Open Data: The information needed to reproduce all of the reported results are not openly accessible.

Open Materials: The information needed to reproduce all of the reported methodology is made available in and can also be requested from the first author. *We have uploaded an annotated version of our R-code with two practice datasets on <https://www.psycharchives.org/>, which will allow other researchers to apply it to these data (and their own data).*

Preregistration of Studies and Analysis Plans: This study was not preregistered.