

## Appendix 2. Criteria for grading PROs based on the guidelines proposed by Khadka et al.<sup>13</sup>

<b>Content Development</b>	
Item identification	A: Comprehensive consultation with the patients, experts (interviews, FGD) and literature review B: Minimal consultation with appropriate patients, experts and literature review C: No consultation with patients
Item selection	A: Pilot instrument developed and tested with Rasch or factor analysis, Statistical justification for reducing items, Items with floor and ceiling effects removed; missing data considered B: Only some of these techniques C: No pilot instrument, or no statistical justification of the items in the final instrument.
<b>CTT based psychometric properties</b>	
Acceptability <sup>170</sup>	A: The percentage of missing data for majority of items : $\leq 5\%$ B: The percentage of missing data for majority of items : $> 5\%$ ; $\leq 40\%$ C: The percentage of missing data for majority of items: $> 40\%$
Targeting <sup>170, 171</sup>	A: End-point responses $\leq 5\%$ for majority of items B: End-point responses $> 5\%$ ; $\leq 40\%$ for majority of items C: End-point responses $> 40\%$ for majority of items
Internal Consistency <sup>174 172</sup>	A: $0.95 \geq \text{Cronbach } \alpha \geq 0.7$ B: $0.7 > \text{Cronbach } \alpha \geq 0.6$ , Or $\text{Cronbach } \alpha > 0.95$ C: $\alpha < 0.6$
Item dependency <sup>170</sup>	A: Inter-item correlations $< 0.3$ B: Inter-item correlations $\geq 0.3$ ; $< 0.6$ C: Inter-item correlations $\geq 0.6$
Dimensionality <sup>12, 173</sup>	A: 1 <sup>st</sup> factor loading $> 0.4$ for all items; Principal component analysis (PCA) - variance explained by the measure $> 60\%$ and eigenvalue of the first contrast $< 2.0$ B: $\text{Cronbach } \alpha : 0.7 < \alpha < 0.9$ ; PCA - variance explained by the measure $\geq 50\%$ to $< 60\%$ ; and eigenvalue $< 2.0$ C: $\text{Cronbach } \alpha : 0.7 > \alpha$ or $\alpha > 0.9$ ; PCA- variance explained by the measure $< 50\%$ ; and eigenvalue $> 2.0$ (indicates multidimensionality)
<b>Rasch based psychometric properties</b>	
Response categories	A: All the categories were ordered or ordering of the categories were obtained after repairing disordered categories and evenly spaced categories B: All the categories were ordered or ordering of the categories were obtained after repairing disordered categories and categories were not evenly spaced. C: Unrepairable disordered categories
Dimensionality	A: PCA of residuals: variance explained by the measure $\geq 60\%$ ; and eigenvalue of the first contrast $< 2.0$ B: PCA of residuals : variance explained by the measure $\geq 50\%$ to $< 60\%$ ; and eigenvalue $< 2.0$ C. PCA of residuals : variance explained by the measure $< 50\%$ ; and eigenvalue $> 2.0$ (indicates multidimensionality)
Measurement Precision	A: Person separation index (PSI) $\geq 2.5$ ; Reliability ( $\alpha$ ) $> 0.85$

	B: $PSI\ 2.0 < 2.50$ ; $0.8 \leq \alpha < 0.85$ C: $PSI \leq 2.0$ ; $\alpha \leq 0.80$
Item fit statistics	A: All items with infit and outfit mean square between 0.70 and 1.30 B: One or two items within the 0.50 and 1.50 limit C: More than two items outside the 0.50 and 1.50 limit
Differential Item functioning (DIF)	A: All items with DIF $< 0.5$ logits B: Some items 0.5 to 1 logits, and one at the most $> 1$ logits C: More than one item $> 1.0$ logit
Targeting	A: Difference between item and person means $\leq 1$ logit B: $> 1$ to $\leq 2$ logits C: $> 2$ logits
<b>Validity</b>	
Convergent validity	A: Tested against appropriate measure, and correlation: (0.3 – 0.9) B: Debatable choice of measure and correlation: (0.3 – 0.9) C: Correlation $< 0.3$ or $> 0.9$
Discriminant validity	A: Tested with an appropriate measure, and correlation: $< 0.3$ B: Debatable choice of measure, and correlation: $< 0.3$ C: Correlation $> 0.3$
Concurrent validity	A: Tested with an appropriate measure and correlation: 0.3 – 0.9 B: Debatable choice of measure, and correlation: 0.3 – 0.9 C: Correlation $< 0.3$ or $> 0.9$
Known group validity	A: Tested between appropriate clinical groups; significant difference between groups. B: Tested between debatable groups; significant difference between groups C: Not tested, or insignificant difference between groups
<b>Repeatability or reproducibility/Responsiveness</b>	
Test retest agreement	A: Intra-class correlation (ICC) $\geq 0.8$ , B: $0.6 \leq ICC < 0.8$ C: $ICC < 0.6$
Interobserver/intermode agreement <sup>12</sup>	A: Limits of agreement (LOA) $<$ Minimally important difference (MID), Weighted kappa $> 0.8$ , Intermodal correlation $> 0.7$ B: LOA broader but still close to MID, Kappa: 0.6 – 0.79, intermodal correlation 0.5-0.7 C: LOA $\gg$ MID, kappa $< 0.6$ , intermodal correlation $< 0.5$ , or incorrect statistical test or inadequate sample ( $< 30$ )
Responsiveness	A: Score changes over time $>$ MID, or changes with intervention; Effect size ( $\geq 1$ ) or responsiveness statistics given B: Changes over time but relationship to MID nor reported, effect size $\geq 0.5$ to $< 1$ ; small sample, and inadequate time frame C: Score changes $\leq$ MID; effect size $< 0.5$

Note: A = Excellent; B = Fair/OK, C = Unsatisfactory