

Factor Validation and Rasch Analysis of the Individual Recovery Outcomes Counter

Journal:	<i>Disability and Rehabilitation</i>
Manuscript ID	TIDS-03-2017-022.R2
Manuscript Type:	Research Paper
Keywords:	Mental health, Recovery, Factor analysis, Rasch measurement theory, Validity

SCHOLARONE™
Manuscripts

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Implications for Rehabilitation

Agencies and services working with people with mental health problems aim to help them with their recovery.

The Individual Recovery Outcomes Counter has been developed and is used widely in Scotland to help service users track their progress to recovery.

Using a large sample of routinely collected data we have demonstrated that a number of modifications are needed if the tool is to adequately measure recovery.

This will involve consideration of the scoring system, item content and inclusion, and theoretical basis of the tool.

Abstract

Objective: The Individual Recovery Outcomes Counter is a 12-item personal recovery self-assessment tool for adults with mental health problems. Although widely used across Scotland, limited research into its psychometric properties has been conducted. We tested its measurement properties to ascertain the suitability of the tool for continued use in its present form. **Materials and methods:** Anonymised data from the assessments of 1,743 adults using mental health services in Scotland were subject to tests based on principles of Rasch measurement theory, principal components analysis and confirmatory factor analysis. **Results:** Rasch analysis revealed that the 6-point response structure of the Individual Recovery Outcomes Counter was problematic. Re-scoring on a 4-point scale revealed well-ordered items that measure a single, recovery-related construct, and has acceptable fit statistics. Confirmatory factor analysis supported this. Scale items covered around 75% of the recovery continuum; those individuals least far along the continuum were least well addressed. **Conclusions:** A modified tool worked well for many, but not all, service users. The study suggests specific developments are required if the Individual Recovery Outcomes Counter is to maximise its' utility for service users and provide meaningful data for service providers.

Keywords: Mental health, recovery, factor analysis, Rasch measurement theory, validity

Introduction

Recovery, the concept that people can successfully negotiate periods of mental distress or diagnosed mental disorder to lead rich and fulfilling lives, is one of the key drivers of contemporary mental health policy [1,2]. The core recovery construct relates to goals that focus on development of personal meaning rather than on symptom eradication or cure, and it is often conceptualised as a personal and social journey [3,4,5]. Specific elements that have been reported as relevant and important include support from others, finding hope, engaging in meaningful activities, incorporating difficulties with illness or disability, overcoming stigma, taking control, managing mental health symptoms, empowerment, and citizenship [6-8]. The extent to which recovery has become mainstream since the publication of Anthony's [9] seminal paper is evidenced by its embedding in national strategy documents [10,11], and in the presence of multiple literature reviews published in international journals [7,8,12-15].

The growing traction of recovery-informed approaches to service provision has been aided by, and has itself fuelled demand for, the development of instruments that can quantify service user progress towards the goal of personal recovery. However, the psychometric properties of individual measures is variable, and there is little consensus about which are the strongest [16,17]. There is a need to further test existing measures of recovery using rigorous techniques in order to provide evidence for their continued use, or information about how they can be improved.

The Individual Recovery Outcomes Counter (I.ROC), has been developed and implemented by the Scottish mental health charity Penumbra [18]. The 12-item tool is predicated on four hypothetical or notional domains, each comprising three items (see Figure 1). The acronym HOPE is used to capture description of the four domains (Home, Opportunity, People, Empowerment). The tool was developed to i) measure service user outcomes following a period of service contact in order to demonstrate effectiveness; ii)

Measuring individual mental health recovery 3

1
2
3 allow the client or service user to self-monitor progress throughout the service contact; and
4
5 iii) facilitate therapeutic conversation between service users and support staff by signposting
6
7 areas key to personal recovery.
8

9
10 Despite its widespread use in Scotland, it is important to clarify that development of
11
12 the Individual Recovery Outcomes Counter was conducted through extensive consultation
13
14 with service users and providers rather than via a formal or traditional process of
15
16 psychometric test development; as a result, the four domains described are notional rather
17

18
19 than empirical constructs. The actual empirical research into the psychometric properties of
20
21 the tool has thus far been limited to a validation study using data collected in its first year of
22
23 use among $N=171$ adults [11,18]. Principal components analysis [11] revealed a two factor
24
25 structure (eight *intrapersonal self-reflection/change items* and four *interpersonal*
26
27 *outward/forward looking items*) rather than a replication of the notional four domains. There
28
29 are now considerably more data available, supported by more robust training, data collection
30
31 methods, and quality control measures, to further assess the tool's properties using two
32
33 statistical techniques which, primarily due to sample size, could not be conducted previously:
34
35 Rasch analysis and confirmatory factor analysis.
36
37

38
39 Traditionally, the testing of the psychometric properties of rating scales in mental
40
41 health practice has been based on classical test theory [19]. One of the limitations of classical
42
43 test theory is that, because tests reflect both the ability of the person taking the test
44
45 (personability) and the difficulty of the test (item difficulty), the information generated is
46
47 sample specific and cannot readily be used to compare scores across groups or individuals.
48
49 However, scales based on item response theory, most notably Rasch measurement
50
51 approaches, are increasingly utilised in test development and validation having previously
52
53 been relatively neglected due to the larger sample sizes needed compared with tests based on
54
55 classical test theory [20]. A Rasch model defines how a set of items should perform to
56
57
58
59
60

Measuring individual mental health recovery

4

1
2
3 generate reliable and valid measurements. Rasch analysis is used to quantify the extent to
4
5 which rating scale data (in this case the service user ratings on tool items) fit with predictions
6
7 of those ratings from the Rasch model [21]. Close fit between the predicted and actual scores
8
9 indicates valid measurement [22], and thus Rasch analysis can provide detailed diagnostic
10
11 information about how a scale can be improved, for example through identification of items
12
13 that fail to fit the Rasch model, and score interpretation [21]. Further, it reveals important
14
15 information about a tool's ability to test the full range of a single latent trait, in this case
16
17 'personal recovery', within the context of use, in this case, people who use the services in
18
19 which the Individual Recovery Outcomes Counter has been implemented. The underlying
20
21 assumption of Rasch analysis, therefore, is that the scale under investigation is
22
23 unidimensional. Further, this assumption is tested during Rasch analysis. Any indication that
24
25 the scale violates unidimensionality can be resolved by further testing items that appear to
26
27 comprise separate underlying latent traits as separate unidimensional scales.
28
29
30

31
32 In addition to the Rasch analysis, confirmatory factor analysis is also warranted. The
33
34 previous study of the Individual Recovery Outcomes Counter [23] used principal components
35
36 analysis, an exploratory procedure, to identify factor structure: i.e. latent variables comprising
37
38 subgroups of individual items. However, the hypothesis-testing confirmatory procedure may
39
40 be more appropriate the tool is predicated on a notional four domain structure resulting from
41
42 considerable user-involvement whose validity as an empirical entity this approach could
43
44 confirm or refute. It may seem counterintuitive to simultaneously use both factor analysis to
45
46 explore the potential multidimensionality of the scale and Rasch that assumes that the same
47
48 scale is univariate. However, we have followed advice that where there is a lack of certainty
49
50 about the structure of the items under study, and where dimensionality needs to be examined
51
52 alongside evidence of local dependence, confirmatory factor analysis alone is insufficient,
53
54 and its use in combination with Rasch is warranted [24,25,26]. More rigorous validation of
55
56
57
58
59
60

Measuring individual mental health recovery

5

1
2
3 the notional four domain structure would allow researchers, practitioners, and service
4
5 providers to have added confidence in self-assessment data such that results for particular
6
7 factors, rather than single items, could be interpreted to signify needs for specific targeting.
8

9
10 The aims of the current study therefore were i) to test the extent to which the
11
12 Individual Recovery Outcomes Counter measures one or more latent recovery-related
13
14 variables; ii) to provide information about the validity of individual tool items and determine
15
16 whether and how the the tool might be improved to better measure recovery; iii) to determine
17
18 whether the notional four domain model is empirically supported.
19

20 21 **Materials and Methods**

22 23 **Participants**

24
25 The study forms part of a larger project examining the development and use of the
26
27 Individual Recovery Outcomes Counter and involves secondary analysis of an existing but
28
29 previously unexamined cross-sectional dataset of self-assessment questionnaires completed
30
31 by clients/ users of mental health and related services in Scotland. All aspects of the study
32
33 were reviewed and approved by the Abertay University Research Ethics Committee. Eligible
34
35 participants were adults over 18 years of age who used third sector community mental health
36
37 services in Scotland, and completed at least one assessment, between January 2012 and
38
39

40
41
42 October 2014. Sample size was essentially arbitrary and used all data meeting inclusion
43
44 criteria at the date of collection; however, the final sample exceeded the most conservative
45
46 guidelines on required absolute sample size [27] and subject-to-variable ratio [28] for use in
47
48 confirmatory factor analysis. Participants were receiving services from one of 35 projects
49
50 (*Mdn* participants per service=10, range 1-268) at the time of assessment. Projects included
51
52 supported living services ($k=9$); self-harm services ($k=5$); homeless services and youth
53
54 projects (both $k=2$). Characteristics of study participants are presented in Table 1.
55

56 >>Insert Table 1 about here<<
57
58
59
60

Measuring individual mental health recovery

6

Procedure

Anonymised baseline Individual Recovery Outcomes Counter self-assessment and demographic data held on the computerised database were retrieved on 17th November 2014.

Scores were cleansed manually to remove duplicate or incomplete iterations.

Measure

The Individual Recovery Outcomes Counter was developed by Penumbra, a Scotland-based mental health charity operating in the third/voluntary sector as a mental health service provider. The tool aims to establish, and subsequently track, an individual's level of personal recovery across four notional domains (*Home, Opportunity, People, Empowerment*). Each domain comprises three items, each described using text (descriptions, related terms) and graphical prompts, an issue important to personal recovery as identified during its extensive development (see Figure 1). For each item, the service user is requested to respond on an ordinal 6-point unipolar scale (1 = *never*; 2 = *almost never*; 3 = *sometimes*; 4 = *often*; 5 = *most of the time*; 6 = *all of the time*) with a higher score intended to represent a greater level of personal recovery. Each rating is intended to refer to the past 3-month period. Scores can be plotted on a radar chart that facilitates visualisation of change over time (see Figure 1). In a preliminary validation study [23], the tool had good internal consistency ($\alpha=.86$), correlated significantly with the Recovery Scale [29], an established measure of recovery [30], and one of mental health outcome (BASIS-32 [31,32]). Additionally, it was favoured over either of these measures by service users as a personal recovery outcomes measure. The tool has adequate readability (Flesch-Kincaid score 6.2, unpublished data). Previous analysis of a sample of $N=171$ adults revealed that the tool comprised two factors accounting for 51.8% of total variance in scores: an intrapersonal factor largely relating to the individual's inner life (*mental health, life skills, safety and comfort, physical health, personal network, valuing myself, participation and control, and self-management*); and an interpersonal factor

Measuring individual mental health recovery 7

1
2
3 comprising four items (*exercise and activity, purpose and direction, social network, hope for*
4
5 *the future*) relating to the individual's ability to participate socially, and play a meaningful
6
7 part both in their own lives and in the wider community [18]. Thus, the notional four domains
8
9 measured by the tool were not reflected by the results of this prior analysis.

10
11 >>*Insert Figure 1 about here*<<

12
13
14 Training is provided to support workers employed by Penumbra (including author 2)
15
16 and is mandatory for all those using the tool in practice; further training addresses 'Recovery
17
18 in Practice', 'Coaching for HOPE', 'Planning for HOPE', and 'Motivational Interviewing'. In
19
20 practice, the support worker introduces the individual service user to the Individual Recovery
21
22 Outcomes Counter, and completes a baseline assessment with them during the first four
23
24 scheduled support sessions. Assessment is then repeated quarterly to support individual work
25
26 and track change. The assessment questions facilitate an outcomes-focused conversation
27
28 during which service users can acknowledge progress and identify priorities. Ratings and
29
30 notes reflecting the conversation are manually entered by the practitioner into a secure online
31
32 database managed by Penumbra. In addition, information about age, gender, ethnicity,
33
34 employment status, source of referral, length in service, and reason for leaving service are
35
36 collected.
37
38
39

40
41 The Individual Recovery Outcomes Counter has been used across Penumbra's
42
43 community based services and within a further 16 third sector organisations since 2012.
44
45 While it is used most commonly within mental health focused community settings, it is also
46
47 used within services focusing on related issues including homelessness and substance misuse.

48 49 **Tests of data quality, distribution, stability, scaling assumptions, reliability, and validity**

50
51 Data were examined for data quality (percent missing data for each item in excluded
52
53 data) and for normality of distribution; Hair et al. [28] suggest that data distribution is
54
55 considered normal if skewness is between -2 and +2 and kurtosis is between -7 and +7. Bond
56
57
58
59
60

Measuring individual mental health recovery

8

1
2
3 and Fox [33] suggest that a minimum of 10 responses per scoring category is required to
4
5 determine an estimation of stable threshold values. Descriptive and correlational analyses
6
7 were conducted to evaluate scaling assumptions (e.g., similar item mean scores and
8
9 variances, scores which span the entire measurement continuum, and the magnitude and
10
11 similarity of corrected item-total correlations). Further analysis was conducted of reliability
12
13 and validity, scale-to-sample targeting (score means and standard deviation [*SDs*]; floor and
14
15 ceiling effects), and internal consistency (Cronbach's alpha). Test-retest reliability was
16
17 measured on a separate sample of $n=70$ staff and students from Abertay University (M age =
18
19 26.4 years, $SD=10.7$ years, range 18 to 65 years; 73.1% female) who completed the tool
20
21 twice with an interval of one week. Intraclass correlation coefficients (r) were calculated for
22
23 items and total scores using 2-way mixed models to test the hypothesis that ratings remained
24
25 stable over the brief interval. A value between 0.75 and 1.00 = *excellent*; 0.60 to 0.74 = *good*;
26
27 0.40 to 0.59 = *fair*; and < 0.40 = *poor* [34].
28
29
30
31

32 **Rasch Measurement Testing of the Individual Recovery Outcomes Counter**

33
34 Rasch measurement methods were employed to better establish whether the tool
35
36 captures the full range of the construct of personal recovery in the context of community
37
38 mental health service use. The Rasch model conceptualizes the measurement scale of a
39
40 construct as a ruler; a scale that defines the full range of a construct along its whole
41
42 continuum will comprise scores ranging from ± 4 logits (equivalent to ± 4 standard
43
44 deviations). Both test items and test-takers (i.e., service users) can be located along the scale
45
46 from left to right in terms of their difficulty (less to more) and ability (less to more)
47
48 respectively (see Figure 2). The Rasch model expresses the likelihood that an item that
49
50 represents a given level of the construct of interest will correspond with the perceived level of
51
52 that construct in people with a given level of the construct as a logistic function of the
53
54 difference between item difficulty and person ability [35].
55
56
57
58
59
60

Measuring individual mental health recovery

9

1
2
3 Rasch measurement provides a choice of two models of parameterization for non-
4
5 dichotomous data. The rating scale model specifies that a set of items share the same rating
6
7 scale structure [36] while the partial credit model specifies that each item has its own rating
8
9 scale structure [37]. Model selection requires consideration of whether the thresholds of the
10
11 rating scale are known in advance of data collection (not in the case of the Individual
12
13 Recovery Outcomes Counter), a condition which supports use of the partial credit model
14
15 [37]. Further, previous research, into recovery-related constructs using secondary data [38-
16
17 41] have found the amount of partial correctness to vary across items, again supporting use of
18
19 the partial credit model. We therefore used this model to guide us in establishing whether five
20
21 important indicators of rigorous measurement were met: fit, targeting, dependency,
22
23 multidimensionality, and reliability.
24
25

26
27 **Fit.** To measure the extent to which items in the Individual Recovery Outcomes
28
29 Counter work together to capture the individuals' level of personal recovery we tested the
30
31 performance of each item by visually inspecting for a monotonic ordering of mean item
32
33 ability, item thresholds, and of the item characteristics curves. Further, to measure the item fit
34
35 to the Rasch model, the unweighted mean square outfit statistic and the weighted mean
36
37 square infit statistic were calculated. The outfit statistic is sensitive to unexpected
38
39 observations by person or item, while infit is sensitive to unexpected patterns where residuals
40
41 are close to estimated person abilities [42]. Expected values are close to 1.0 with greater
42
43 values indicating underfit between the items and the model, and lower values indicating
44
45 overfit (i.e., that the data predict the model too well) and hence item redundancy. Scale
46
47 validity is more greatly affected by underfit than overfit. Mean square of 0.6–1.4 represents
48
49 the ideal range [43] items with a value >2.0 are likely to distort or degrade the scale causing
50
51 inaccurate measurement, while those of 1.4 – 2.0 or <.5 are potentially unproductive for the
52
53 measurement but not degrading [44]. Finally, we inspected the indices of person and item
54
55
56
57
58
59
60

Measuring individual mental health recovery

10

1
2
3 reliability and separation which are used to classify people. Low person separation (< 2) or
4
5 person reliability < 0.8) implies that the instrument may be insufficiently sensitive to
6
7 distinguish between high and low performers and hence more items may be needed. Item
8
9 separation and reliability are used to verify the item hierarchy. Low item separation (< 3) or
10
11 item reliability (< 0.9) implies that the person sample is insufficiently large to confirm the
12
13 item difficulty hierarchy of the instrument; this is equivalent to a measure of construct
14
15 validity. In the event of inadequate fit then collapsing of categories is recommended [33].
16
17

18
19 **Targeting.** We examined how people and items were distributed along the proposed
20
21 latent personal recovery continuum, and whether the 12 items covered the full range of the
22
23 continuum and targeted the sample under investigation. This allowed us to gauge the
24
25 calibration of the instrument to the population by comparing graphically how closely the
26
27 amount of personal recovery orientation displayed by the respondents was adequately
28
29 measured by the items on the scale [45]. We also flagged items in similar locations as in need
30
31 of further investigation because of their potential redundancy.
32
33

34
35 **Dimensionality.** Principal component analysis of the residuals is used in Rasch
36
37 measurement theory to test its underlying assumption that all of the data can be explained by
38
39 the latent measures (in this instance personal recovery). This differs from the principal
40
41 components analysis used in classical test theory which is a correlational model that aims to
42
43 identify factors *within* a scale; principal component analysis of the residuals is a hierarchical
44
45 implication model where positive responses to difficult items imply positive responses to
46
47 easy items, but not vice versa [46]. The Rasch model focuses on the unexplained part of the
48
49 data, the residuals, by extracting the common factor that explains the most residual variance.
50
51 Following standardisation of each residual, the noise should, if there is no meaningful
52
53 structure to the residuals and the scale is most likely unidimensional, follow a random normal
54
55 distribution. Identification of a meaningful structure, indicated by Eigenvalues ≥ 2 , suggests
56
57
58
59
60

Measuring individual mental health recovery 11

1
2
3 that the presence of another dimension to the original factor or scale should be investigated
4
5 [47,48,49].
6

7 **Item invariance.** In order to test whether particular groups of people respond to the
8 scale in a systematically different way to others the Differential Item Functioning (DIF)
9 statistic using the Mantel-Haenszel approach was employed [50]. It could, for example, be
10 predicted that females might respond differently to items such as social network. Testing DIF
11 by gender will facilitate conclusions about whether, and how, this manifests. We tested DIF
12 by gender for all 12 scale items. To analyse DIF, item parameters are held constant while
13 person measures are estimated separately for each group. The effect size, the DIF contrast, is
14 reported in logits and is the difference between the two DIF measures; a substantive DIF is
15 ≥ 0.64 logits. The statistical significance is computed using t-tests.
16
17
18
19
20
21
22
23
24
25
26

27 **Reliability.** We assessed reliability using the Person Separation Index [51] which is
28 analogous to the Cronbach's alpha [52]. A value of 0.70 and above is considered acceptable
29 as an indicator for group use, and 0.70 through 0.85 for individual use [51].
30
31
32
33

34 **Construct validity.** Point-measure correlations were calculated to investigate whether
35 all the items within the scale were measuring the same construct. A fundamental concept in
36 Rasch is that higher person measures lead to higher ratings on items and vice versa [44]. The
37 accuracy of this concept is reported by point-measure correlations which should be noticeably
38 positive ($>.50$). All Rasch analyses were conducted using Winsteps ® 3.81.0 software.
39
40
41
42
43
44

45 **Factor structure of the Individual Recovery Outcomes Counter**

46
47 Factors comprise multiple variables whose responses are correlated; that covariance is
48 used to infer the presence of latent variables, also known as factors. Factors may have
49 practical value, for instance their presence may suggest interventions that might be targeted at
50 a particular latent trait. Factor analysis facilitates parsimony since only items related to the
51 overall construct and to one constituent factor need be retained in a scale. The potential
52
53
54
55
56
57
58
59
60

Measuring individual mental health recovery

12

1
2
3 importance in the current case is that, should different factors exist, then low scores on a
4
5 particular factor (e.g., Home items) but not another (e.g., Opportunity items) might help
6
7 services to target resources at issues which are most problematic and which could have
8
9 greatest impact on outcome. In addition, the tool is predicated on a hypothesised model based
10
11 on considerable collaborative development and it is desirable to test that model empirically.
12
13

14 To test for scale-item redundancy Pearson correlations were conducted between all 12
15
16 item scores. Correlations between 0.3 and 0.7 indicate that items are sufficiently related to
17
18 form part of the same latent construct but not so related as to be redundant. Next, a procedure
19
20 similar to exploratory factor analysis, principal components analysis, was conducted to
21
22 determine whether the previously reported two factor structure [18] was replicated in this
23
24 considerably larger sample, or whether the notional four domain structure would now be
25
26 revealed. To assess the appropriateness of the data for factor analysis, a Kaiser-Mayer-Olkin
27
28 measure of sampling adequacy was conducted. A score of $\geq .90$ is described as *excellent* while
29
30 scores $< .50$ are *unacceptable* [53,p.58]. To determine the number of factors to be extracted,
31
32 guidelines described by Costello and Osborne [54] were followed. Multiple analyses were
33
34 conducted; first by setting the number of factors to be extracted as all those with Eigen values
35
36 greater than one; second, analyses were run with number of extracted factors manually set i)
37
38 equal to the number of factors in previously demonstrated analyses (two factors), ii) to the
39
40 number of factors hypothesised by the tool developers (four factors), and iii) to the number of
41
42 factors identified from inspection of the 'elbow' on the accompanying scree plot (also two
43
44 factors). Analyses were also run for the number of factors between one above and one below
45
46 those numbers. As a result, possible solutions ranging from one to five factors were
47
48 considered. Oblique rotation was conducted where extracted factors were significantly
49
50 correlated ($> .32$), and orthogonal rotation where factor correlation was less evident ($< .32$;
51
52 [55]). The most satisfactory factor structure was decided according to i) the smallest number
53
54
55
56
57
58
59
60

Measuring individual mental health recovery 13

1
2
3 of cross-loading items, ii) structure comprising factors of three or more items [54], and iii)
4 acceptable internal reliability of factors indicated by Cronbach's alpha. George and Mallery
5 [56,p.231] suggest $>.9 = excellent$; $>.8 = good$; $>.7 = acceptable$; $>.6 = questionable$; $>.5 =$
6 *poor*; and $< .5 = unacceptable$. Because internal reliability tends to increase with test length
7 [57] we used the Spearman-Brown prophecy formula to calculate the likely internal reliability
8 of each factor scale in the event that it was increased to 10 items of similar quality. Principal
9 components analysis was conducted in IBM SPSS Statistics (V.22.0.0.1).
10
11
12
13
14
15
16
17

18 Finally, confirmatory factor analysis was conducted. While large sample size reduces
19 the problem of multivariate non-normality which might undermine the assumptions of
20 confirmatory factor analysis, we tested this anyway by calculating the Mahalanobis distance
21 in order to identify data outliers and exploring whether they had significant effects on the
22 data structure. Maximum Likelihood Estimation was used to estimate the models' fit using
23 the following indices: Root Mean Square Error of Approximation, the Comparative Fit Index,
24 the Normed Fit Index, the Goodness of Fit Index, and the Adjusted Goodness of Fit Index.
25 For the Root Mean Square Error of Approximation, values < 0.08 and <0.05 reflect
26 reasonable and excellent fits respectively. For the fit indices, values vary along a continuum
27 of 0 to 1 with those >0.9 and >0.95 considered *satisfactory* and *excellent* respectively [58].
28 The Chi-square difference between the model and the data is routinely reported and should be
29 small and non-significant, but is sample size dependent. It should, however, decrease in better
30 fitting models. Since confirmatory factor analysis is intended to be theory-driven rather than
31 exploratory, we proposed to test the four domain model on which the Individual Recovery
32 Outcomes Counter is predicated, and the two factor model, or any close approximation of it
33 revealed in the current study, reported in a previous study [23]. Confirmatory factor analysis
34 was conducted using SSI International LISREL (V.9.10).
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55

56 Results

57
58
59
60

Measuring individual mental health recovery

14

Tests of data quality, distribution, stability, scaling assumptions, reliability, and validity

Analysis revealed a significant amount of missing data. Of $N=2,680$ baseline assessments $n = 937$ (34.9%) had missing data. There was no over-representation of any item among the missing data and, since sample size was not problematic, cases with missing values were deleted listwise leaving a final sample of $N = 1,743$ (see Table 1 for sample characteristics). The most highly scored item, indicating least recovery-related need, was safety and comfort, and the least highly scored item, indicating most need, was social network. Scaling assumptions were verified (see Table 2). Scale scores spanned the measurement continuum; mean item scores were largely dissimilar: repeated measures ANOVA revealed that 57/66 [86.4%] possible pairwise comparisons differed significantly; this is deemed acceptable when the intent of the tool is to extend the range of measurement to cover a wide range of health states [59]. Internal consistency was good (Cronbach's $\alpha = 0.85$). There was little evidence of floor- or ceiling- effects with very few participants reporting either a maximum or minimum scale score. Data for all items, with the single exception of social network, which had a marked positive skew, were normally distributed. Transformation using Log 10 and Square Root methods did not resolve the issue. Mean r for test-retest reliability was .73 ($SD=.06$, range .61 to .82); six items fell within the good ($r=.60$ to .74 range) and six in the excellent ($r>.75$) range ($p<.001$ for all intraclass correlation coefficients). For the total score, $r=.90$ (95% CI .84, .94, $p<.001$).

>>Insert Table 2 about here<<

Rasch measurement testing of the Individual Recovery Outcomes Counter

Fit. At this stage all items fit the model with an overall standardized mean square item fit of .00 ($SD=1.01$), where values of 0 and 1 are expected. However, visual inspection of probability category curves and threshold scores revealed that all 12 items had a similar problem with fit. While the ordinal numbering of the response categories (i.e., 1 through 6)

Measuring individual mental health recovery

15

1
2
3 were congruent with their imputed meaning (i.e., stronger endorsement of an item was
4
5 associated with a higher total score and lower endorsement with the probability of a lower
6
7 total score), Figure 3(A) shows that a response score of 4 was at no point the most probable
8
9 outcome. In addition, for all items except social network, the Andrich threshold values for a
10
11 response of 4 was greater than that for a response of 5. This suggested that a rating of 4 adds
12
13 little of meaning to the item response categories, that they should be discarded, and items re-
14
15 scored. However, a number of re-scoring options were possible; therefore, we calculated fit
16
17 statistics and inspected probability curves for a range of alternative models (Figure 3B-F) to
18
19 determine which, if any, scenario provided the most meaningful information and improved
20
21 the quality of measurement taking place with these data. In detail, Figure 3(A) contains six
22
23 hills, each indicating an original response option of 1-6 (1-2-3-4-5-6). Figure 3(B) illustrates
24
25 five hills representing five category scoring; Figure 3(C) has four hills representing scale
26
27 categories (1[234]56); Figure 3(D) shows responses to a tripartite model ([12][34][56]);
28
29 Figure 3(E) illustrates four responses in the categories (1[23][45]6). Close inspection
30
31 revealed that the four response model in Figure 3(D) depicts response categories that are
32
33 ordered and working as intended. Inspection of infit and outfit statistics for the alternate
34
35 solutions revealed that the 4-category model (Figure 3E) had the least underfit or overfit
36
37 overall. Given that person (>2) and item (>3) separation; and person ($>.8$) and item ($>.9$)
38
39 reliability were acceptable in all variations we therefore decided to re-score using the 4-
40
41 category model depicted in Figure 3(E). Figure 4 shows the ordering of the threshold
42
43 categories for each item using the 4-category model. The item map shows a person's expected
44
45 score for each item as a function of the measure of personal recovery. The x -axis represents
46
47 the theoretical continuum of the latent construct (less to more personal recovery) and the y -
48
49 axis lists the items included. In this case, the six item response categories were collapsed to
50
51 four response categories that are ordered and working as intended, thus further supporting the
52
53
54
55
56
57
58
59
60

Measuring individual mental health recovery

16

1
2
3 case for re-scoring. As a result we have conducted all further analyses of the data based on
4
5 the re-scored data. Normality of distribution was re-examined following re-scoring; data for
6
7 all variables, including social network, now fell within acceptable limits.
8

9
10 With one exception, all scale items in the Individual Recovery Outcomes Counter fell
11
12 within the reasonable range for infit and outfit (mean square =0.6-1.4). Social network was
13
14 marginally outside this range (1.44 and 1.41 respectively) but not to a sufficient amount to
15
16 degrade the scale and were therefore retained. Graphs showing the intraclass correlation
17
18 coefficients were created for all items; observed values were located close to expected values
19
20 with no marked deviation, and within, or at least very close to, 95% confidence intervals.
21

22 >>Insert Figure 2 about here<<

23
24 >>Insert Figure 3 about here<<

25
26
27 **Targeting.** Figure 4 shows the targeting of the sample to the 12 items and reveals that
28
29 they capture just over three quarters (75.8%) of the sample. In particular, items did not
30
31 adequately capture the people who report the lowest levels of personal recovery.
32

33
34 >>Insert Figure 4 about here<<

35
36 **Dependency.** High residual correlations (>.7) [.60-.61] may indicate local item
37
38 dependency between pairs of items or persons. Residual correlations approached but did not
39
40 exceed .7 (*Mdn* = .59, range .47 - .66) suggesting that up to half of random variance between
41
42 items is shared.
43

44
45 **Multidimensionality.** Principal components analysis of the residuals revealed an
46
47 Eigenvalue of 1.65 at the first contrast suggesting that the 12-item scale is unidimensional.
48

49
50 **Item invariance.** The DIF contrasts, comparing male and female responses ranged
51
52 from -0.19 - 0.23, all well below the cut-off of 0.64 logits. Table 3 shows infit and outfit
53
54 statistics. Higher patient measures were associated with higher item ratings; all point measure
55
56
57
58
59
60

Measuring individual mental health recovery

17

1
2
3 correlations were in excess of .5 with the exception of social network (.46) which was
4
5 sufficiently close to the expected correlation (.55) to be non-problematic.
6

7 **Reliability.** Possible scores on the item reliability index lie between 0 and 1 ($>0.8 =$
8 *strongly acceptable*; [62]). Table 3 shows that item reliability on the 4-item scale used was
9
10 1.0. This indicates that items are adequately separating this sample along the measurement
11
12 continuum.
13

14
15
16 >>Insert Table 3 about here<<
17

18 **Factor structure of the Individual Recovery Outcomes Counter**

19
20 **Principal components analysis.** Kaiser-Mayer-Olkin measure of sampling adequacy
21 score was .907 indicating excellent adequacy of the data for factor analysis. The internal
22 reliability for the 12-item scale ($M = 35.6$, $SD = 7.99$, $\alpha = .841$) was good; in a single factor
23
24 'recovery' model 11 items loaded onto that factor $>.5$ (range .44-.72) with only social
25
26 network failing to reach this threshold. Communalities in the single factor model ranged from
27
28 .30 (social network) to .60 (exercise and activity); total communality was 5.60 and the
29
30 percentage of variation explained by the model was .38 (38%). The best fitting model to
31
32 emerge from principal components analysis was very similar to that previously extracted
33
34 from a smaller sample [23] comprising two factors, one of eight and one of four items and
35
36 accounted for 46.63% of total variance. There was no cross loading of factor items following
37
38 rotation. Internal reliability of factor 1 ($\alpha = .809$) and factor 2 ($\alpha = .636$) was good and
39
40 questionable respectively; application of the Spearman-Brown prophecy formula did not
41
42 result in improved reliability. A 3-factor solution explained 55.89% of variance but there was
43
44 significant cross-loading on a number of items. Four and five factor solutions both produced
45
46 some factors with <3 items, and cross-loading items. Internal reliability of the four factors
47
48 hypothesised by the tool's designers after application of the Spearman-Brown prophecy
49
50
51
52
53
54
55
56
57
58
59
60

Measuring individual mental health recovery

18

1
2
3 formula was good (Home $\alpha=.835$; Opportunity $\alpha=.862$; Empowerment $\alpha=.895$); and
4
5 acceptable (People $\alpha=.798$).
6

7 **Confirmatory factor analysis.** Calculation of the Mahalanobis distance revealed that
8 assumption of multivariate normality was not supported because data for 30 (1.7%)
9 individuals could be considered outliers. Subsequent inspection of the Chi-square-
10 Mahalanobis scatter plot and distribution of univariate item data indicated that extreme scores
11 were not due to an anomalous score on one variable. Inspection also indicated that none of
12 the outliers were *error* outliers thus there is a case for removal of outlier data prior to
13 confirmatory factor analysis since failure to meet the assumptions of multivariate normality
14 can lead to overestimation of the Chi-square statistic and thus to increased chance of a Type I
15 error [63]. Following guidelines [64] we checked whether removal of cases changed the
16 resulting model(s) and report findings for the data both with and without removal of outliers.
17 Further, we conducted confirmatory factor analysis on both models using the re-scored 4-
18 category data and the original 6-category data to examine whether model fit worsened under
19 the re-scored data condition. Inspection of 90% confidence intervals for RMSEAs revealed
20 that no model was a significantly better fit than any other. However, under the two different
21 scoring systems, the re-scored 4-category data always provides a better model fit than its
22 equivalent scored on the original 6-category responses. Further, the two factor model
23 revealed in the current study always provides a marginally better fitting model than the
24 notional four domain model. This ultimately culminates in the re-scored 2-factor model
25 which achieves an RMSEA of 0.051 (0.045-0.0565) verging on the margin of the threshold of
26 an excellent fit (RMSEA=0.05). However, four domain model envisaged by the tool's
27 developers was only marginally a less good fit (RMSEA=0.0536, 90%; CI 0.0477-0.0597).
28 The very high correlation (.79) between the factors in the optimum model suggests that they
29 may in fact comprise part of a super-ordinate personal recovery factor and not unrelated
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Measuring individual mental health recovery

19

1
2
3 latent variables. We tested this hypothesis by re-running the confirmatory factor analysis with
4
5 covariance between the two factors set to zero. This resulted in a poorer fitting model
6
7 (RMSEA=0.069) lending support to the idea that the items are all related to a single
8
9 underlying construct.
10

11 Discussion

14 The current study has examined, in detail and with a considerably larger sample than
15
16 previously conducted, the properties of the Individual Recovery Outcomes Counter, an
17
18 outcomes and key-working tool for working with people with mental health problems
19
20 towards their personal recovery. Results can inform the future development of the tool, and
21
22 selection of the most appropriate model for use in clinical practice and outcomes reporting.
23
24 The first important finding, revealed by Rasch analysis, is that the current 6-category scoring
25
26 structure is problematic in the same way for each of the tool's 12 constituent items. A self-
27
28 reported score of 4, representing a response of 'often', was at no point the most likely
29
30 response for any item, thus rendering it redundant. Consideration of the reasons for this are
31
32 warranted since one implication may be that future versions of the tool should use an
33
34 amended scoring structure if it is the original scoring structure itself which is problematic.
35
36 Analysis of alternative re-scoring methods revealed the best fitting solution involved
37
38 collapsing categories 2 and 3, and 4 and 5, resulting in a 4-category structure: 1= *Never*; 2 =
39
40 *Almost never/Sometimes*; 3= *Often/Most of the time*; 4 = *All of the time*. The solution itself
41
42 suggests that the problem lies in the language of the original 6-category structure; while *never*
43
44 and *all of the time* are clear and discrete categories, the others are less so and not obviously
45
46 ordered. More importantly, it is not immediately obvious why the personal recovery construct
47
48 should be rated solely along a continuum of frequency. While it is in keeping with the
49
50 underlying recovery philosophy to concentrate on positive rather than negative aspects, e.g.,
51
52 by measuring in terms of possession of positive attributes or strengths rather than on
53
54
55
56
57
58
59
60

Measuring individual mental health recovery

20

1
2
3 limitations, this should also encompass the intensity of the experience [65]. In other words,
4
5 while one might 'often' experience good mental health and wellbeing, it does not follow that
6
7 the relatively infrequent experience of less good mental health is not deeply distressing if the
8
9 episode is characterised by great intensity. Hence, while re-scoring items in the current study
10
11 appears to have improved model fit for a frequency-based model, the tool could potentially
12
13 be improved further by development of a scoring structure which also requires consideration
14
15 of personal recovery-related intensity. Nevertheless, in the current study the 4-category re-
16
17 scoring produced a well-ordered tool with good fit statistics.
18
19

20
21 The second important finding is that the re-scored items measured only a portion of
22
23 the personal recovery continuum. In particular, the Individual Recovery Outcomes Counter
24
25 was unable to adequately target around a fifth of the sample who had total scores below item
26
27 thresholds. More positively, the tool's items do successfully target around 75% of the
28
29 relevant service user population and, therefore, has the potential to successfully track change
30
31 for those who are furthest on their recovery journey. Nevertheless, the implication is that, for
32
33 individuals in the bottom 20%, their true level of recovery is not captured; this means that
34
35 those with most recovery-related needs are unlikely to be identified through use of the tool
36
37 and therefore cannot be targeted for more intensive interventions. Further, those well below
38
39 threshold scores are unlikely to demonstrate progress along the scale compared with those
40
41 just below those thresholds; this might prove demoralising both for service users and workers
42
43 when apparent progress fails to be captured. The conclusion to be drawn is that individual
44
45 items may need to be amended in order to better target the full range of service users for
46
47 whom it intends to have relevance. At this point, a note of caution is warranted since an
48
49 amended scoring category-structure system that allows expression of recovery-intensity
50
51 might, in itself, solve this problem and, hence, we suggest it will be beneficial to progress any
52
53 tool redevelopment incrementally. Considerations to be made in adjusting, deleting or adding
54
55
56
57
58
59
60

Measuring individual mental health recovery 21

1
2
3 items should include whether specific types of item will improve targeting, and whether there
4
5 is item redundancy. In respect of the former, priority should be given to development of items
6
7 that are sensitive to reduced levels of personal recovery.
8

9
10 While we conducted exploratory and confirmatory factor analysis ostensibly to test
11
12 competing theories about the factor structure of the Individual Recovery Outcomes Counter ,
13
14 we found only partial support for either. Further, it might be considered that even partial
15
16
17 support contradicts indications from the Rasch analysis that the tool is unidimensional. While
18
19 a 2-factor intrapersonal/interpersonal structure was supported to an extent, the component
20
21 factors were highly correlated suggesting that both represent a single, super-ordinate factor.
22
23 Further, while the structure was very similar to that revealed in Monger et al.'s [23] analysis
24
25 it was not identical. The relocation of the hope for the future item from what had been
26
27 interpreted as an interpersonal, outward/forward looking factor to an intrapersonal self-
28
29 reflection/change factor and that of the physical health item in the opposite direction is
30
31 worthy of consideration. The former change, in particular, brings into question the 'forward
32
33 looking' aspect of this interpretation; indeed, no obvious intuitive solution has occurred to us
34
35 and this in our view further strengthens the case for a unidimensional scale.
36
37

38 We suggest that a pragmatic rather than prescriptive approach is warranted;
39
40 essentially, rather than make categorical statements about the tool as definitively
41
42 unidimensional or multidimensional, findings should be used as a diagnostic aide and guide
43
44 for its future development. Currently, the weight of evidence supports the unidimensionality
45
46 of the tool. Nevertheless, the notional four domain structure is user-friendly and provides an
47
48 intuitively appealing approach which provides an aide-memoire for both workers and service
49
50 users. However, statistically it does not provide the best explanation of the data. The 2-factor
51
52 structure, similarly, does not provide the best empirical solution but does suggest future
53
54 routes of development. Given that development of the tool has not included statistical testing
55
56
57
58
59
60

Measuring individual mental health recovery

22

1
2
3 of a larger pool of potential constituent items it is possible that the emergence of a 2-factor
4
5 structure in the principal components analysis might be strengthened by the addition of more
6
7 items similar to those in factor 2. However, while not statistically significant, those items
8
9 found to constitute the previously-titled interpersonal factor in principal components
10
11 analysis are definitely positioned to the left of the personal recovery continuum (the fifth,
12
13 seventh, ninth and twelfth most 'difficult' items) relative to the previously-titled intrapersonal
14
15 items. We can conclude, then, that inclusion of more of this type of item is unlikely to extend
16
17 the targeting range of the tool. We also note at this juncture that two pairs of items (valuing
18
19 myself and mental health, and life skills and personal network) sit at the same part of the
20
21 personal recovery continuum as one another and there may be some redundancy.
22
23

24
25 Since the Individual Recovery Outcomes Counter is intended to inform the
26
27 therapeutic key-worker - service user dialogue as well as functioning as an outcome tool, any
28
29 added or amended items will need to be of true clinical as well as of statistical value;
30
31 conversely, removal of items should be done cautiously where they retain clinical value.
32
33

34 Simply adding items to the tool for statistical expedience may not add practical utility and
35
36 should be avoided since this would simply increase item redundancy [66]. A clear corollary
37
38 of this is that further generation of ideas for new items should be led by experts by experience
39
40 once they have been apprised of the findings of the current study. Nevertheless, examination
41
42 of the wider mental health recovery outcomes literature suggests some potentially fruitful
43
44 issues to consider for inclusion might be an item related to work, since the current item
45
46 purpose and direction may not sufficiently capture the unique value to potentially be made by
47
48 gainful paid employment. A potential barrier is the relatively small proportion (12%) of
49
50 service users engaged in paid employment which might exert a ceiling effect on responses. It
51
52 is also necessary to consider that, while many people with mental health conditions desire to
53
54 be employed, it may be counterproductive to set expectations that are perceived to exert
55
56
57
58
59
60

Measuring individual mental health recovery 23

1
2
3 pressure to take unsuitable employment. Candidate variables might also include items such as
4
5 'Giving it back' [67]. Developers should aim not to make any new version unwieldy and,
6
7 given the status of the tool in shaping the therapeutic conversation, should first and foremost
8
9 be guided by what is helpful to the patient and the relationship. We suggest that evidence for
10
11 the notional four domain structure is sufficiently strong for it to be retained in key-working
12
13 materials and outcomes tools. However, it may be more useful from a development
14
15 perspective to consider the tool as comprising two recovery-related strands that requires some
16
17 balancing in favor of intrapersonal-type items. The challenge, of course, will be in integrating
18
19 these approaches when these and scoring-related modifications are made. Further
20
21 examination of factor structure will be warranted following any modifications made to the
22
23 Individual Recovery Outcomes Counter.
24
25
26

27
28 Given the changes required, it may seem moot at present to address whether the
29
30 current findings on factor structure can inform clinical use of the tool. For example, could
31
32 either the two factor intrapersonal/interpersonal or scores derived from the notional four
33
34 domain structure be used to inform interventions or to be meaningfully reported in outcomes
35
36 data? Clearly this would be an aim of further development but at present it is not possible to
37
38 recommend this approach until issues of targeting, scoring, and potential multidimensionality
39
40 have been resolved. Finally, since the scale-items accounted for less than half (47%) of the
41
42 variance in the sample scores, it should be fruitful to listen carefully to what users say about
43
44 other issues that affect their personal sense of recovery and help them with these issues where
45
46 possible. Recording of additional issues might aid identification of potential new items.
47
48

49 50 **Limitations**

51
52 The study depended upon routinely collected assessment data gathered in the day-to-
53
54 day work milieu. There was a considerable number of incomplete assessments; it may be
55
56 unavoidable in the context of routine assessment with people with mental health problems
57
58
59
60

Measuring individual mental health recovery

24

1
2
3 that a considerable proportion do not initially engage with services. Nevertheless, we could
4
5 detect no systematic reason for non-engagement in terms of demographic variables. The data
6
7 collected here was, in fact, subject to some quality control in terms of training. A further
8
9 potential limitation is the patchiness of some of the descriptive demographic and clinical data
10
11 whose inclusion would have allowed us to better describe the sample and to test DIF for
12
13 variables other than gender: it is possible that responses may have differed between, for
14
15 example, people with psychotic disorders and those with anxiety and non-psychotic
16
17 depressive episodes. Future work should aim to gather more complete demographic and
18
19 clinical data to better address this question.
20
21

22 **Conclusion**

23
24
25 The current study provides further support for the use of the Individual Recovery
26
27 Outcomes Counter as a unidimensional measurement of recovery; this is qualified support,
28
29 however, and further development is required in order for the tool to more adequately capture
30
31 the recovery of those who are at most disadvantage. This paper sets out a programme of work
32
33 to achieve this.
34
35
36
37
38

39 Acknowledgments: The authors would like to thank Nigel Henderson and Jane Cumming
40
41 (Penumbra)
42
43

44 Declaration of interest: 'The authors report no conflicts of interest'.
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

References

- [1] Slade M. Policy rationale. In: M. Slade. *Personal Recovery and Mental Illness* (p.74-77). Cambridge (UK): Cambridge University Press; 2009.
- [2] Tomes N. The patient as a policy factor: A historical case study of the consumer/survivor movement in mental health. *Health Aff.* 2006;25:720-729. doi: 10.1377/hlthaff.25.3.720
- [3] Deegan PE. Recovery and empowerment for people with psychiatric disabilities. *Soc Work Health Care.* 1997;25:11-24. doi:10.1300/J010v25n03_02
- [4] Jacobson N, Greenley D. What Is Recovery? A conceptual model and explication. *Psychiatr Serv.* 2001;52:482-485. doi: 10.1176/appi.ps.52.4.482
- [5] Topor A, Borg M, Di Girolamo S, et al. Not just an individual journey: social aspects of recovery. *Int J Soc Psychiatry.*2011;57:90-99. doi: 10.1177/0020764009345062
- [6] Davidson L, O'Connell MJ, Tondora J, et al. (2005). Recovery in serious mental illness: a new wine or just a new bottle? *Prof Psychol Res Pract.* 2005;36:480-487. doi:10.1037/0735-7028.36.5.480
- [7] Stickley T, Wright N. The British research evidence for recovery, papers published between 2006 and 2009 (inclusive). Part One: a review of the peer-reviewed literature using a systematic approach. *J Psychiatr Ment Health Nurs.* 2011a;18:247-256. doi:10.1111/j.1365-2850.2010.01662.x
- [8] Stickley T, Wright N. The British research evidence for recovery, papers published between 2006 and 2009 (inclusive). Part Two: a review of the grey literature including book chapters and policy documents. *J Psychiatr Ment Health Nurs.* 2011b;18: 297-307. doi: 10.1111/j.1365-2850.2010.01663.x

Measuring individual mental health recovery

26

- 1
2
3 [9] Anthony WA. Recovery from mental illness: the guiding vision of the mental health
4 service system in the 1990s. *Psychosoc Rehabil J.* 1993;16:11-23. doi:
5 <http://dx.doi.org/10.1037/h0095655>
6
7
8
9
10 [10] Carson J, McManus G, Chander A. Recovery: a selective review of the literature and
11 resources. *Ment Health Soc Incl.* 2010;14:35-44. doi:
12 <http://dx.doi.org/10.5042/mhsi.2010.0068>
13
14
15
16 [11] Ion R, Monger B, Hardie S, et al. A tool to measure progress and outcome in recovery.
17 *Br J Ment Health Nurs.* 2013;2:211-215. doi:10.12968/bjmh.2013.2.4.211
18
19
20 [12] Le Boutillier C, Chevalier A, Lawrence V, et al. (2015). Staff understanding of
21 recovery-orientated mental health practice: A systematic review and narrative
22 synthesis. *Implement Sci.* 2015;10:87-87. doi:10.1186/s13012-015-0275-4
23
24
25
26
27 [13] Salzman-Erikson M. (2013). An integrative review of what contributes to personal
28 recovery in psychiatric disabilities. *Issues in Mental Health Nursing*, 34, 185-191.
29 doi:10.3109/01612840.2012.737892
30
31
32
33 [14] Shepherd A, Doyle M, Sanders C, et al. Personal recovery within forensic settings:
34 Systematic review and meta-synthesis of qualitative methods studies. *Crim Behav*
35 *Ment Health.* 2016;26:59-75. doi:10.1002/cbm.1966
36
37
38
39
40 [15] Van Lith T, Schofield MJ, Fenner P. Identifying the evidence-base for art-based
41 practices and their potential benefit for mental health recovery: A critical review.
42 *Disabil Rehabil.* 2013;35:1309-1323. doi:10.3109/09638288.2012.732188.
43
44
45
46 [16] Scheyett A, DeLuca J, Morgan C. Recovery in severe mental illnesses: A literature
47 review of recovery measures. *Soc Work Res.* 2013;37:286-303.
48 doi:10.1093/swr/svt018
49
50
51
52
53 [17] Shanks V, Williams J, Leamy M, et al. Measures of personal recovery: A systematic
54 review. *Psychiatr Serv.* 2013;64:974-980. doi:10.1176/appi.ps.005012012
55
56
57
58
59
60

Measuring individual mental health recovery

27

- 1
2
3 [18] Monger B, Hardie SM, Ion R, et al. The Individual Recovery Outcomes Counter:
4 Preliminary validation of a personal recovery measure. *Psychiatri*.2013;37:221-227.
5 doi:10.1192/pb.bp.112.041889
6
7
8
9
10 [19] Tractenberg RE. Classical and modern measurement theories, patient reports, and
11 clinical outcomes. *Contemp Clin Trials*. 2011;31:1-3. doi: 10.1016/S1551-
12 7144(09)00212-2
13
14
15
16 [20] Urbina S. *Essentials of psychological testing*. Hoboken (NJ): John Wiley & Sons Inc;
17 2004.
18
19
20
21 [21] Cano SJ, Hobart JC. The problem with health measurement. *Patient Preference*
22 *Adherence*. 2011;5:279 –290. doi: <http://dx.doi.org/10.2147/PPA.S14399>
23
24
25 [22] Cano SJ, Mayhew A, Glanzman AM, et al. Rasch analysis of clinical outcome
26 measures in spinal muscular atrophy. *Muscle Nerve*. 2014;49:422–430. doi:
27 <http://dx.doi.org/10.1002/mus.23937>
28
29
30
31
32 [23] Monger B, Hardie SM, Ion R, et al. The Individual Recovery Outcomes Counter:
33 Preliminary validation of a personal recovery measure. *Psychiatri*. 2013;37:221-
34 227. doi:10.1192/pb.bp.112.041889
35
36
37
38 [24] Christensen KB, Engelhard Jr, J. Salzberger, T. Ask the experts: Rasch vs. factor
39 analysis. *Rasch Measur Trans*. 2012;26:1373-1378.
40
41
42
43 [25] Waugh RF, Chapman ES. An analysis of dimensionality using factor analysis (true-score
44 theory) and rasch measurement: What is the difference? Which method is better? *J*
45 *App Measur*.2005;6:80–99.
46
47
48
49 [26] Yu CH, Popp SO, Digangi S, Jannasch-Pennell A. Assessing unidimensionality : a
50 comparison of Rasch Modeling,Parallel Analysis, and TETRAD. *Pract Assessment*
51 *Res Eval*. 2007;12:1–19.
52
53
54
55
56
57
58
59
60

Measuring individual mental health recovery

28

- 1
2
3 [27] Comrey AL, Lee HB. A first course in factor analysis. Hillsdale (NJ): Erlbaum;
4
5 1992.
6
7 [28] Hair JFJ, Anderson RE, Tatham RL, et al. Multivariate data analysis. Cambridge (UK):
8
9 Pearson; 2010.
10
11 [29] Giffort D, Schmook A, Woody C. Recovery Assessment Scale. Chicago (IL):
12
13 Department of Mental Health; 1995.
14
15 [30] Corrigan PW, Salzer M, Ralph RO, et al. Examining the factor structure of the
16
17 Recovery Assessment Scale. Schizophr Bull. 2004;30:1035–1041.
18
19 doi:10.1093/oxfordjournals.schbul.a007118
20
21
22 [31] Sederer LI, Dickey B, Eisen SV. Behavior and Symptom Identification Scale
23
24 (BASIS-32). In: LI. Sederer & B. Dickey (Eds). Outcomes Assessment in Clinical
25
26 Practice (p.65–69). Baltimore,(MD): Williams & Wilkins; 1996.
27
28
29 [32] Eisen SV, Wilcox M, Leff HS. Assessing behavioural health outcomes in outpatient
30
31 programs: reliability and validity of the BASIS-32. J Behav Health Serv Res.
32
33 1999;26:5-17. doi: 10.1007/BF02287790
34
35
36 [33] Bond TG, Fox CM. Applying the Rasch model. fundamental measurement in the human
37
38 sciences (2nd Edition). New York (NY): Routledge; 2007.
39
40
41 [34] Cicchetti DV. Guidelines, criteria, and rules of thumb for evaluating normed and
42
43 standardized assessment instruments in psychology. Psychol Assess.1995;6:284–290.
44
45 doi:10.1037/1040-3590.6.4.284 doi:10.1093/oxfordjournals.schbul.a007118
46
47 [35] Rasch G. Probabilistic models for some intelligence and attainment. Chicago (IL):
48
49 University of Chicago Press; 1980.
50
51 [36] Massof DW. Is the partial credit model a Rasch model? J Appl Measur. 2012;13:114-
52
53 131.
54
55
56
57
58
59
60

Measuring individual mental health recovery

29

- 1
2
3 [37] Linacre JM. Partial Credit Models" (PCM) and "Rating Scale Models" (RSM).
4
5 Rasch Measur Trans. 2000;14:768.
6
7 [38] Barbic SP, Bartlett SJ, Mayo NE. Emotional vitality in caregivers: application of
8
9 Rasch Measurement Theory with secondary data to development and test a new
10
11 measure. Clin Rehabil. 2015;29:705-716. doi:10.1177/0269215514552503
12
13 [39] Barbic SP, Kidd SA, Davidson L, et al. Validation of the brief version of the Recovery
14
15 Self-Assessment (RSA-B) using Rasch measurement theory. Psychiatr Rehabil J. 2015
16
17 38:349-358. doi: 10.1037/prj0000139
18
19 [40] Covic T, Pallant J, Conaghan P, et al. A longitudinal evaluation of the Center for
20
21 Epidemiologic Studies-Depression scale (CES-D) in a rheumatoid arthritis
22
23 population using Rasch analysis. Health Qual Life Outcomes. 2007;5:41. doi:
24
25 10.1186/1477-7525-5-41
26
27 [41] Pallant JF, Tennant A. An introduction to the Rasch measurement model: An
28
29 example using the Hospital Anxiety and Depression Scale (HADS). Br J Clin
30
31 Psychol. 2007;46:1–18. doi: 10.1348/014466506X96931
32
33 [42] Linacre JM. Understanding Rasch measurement: Optimizing category
34
35 effectiveness. J Appl Measur. 2002;3:85-106 .
36
37 [43] Wright BD, Linacre JM, Gustafson J, et al. Reasonable mean-square fit values.
38
39 Rasch Measur Trans. 1994; 8:370.
40
41 [44] Linacre JM. Winsteps Rasch tutorial 2 [Internet]. 2012 . [cited 2015 Jan 5]. Available
42
43 from: <http://www.winsteps.com/a/winsteps-tutorial-2.pdf>
44
45 [45] Wright BD, Masters GN. Rating scale analysis. Chicago (IL): MESA Press; 1982.
46
47 [46] Sick J. Rasch measurement in language education part 6: Rasch measurement and
48
49 factor analysis. JALT Test Evaluation SIG Newsl. 2011;15:15-17.
50
51
52
53
54
55
56
57
58
59
60

Measuring individual mental health recovery

30

1
2
3 [47] Linacre JM. Structure in Rasch residuals: why principal components analysis

4
5 (PCA)? Rasch Measur Trans. 1998; 12:636.

6
7 [48] Linacre JM. Data variance explained by Rasch measures. Rasch Measur Trans.

8
9 2006;20:1045.

10
11 [49] Smith Jr. EV. Detecting and evaluating the impact of multidimensionality using

12
13 item fit statistics and principal component analysis of residuals. J Appl Measur.

14
15 2002;3:205-231.

16
17 [50] Linacre JM. DIF - DPF - bias - interactions concepts [Internet]. no date. [cited 2015 Jan

18
19 15]. Available from: <http://www.winsteps.com/winman/difconcepts.htm>

20
21 [51] Andrich D. An index of person separation in latent trait theory, the traditional KR. 20

22
23 index, and the Guttman scale response pattern. Educ Res Perspect. 1982;9:95-104.

24
25 [52] Cronbach LJ. Coefficient alpha and the internal structure of tests. Psychom.

26
27 1951;16:297-334.

28
29 [53] Stewart DW. The application and misapplication of factor analysis in marketing

30
31 research. J Mark Res, 1981;18:51-62.

32
33 [54] Costello A, Osborne J. Best practices in exploratory factor analysis: Four

34
35 recommendations for getting the most from your analysis. Pract Assess Res

36
37 Evaluation, 2005;10:173-178.

38
39 [55] Tabachnick BG, Fidell LS. Using multivariate statistics. Upper Saddle River (NJ):

40
41 Pearson, Allyn, and Bacon; 2007.

42
43 [56] George D, Mallery P. SPSS for Windows step by step: A simple guide and reference.

44
45 11.0 update (4th ed.). Boston (MA): Allyn & Bacon; 2003.

46
47 [57] Wells C, Wollack J. An instructor's guide to understanding test reliability [Internet].

48
49 2003. [cited 2016 December 2]. Available from: <http://testing.wisc.edu/Reliability.pdf>

50
51 [58] Byrne B. Structural Equation Modeling with AMOS. New York (NY): Routledge; 2010.

Measuring individual mental health recovery 31

- 1
2
3 [59] Ware JE, Gandek B. (1998). Methods for testing data quality, scaling assumptions, and
4 reliability: The IQOLA project approach. *J Clin Epidemiol.* 1998;51:945-952. doi:
5 [http://dx.doi.org/10.1016/S0895-4356\(98\)00085-7](http://dx.doi.org/10.1016/S0895-4356(98)00085-7)
6
7
8
9
10 [60] Yen, WM. Effects of local item dependence on the fit and equating performance of the
11 three-parameter logistic model. *Appl Psychol Measur.* 1984;8:125-145. doi:
12 10.1177/014662168400800201
13
14
15
16 [61] Yen WM. Scaling performance assessments: Strategies for managing local item
17 dependence. *J Educ Measur.* 1993;30:187-213. doi: 10.1111/j.1745-
18 3984.1993.tb00423.x
19
20
21
22
23 [62] Fox CM, Jones JA. Uses of Rasch modeling in counseling psychology research. *J Couns*
24 *Psychol.* 1998;45:30-45. doi: <http://dx.doi.org/10.1037/0022-0167.45.1.30>
25
26
27
28 [63] Curran PJ, West SG, Finch JF. The robustness of test statistics to nonnormality and
29 specification error in confirmatory factor analysis. *Psychol Methods.* 1996;1:16-29.
30 doi: <http://dx.doi.org/10.1037/1082-989X.1.1.16>
31
32
33
34 [64] Aguinis H, Gottfredson RK, Joo H. Best-practice recommendations for defining,
35 identifying, and handling outliers. *Organ Res Methods.* 2013;16:270-301. doi:
36 10.1177/1094428112470848
37
38
39
40 [65] Dela Cruz AM, Bernstein IH, Greer TL, et al. Self-rated measure of pain frequency,
41 intensity, and burden: psychometric properties of a new instrument for the assessment
42 of pain. *J Psychiatr Res.* 2014;59:155–160. doi:
43 <http://doi.org/10.1016/j.jpsychires.2014.08.003>
44
45
46
47
48
49 [66] Streiner DL, Norman GR. Health measurement scales: a practical guide to their
50 development and use. New York (NY): Oxford University Press; 1989.
51
52
53
54
55
56
57
58
59
60

[67] Ridgway P. Re-storying psychiatric disability: learning from first person recovery

narratives. *Psychiatr Rehabil J.* 2001;24:335-343. doi:

<http://dx.doi.org/10.1037/h0095071>

Figure captions

Figure 1: Individual Recovery Outcomes Counter notional four domain (HOPE) model with example scoring and item descriptors

Figure 2: Service user-Item Threshold Distribution

Accompanying text:

Distribution of IROC items obtained by converting raw scores into logits. The x -axis represents the level of personal recovery continuum from low to high. The top bars (above the x -axis) represent the distribution of people in the sample, whereas the bottom bars (below the x -axis) represent items. Ideal targeting would depict a range of item and people covering the whole breadth of the scale. The figure shows that there are few items (bottom bars) that are covering the people (top bars) at the low end of the continuum (those lowest on the continuum of personal recovery). The scale is well-targeted at those scoring between -1.2 and +1.8 logits. The measurement gaps in the personal recovery continuum are shown by the 2-way block arrows. Those scoring above the measurement threshold represent 2.9% of the total sample while those scoring below represent 21.2% of the sample. Therefore, a solution is required to capture the lower end of the personal recovery construct for this sample. In addition, several items are capturing the level of personal recovery of the same subgroup of patients.

Key: ¹. Social Network; ². Valuing myself; ³. Mental health ⁴. Purpose & direction; ⁵. Hope for the future; ⁶. Physical health; ⁷. Self management; ⁸. Exercise & activity; ⁹. Participation & control; ¹⁰. Personal network; ¹¹. Life skills; ¹². Safety & Comfort.

Figure 3: Probability category curves of item 1 (Mental Health)

Accompanying text:

In all examples the responses are correctly ordered, but in the original scoring model (A) a score of 4 is at no time most probable. Models B to F represent alternative re-scoring methods. Collapsing scores of 4 and 5 (Model B), 3,4, and 5 (C), 2 and 3 and 3 and 4 (D) and 1 and 2, 3 and 4, and 5 and 6 (E) all lead to ordered solutions. Visual inspection suggests Model D to be the most satisfactory and is supported by inspection of fit statistics (See Table 3).

Figure 4: Item map

Accompanying text:

Item map showing an individual's expected score to each item as a function of the measure of personal recovery. The x -axis represents the theoretical continuum of the latent construct (less to more personal recovery) measured in logits. The y -axis lists the IROC items in terms of more (Safety and comfort) to less (Social network) personal recovery. In this case the 6-item response categories were collapsed into 4 (i.e., collapse responses 2 & 3 and 4 & 5 into one category each (2 and 3) and transform response 6 to 4). The figure depicts response

Measuring individual mental health recovery

33

categories that are ordered and working as intended, suggesting a 4 item category response may be more favourable for this sample.

Table 1

Participant Characteristics

		<i>n</i>	(%)
Gender	Male	787	(45.1)
	Female	957	(54.9)
Ethnicity	White British	927	(53.2)
	Asian/ Asian mixed/ Asian other	17	(1.0)
	White other	14	(0.8)
	African-Caribbean	4	(0.2)
	Not known/prefer not to say	782	(44.8)
Referral from	NHS	316	(18.1)
	Self/private	287	(16.5)
	Social work	273	(15.7)
	Housing	239	(13.7)
	Community/Independent service	180	(10.3)
	Education authority	137	(7.9)
	General Practitioner	81	(4.6)
	Self	1	(0.1)
	Other	200	(11.5)
Reason for leaving service	Not recorded	70	(4.0)
	Moved	548	(31.4)
	Did not engage	433	(24.8)
	Still in service	412	(23.6)
	No longer meeting client needs	213	(12.2)
	Disengaged with service	134	(7.7)
Employment status	Died	4	(0.2)
	Unemployed	1138	(65.3)
	Student	340	(19.5)
	Employed	209	(12.0)
Length in service (Days)	Other	57	(3.3)
	<i>M</i>	240.3	
	<i>SD</i>	290.1	
Age at assessment (Years)	Range	1-3416	
	<i>M</i>	37.9	
	<i>SD</i>	14.9	
	Range	18-81	

Measuring individual mental health recovery

34

Table 2

Analysis of data quality, scaling assumptions, targeting, and reliability

Psychometric property	Total
Data quality:	
Missing data (%)	34.9
Computable scale scores	1743
Scale assumptions:	
Item scores: <i>M</i> (range)	3.67 (2.17-4.09)
Item <i>SD</i> range	1.16-1.55
Targeting:	
Mean score (<i>SD</i>)	38.24 (10.17)
Possible score range ^a	12-72
Observed score range	12-72
Floor/ ceiling effect ^b	<1/<1
Rating scale score:	Of 20,916 observations:
1 = <i>Never</i>	14%
2 = <i>Almost never</i>	20%
3 = <i>Sometimes</i>	29%
4 = <i>Often</i>	15%
5 = <i>Most of the time</i>	14%
6 = <i>All the time</i>	8%
Reliability:	
Cronbach's alpha	.85
<i>M</i> (<i>SD</i> , Range) inter-item correlation byitem	.32 (.07, .14-.43)
Item-total correlation <i>M</i> (<i>SD</i> , Range)	.61 (.08, .47 - .70)

^a Higher scores represent higher personal recovery ^b Floor effect = % receiving a score of 12 (lowest personal recovery); ceiling effect = % receiving a score of 72 (highest possible personal recovery orientation total score on the original Individual Recovery Outcomes Counter 6-point scale)

Measuring individual mental health recovery

35

Table 3

Item fit statistics

Item	Total Score	Total Count	Measure	SE	Infit		Outfit	
					MNSQ ^a	ZSTD ^b	MNSQ ^c	ZSTD ^d
Mental Health	3,473	1,743	.39	.04	.67	-9.9	.68	-9.9
Life Skills	4,461	1,743	-.66	.04	.82	-6.1	.81	-6.2
Safety & Comfort	4,857	1,743	-1.21	.04	1.30	8.7	1.31	8.8
Physical Health	4,000	1,743	.00	.04	.92	-2.5	.93	-2.1
Exercise & Activity	4,072	1,743	-.10	.04	1.27	7.5	1.27	7.4
Purpose & Direction	3,755	1,743	.37	.04	.99	-.3	.99	-.2
Personal Network	4,437	1,743	-.62	.04	1.23	6.7	1.24	6.9
Social Network	3,043	1,743	1.55	.04	1.46	9.9	1.41	9.9
Valuing Myself	3,548	1,743	.70	.04	.88	-3.5	.88	-3.8
Participation & Control	4,424	1,743	-.61	.04	.95	-1.6	.96	-1.4
Self-management	4,020	1,743	-.03	.04	.67	-9.9	.67	-9.9
Hope For The Future	3,857	1,743	.22	.04	.86	-4.3	.87	-4.2

^a Infit MNSQ (mean square) = outlier sensitive fit statistic; it is sensitive to unexpected observations by service users on items that are relatively very hard or very easy for them (and vice versa); ^c Outfit MNSQ = Inlier-pattern-sensitive fit statistic, sensitive to unexpected responses by service users on items that are roughly targeted on them (and vice versa). Infit and Outfit ZSTD^{b,d} (standardised Z-score or 't-statistic') report statistical significance ($1.96 = p < .05$) of MNSQ values occurring by chance when the data fit the Rasch model. Infit and Outfit MNSQ values close to 1.0 indicate acceptable fit and that items are productive for measurement (1.5-2.0 unproductive for measurement and >2.0 are degrading to the scale measurement). Significant t-statistics can be ignored if MNSQ is acceptable [60]

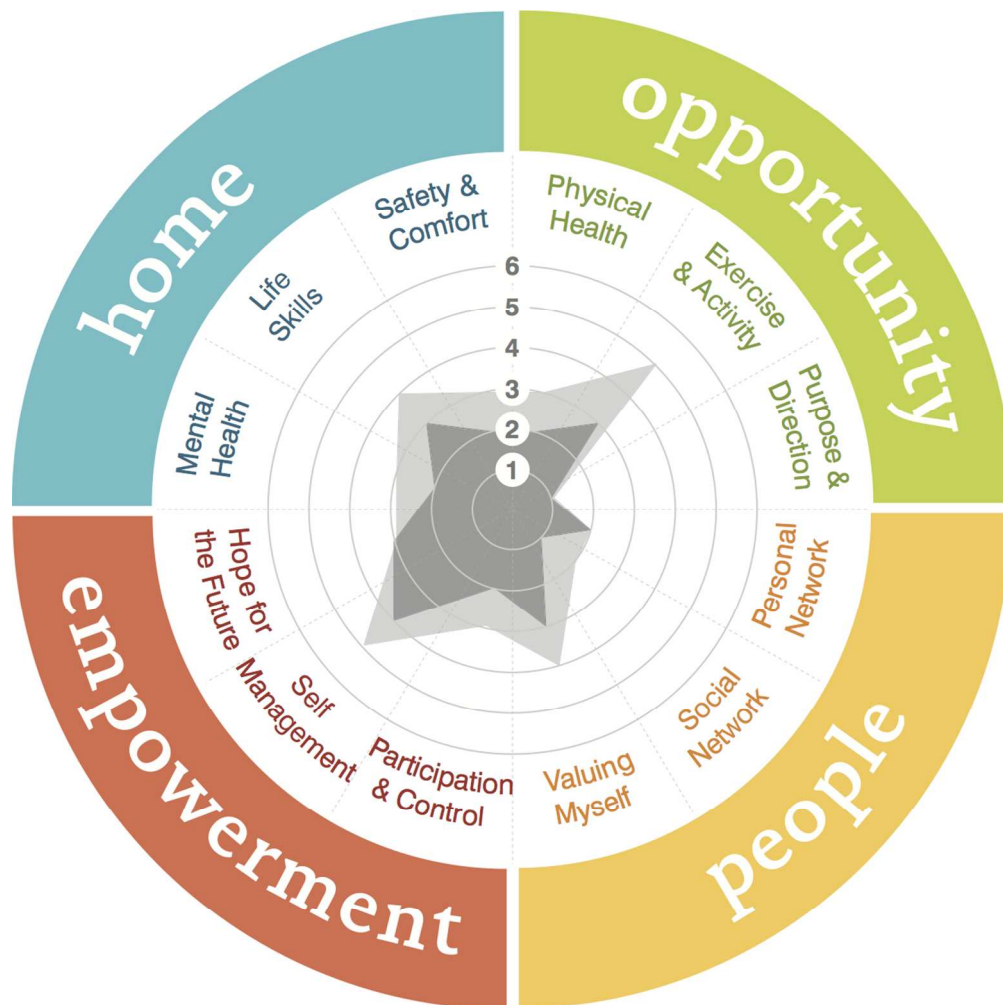


Figure 1: Individual Recovery Outcomes Counter notional four domain (HOPE) model with example scoring and item descriptors

- 1 Mental health: the balance of our physical, emotional, social and spiritual needs: emotions, feelings, optimism, attention, thoughts, beliefs and sense of well-being
- 2 Life skills. The range of skills that we use to cope with the demands of everyday life
- 3 Safety and comfort. Our home should be a place that provides us with safety and comfort, somewhere that we can relax. We should also be able to live in a home that is suitable for us, that we can afford, and that is manageable. We should also feel safe in the area in which we live
- 4 Physical health. Diet, exercise, rest, sleep, illness, pain, if/what we smoke, drink, how well we recover, medication we take, and generally how we look after ourselves
- 5 Exercise and activity. Regularity of exercise or physical activity undertaken
- 6 Purpose and direction. Sense of purpose, of having things to do during the day, a structure
- 7 Personal network. The family/friends/loved ones that are in our lives. People that we can talk to, who are there for us and people who we support (It does not refer to professionals that are paid to support us, including support workers.
- 8 Social network. The connections we have with other people, e.g., groups/clubs we belong to, interests we share with other people, community events/activities we take part in
- 9 Valuing myself. The degree to which we respect ourselves and how we feel about ourselves as a person
- 10 Participation and control. Relates to the degree that we feel we have a say in the decisions that affect our lives
- 11 Self-management. The degree to which we feel able to manage our own health and well-being

12 Hope for the future. How optimistic we feel for our future and how much we are able to look forward.
How positive we are about ourselves and the plans we make

111x111mm (300 x 300 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

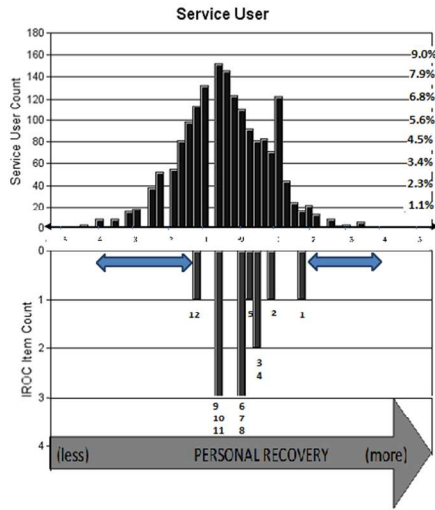


Figure 2: Service user-Item Threshold Distribution

Accompanying text:

Distribution of I.ROC items obtained by converting raw scores into logits. The x-axis represents the level of personal recovery continuum from low to high. The top bars (above the x-axis) represent the distribution of people in the sample, whereas the bottom bars (below the x-axis) represent items. Ideal targeting would depict a range of item and people covering the whole breadth of the scale. The figure shows that there are few items (bottom bars) that are covering the people (top bars) at the low end of the continuum (those lowest on the continuum of personal recovery). The scale is well-targeted at those scoring between -1.2 and +1.8 logits. The measurement gaps in the personal recovery continuum are shown by the 2-way block arrows. Those scoring above the measurement threshold represent 2.9% of the total sample while those scoring below represent 21.2% of the sample. Therefore, a solution is required to capture the lower end of the personal recovery construct for this sample. In addition, several items are capturing the level of personal recovery of the same subgroup of patients.

Key: 1. Social network; 2. Valuing myself; 3. Mental health 4. Purpose & direction; 5. Hope for the future; 6. Physical health; 7. Self management; 8. Exercise & activity; 9. Participation & control; 10. Personal network;; 11. Life skills; 12. Safety & comfort.

272x225mm (96 x 96 DPI)

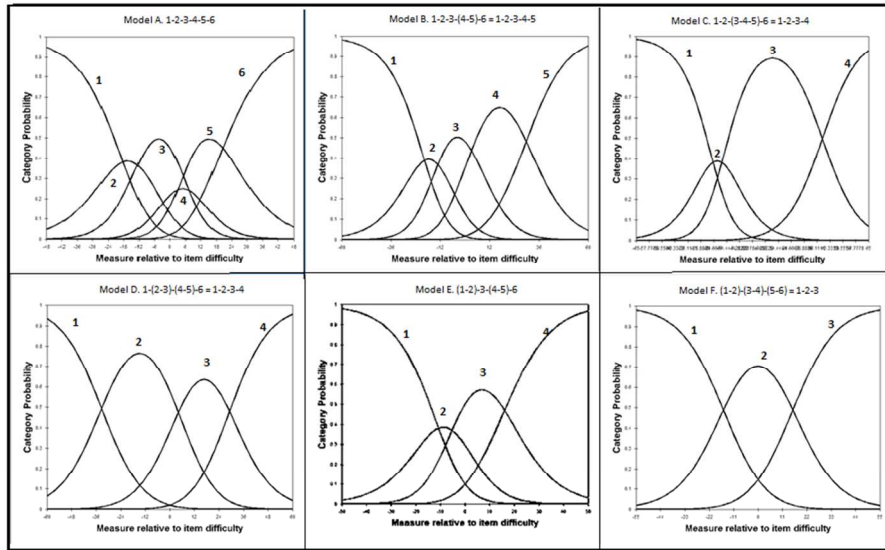


Figure 3: Probability category curves of item 1 (Mental Health)
 Accompanying text:

In all examples the responses are correctly ordered, but in the original scoring model (A) a score of 4 is at no time most probable. Models B to F represent alternative re-scoring methods. Collapsing scores of 4 and 5 (Model B), 3,4, and 5 (C), 2 and 3 and 3 and 4 (D) and 1 and 2, 3 and 4, and 5 and 6 (E) all lead to ordered solutions. Visual inspection suggests Model D to be the most satisfactory and is supported by inspection of fit statistics (See Table 3).

272x200mm (96 x 96 DPI)

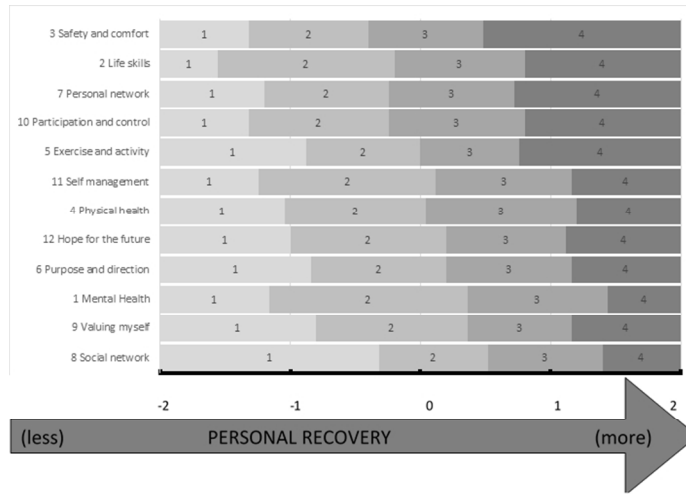


Figure 4: Item map
Accompanying text:

Item map showing an individual's expected score to each item as a function of the measure of personal recovery. The x-axis represents the theoretical continuum of the latent construct (less to more personal recovery) measured in logits. The y-axis lists the I.ROC items in terms of more (Safety and comfort) to less (Social network) personal recovery. In this case the 6-item response categories were collapsed into 4 (i.e., collapse responses 2 & 3 and 4 & 5 into one category each (2 and 3) and transform response 6 to 4). The figure depicts response categories that are ordered and working as intended, suggesting a 4 item category response may be more favourable for this sample.

272x225mm (96 x 96 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60