



University of Dundee

Introduction to the special section

Trigg, Andrew; Lenderking, William R; Boehnke, Jan R

DOI:
[10.31234/osf.io/75a8c](https://doi.org/10.31234/osf.io/75a8c)

Publication date:
2023

Licence:
CC BY

Document Version
Publisher's PDF, also known as Version of record

[Link to publication in Discovery Research Portal](#)

Citation for published version (APA):
Trigg, A., Lenderking, W. R., & Boehnke, J. R. (2023). *Introduction to the special section: "Methodologies and considerations for meaningful change"*. PsyArXiv. <https://doi.org/10.31234/osf.io/75a8c>

General rights

Copyright and moral rights for the publications made accessible in Discovery Research Portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from Discovery Research Portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain.
- You may freely distribute the URL identifying the publication in the public portal.

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

1 **Editorial**

2 **Introduction to the special section**

3 **"Methodologies and Considerations for Meaningful Change"**

4
5 Andrew Trigg¹, William R. Lenderking², Jan R. Boehnke³

6
7 ¹ Medical Affairs Statistics, Bayer plc, Reading, United Kingdom

8 ² Patient-Centered Research, Evidera, Bethesda, MD, USA

9 ³ School of Health Sciences, University of Dundee, Dundee, UK

10
11 **Conflicts of interest**

12 AT declares current employment by Bayer plc and prior employment by Adelphi Values and DRG.
13 AT was purposefully not involved in editorial decisions regarding papers submitted by any employees
14 of these companies, and declares no other financial or non-financial conflicts of interest.

15 WRL declares current employment by Evidera, formerly employed at Pfizer (2000-2007). WRL was
16 not involved in editorial decisions involving papers submitted by current Evidera or Pfizer colleagues,
17 and declares no other financial or non-financial conflicts of interest.

18 JRB is Co-Editor in Chief of Quality of Life Research and declares no financial or non-financial
19 conflicts of interest.

20
21
22 **This preprint has not undergone any post-submission improvements or corrections. The**
23 **Version of Record of this article is published in Quality of Life Research, and is available online**
24 **at <https://doi.org/10.1007/s11136-023-03413-1>**

25
26
27 **Corresponding author:**

28 Jan R. Boehnke

29 University of Dundee

30 School of Health Sciences

31 11 Airlie Place

32 Dundee DD1 4HJ

33 United Kingdom

35 The determination of what constitutes a 'meaningful change' on a health outcome measure remains
36 controversial in both methodological and applied research. Motivated by the question of how to
37 understand the efficacy and effectiveness of interventions or the natural history of conditions better
38 (1,2), the concept builds on the widely held belief that statistical significance in itself is not sufficient
39 to establish a treatment benefit (3,4). Since health-related quality of life (HRQL) research should reflect
40 patients' perceptions and evaluations, the topic is of immense theoretical, statistical, and practical
41 relevance. It was therefore timely to offer a space to present discussions, methods, and questions related
42 to this topic, even as new methods and interpretive standards emerge.

43 In collaboration with the Psychometrics Special Interest Group of the International Society for Quality
44 of Life Research (ISOQOL) the editor-in-chief (JRB) developed a call for papers and selected the
45 editorial team (AT and WRL) to take this special issue forward. Submissions closed in April 2021 and
46 invited submissions exploring existing and novel methods for defining meaningful change thresholds
47 for clinical outcome assessments such as patient- or clinician-reported outcome measures. A simulation
48 dataset, described below, was also provided to encourage researchers to evaluate different methods
49 using the same data. The main aim of the special section was to collate a series of methodological and
50 applied articles reflecting current thinking and developments in meaningful change research. And we
51 also wanted to encourage the practice of explicitly stating whether thresholds are intended to support
52 between-group, within-group or within-individual interpretations (3,5–7).

53 For this special section, we broadly define "meaningful change research" as the determination of
54 guidelines for interpretation of the perceived meaning of health outcome score changes or differences
55 based on the patients' (or: the target population's) perception. For a particular score difference (often
56 described as a "threshold") to indicate a "meaningful change" over time, (i) patients (or an appropriate
57 proxy) need to have described this score difference as directional (e.g., improved or deteriorated); and
58 (ii) to a degree that reflects in their eyes a meaningful difference from the previous state (see for
59 example) (3,4). A variety of methods are used to operationalize this, including anchor-based methods
60 or qualitative evaluations of score differences that are perceived as meaningful (8).

61 When working towards concrete operationalizations, the level, type, and magnitude of change need to
62 be specified. For example, it is likely inadmissible to use change thresholds based on group differences
63 to interpret differences between individuals or within individuals over time (7), although this may be a
64 common practice. Table 1 provides an overview of these three key considerations when classifying
65 change and we point out three examples:

- 66 • Minimal within-individual change over time: the smallest amount of change over time a given
67 person must show on an individual level in order to be regarded as having a meaningful
68 change (1B, 2B, 3A);
- 69 • Minimal between-group difference in change over time: the smallest difference between the
70 changes of one group versus another group that are considered meaningful (1A, 2C, 3A);
- 71 • Minimal within-group change over time: the smallest amount of change over time a group of
72 people must show in order to be regarded as having had a meaningful change (1A, 2B, 3A).

73 Other combinations such as cross-sectional between-individual differences are also made in practice
74 (3), in addition to 'larger than minimal' thresholds (9). Similarly, while some definitions focus on
75 changes that 'warrant a change in a patient's management' (12), we do not consider this to be a
76 necessity, as some studies (natural history) do not involve treatment evaluations, yet still must establish
77 a meaningful change. Finally, we consciously avoid the use of specific terms such as 'minimal clinically
78 important difference' (13) or 'minimally important change' within this editorial (4), given these terms
79 have been used interchangeably to describe a range of the combinations arising from Table 1.
80 Standardized terminology is more likely to be achieved through a consensus-based approach in a large
81 group such as the SISAQOL-IMI (14). Until consensus is achieved, it is essential for clarity of
82 communication that all dimensions in Table 1 are clarified in the description of a threshold, e.g.,
83 "minimal within-individual change over time".

84 -----

85 Insert Table 1 about here

86 -----

87 The special section is split into two parts: the first focuses on meaningful change using clinical anchors,
88 the second one presents papers based on what are often called "distribution-based" approaches.
89 Distribution-based approaches are typically described as (i) using measures of cross-sectional or
90 longitudinal (often inter-individual) variability in order to define (ii) a minimal score difference that
91 would be seen as exceeding the level of measurement error (or otherwise nuisance or negligible
92 variability) given a particular psychometric model (15). These thresholds have no connection to
93 (external) evaluations of "meaningfulness" of that particular score difference. It is for this reason that
94 regulators such as FDA have historically stated that distribution-based approaches cannot be used as
95 the sole basis for establishing a responder definition (16). Instead, the assumption is that score
96 differences that are greater than measurement error are due to a more systematic factor or factors, hence
97 the inference of meaning. Their singular advantage in this context is that they do not depend on finding
98 a suitable external clinical anchor, which can be challenging for some applications, but can be calculated
99 solely using data from the measure being evaluated. In contrast, an index of meaningful change would
100 offer information about 'meaningfulness' by either providing information about the connection to a
101 criterion of change or by offering a clear content-based operationalization of meaningfulness (be it
102 qualitative or quantitative). However, when such a criterion is not available, distribution-based methods
103 can be useful. Furthermore, in this special section, the submissions were of high quality, and their
104 inclusion offers the opportunity to contrast the approaches, and the contribution of these methods is too
105 important to leave out of a special section such as this. Additionally, they have an established history
106 of use for the study of individuals over time (i.e., idiographic research) to complement trends at the
107 group level (i.e., nomothetic research) (5,17,18).

108 Finally, we want to thank Pip Griffiths (Digital Medicine Society; IQVIA; SeeingTheta) for providing
109 the simulated dataset that two articles used to illustrate their approaches (19,20), and which could be
110 interesting for readers to explore some of the issues raised in this special section further. The simulated
111 dataset comprises responses to the twelve-item 'Simulated Disease Questionnaire' for 2,000 individuals
112 at four time points. The items have four response categories where higher scores indicate worse health
113 (graded response model). Responses to a seven-category transition rating (i.e., global impression of

114 change) were also simulated at the follow-up time points (for more details please refer to
115 <https://osf.io/khmzg/>).

116

117 THE SPECIAL SECTION

118 The response to the call was enthusiastic, with twenty-seven submissions exploring a range of
119 conceptual and practical issues, of which fifteen are now brought together in this special section. Ten
120 of these papers focus on *meaningful change*, and five papers and two letters address *distribution-based*
121 *indices*. The focus of each paper, in terms of meaningful change versus distribution-based indices, and
122 further classification on the level and type of threshold, is provided in Table 2. Two things are clear
123 from this table. First, most papers focus on within-individual change over time. Second, several papers
124 on meaningful change did not precisely specify the magnitude of change (minimal versus greater). For
125 one of these cases, meaningful change was instead conceptualized in terms of hypothetical patient-
126 perceived treatment success (21). For another paper (22) specifying the magnitude, authors used the
127 terms minimal to reflect ratings of ‘a little better’ and meaningful to reflect ratings of ‘better’ and ‘much
128 better’. We recommend future papers are clearer in terms of the intended magnitude, but note that the
129 two options for the magnitude dimension in Table 1 are not exhaustive where options such as patient-
130 perceived treatment success can be of interest.

131

132

133

Insert Table 2 about here

134

135

136 Setting the scene for the first part of the special section is a report of an online survey regarding how
137 clinicians from different disciplines determine individual-level meaningful change on patient reported
138 outcome measures (PROMs) (23). The authors investigated how oncology or mental health clinical care

139 providers who used PROMs in the USA determine whether a patient's symptoms have changed. Most
140 commonly, clinicians compared two consecutive scores, without a visual aid; the use of normative
141 scores was uncommon. This research highlights the importance of aligning meaningful change research
142 with current practice, but also that education in the value of interpretative tools is warranted.

143 The papers in this section investigate the use of anchors for the derivation of meaningful change
144 thresholds. Anchor-based methods are the most widely applied method for estimating meaningful
145 change, but this does not mean they are without problems. In the second paper of this section (24), the
146 authors highlight and discuss five important issues with anchors that should be kept in mind, rather than
147 viewing anchor-based approaches as a perfect gold standard. This article serves as a helpful collection
148 of methodological issues to consider when reading the collected papers. The third paper illustrates a
149 fundamental practical question when determining meaningful change thresholds, but likely also for any
150 threshold determination (10): how scoring rules and ranges limit the usability of group-level minimal
151 important differences in individual-level responder definitions. Based on the example of the EORTC
152 QLQ-C30 subscales, the authors illustrate how the commonly used 10-point change may be misleading,
153 as due to scaling, an individual cannot actually be measured with a 10-point change on any scale. They
154 present considerations (their Figure 2) to further support responder threshold selection.

155 Moving to investigations of the effectiveness of study design and analysis approaches, the
156 fourth paper (25) reports the results of a simulation study to evaluate the importance of the strength of
157 the correlation between the anchor and the clinical outcome assessment measuring change, varying the
158 impact of sample size, change score variability, and anchor correlation strength on the estimation of the
159 meaningful change threshold at the individual and group level. Using receiver operator characteristics
160 and logistic regression analyses, they show that sample size and change score variability are key factors
161 impacting the required anchor correlation, but using an 'acceptable' cut-off of > 0.30 was often
162 insufficient for accurate estimates of individual meaningful change thresholds, and always insufficient
163 for group changes. The fifth paper (19) builds on the simulation dataset that accompanied this call to
164 address the problem that traditional methods of evaluating within-individual change ignore the effects
165 of floor/ceiling effects and measurement error in PROM scores and global (transition) ratings. The team

166 combined the use of a longitudinal graded response model with a transition item to measure latent
167 change. The method produced tighter estimates of meaningful change when compared to traditional
168 methods, with the methods overlapping most when the proportion of responders was about 50% of
169 participants. Extensions of this approach show promise for a range of applications (26,27). The final
170 simulation study in the first part (28) casts a view forward to the papers on distribution-based thresholds,
171 as the team evaluated the effects of sample characteristics commonly observed in clinical trials on four
172 anchor-based threshold selection procedures and two distribution-based ones. In a large simulation
173 design, they found that both methodological choices and clinical characteristics exert influence on the
174 results and conclusions, and they suggest prioritising study designs with strongly responsive endpoints
175 in settings with about 50% anchor-based responders.

176 Moving to empirical papers exploring questions of meaningful change, one team explored if,
177 how and when meaningful change in depressive symptoms occurred during a period of four months
178 through three data sources (18): weekly questionnaires, qualitative reports, and ecological momentary
179 assessment (EMA; five prompts per day). The ‘if’ was assessed in terms of measurement error (weekly
180 level), perceived meaningfulness (qualitative), and statistically significant changes in the modelled
181 trajectory of symptoms. The distinction between sudden and gradual change (the how) and when this
182 occurred varied considerably between methods. This research will help others evaluate what
183 information each method can provide, alone or in combination, when designing studies to assess health
184 changes. It also points to the potential of EMA and experience sampling to increase patient-centeredness
185 and granularity when collecting HRQL data (29). The use of multiple data sources also plays a key role
186 in the three papers concluding this section. One team (22) sought to evaluate the validity of a rheumatoid
187 arthritis flare questionnaire by examining minimal and meaningful within-individual change using three
188 anchors: patient global ratings, physician global ratings, and using a disease activity index in patients
189 with rheumatoid arthritis. They found that patients were most likely to report meaningful improvement,
190 physicians were most likely to report meaningful worsening, with changes in either direction on the
191 disease activity index least likely to be classified as meaningful. Another team (21) utilized a clinicians-
192 then-patients qualitative interview methodology to understand patient priorities for treatment and a

193 threshold to declare treatment success for adult and adolescent patients with alopecia areata and $\geq 50\%$
194 scalp hair loss. This paper details the novel qualitative method of explicitly incorporating patient input
195 into the definition of an individual change threshold and the endpoint of %hair loss. The authors
196 documented that due to extensive discussions online by patients about hair loss issues, they were able
197 to make appropriate ratings of their hair loss that were largely consistent with values provided by
198 clinical experts. The first part ends with a qualitative study to define meaningful change in physical
199 function after weight-loss (30). The team conducted a qualitative study to evaluate how much weight
200 loss would be meaningful hypothetically for overweight and obese individuals, if they were to lose
201 weight. These individuals all agreed that a $\geq 10\%$ weight loss would be associated with a meaningful
202 improvement in their physical functioning, and that a one-point change at the item level of two HRQL
203 instruments would represent a noticeable change.

204 The papers in the second part of the special section focus on *distribution-based indices*. The
205 papers explore how these indices and precision of their recovery are affected by different definitions of
206 the error variance, distributions, and level of uncertainty. The first paper (31) builds upon previous work
207 by the authors (5) proposing approaches for the identification of treatment responders, providing further
208 justification and elaboration for the use of the coefficient of repeatability (also known as the ‘smallest
209 real difference’ (32) or ‘minimally detectable change’ (15)) for within-individual interpretations of
210 statistically significant change. However, rather than focusing on the conventional $p < 0.05$ threshold,
211 the authors explore more liberal thresholds. This article serves as a helpful reminder that significance
212 levels are not fixed, where less strict (i.e., smaller) thresholds will be sufficient in some scenarios. In
213 addition, the paper has two letters attached to it in this same issue, which discuss the interpretation of
214 the attached statistical significance level and the applicability of the index to individual change
215 classification, which are also of interest for other indices and their interpretation. The second paper (33),
216 focuses also on a version of the reliable change index and compares its use based on classical test theory
217 and item response theory. Classical test theory assumes measurement error is constant across the scale
218 range, but item response theory relaxes this assumption. The authors compare these approaches to detect
219 change beyond measurement error, where the item response theory-based thresholds fluctuate above or

220 below the fixed classical test theory threshold in accordance with baseline score. Their Table 4 presents
221 an overview of thresholds for PROMIS shortform users within oncology. Using item response models,
222 another team (34) proposes a method for increasing the precision of measurement of within-individual
223 change. They build on existing approaches to quantify the error associated with individual scores
224 derived from item response theory analyses: using plausible values, the precision of scores across the
225 spectrum of theta (severity of underlying trait) can be incorporated. This can increase the accuracy of
226 measuring intra-individual changes, which is very useful in individuals (for example) with chronic
227 illness who need to be monitored repeatedly over time and provides an extension to more typical
228 distribution-based methods.

229 All PROM scores are subject to measurement error and using raw individual change scores
230 does not account for this fact. The last two papers in the special section use regression and predictive
231 frameworks to derive change metrics that also allow to quantify the uncertainty associated with the
232 estimate. One team (35) presents alternatives to the raw change scores that were developed over 50
233 years ago (36,37), but have so far not been widely used or explored within patient-reported outcome
234 research. The two approaches provide estimates of an individual's true gain after incorporating
235 measurement error, which have both conceptual advantages and greater sensitivity compared to raw
236 change scores. The final paper of the special section (20), compares three distribution-based methods:
237 the reliable change index, one of its variants, and Bayesian regression models that regress post-scores
238 on pre-scores to identify group-level change over time. The article shows that there are only small
239 differences between the methods in detecting change when PROM reliability is high, but none of them
240 outperforms all others if that is not the case. The article offers a technical discussion that compares
241 advantages and disadvantages of these approaches.

242

243 EDITORIAL COMMENTARY

244 In closing, we want to take the opportunity to highlight three topics that struck us when reading and
245 editing the papers. A first observation is that anchor-based methods for within-individual guidelines

246 should be based on finding a threshold separating 'no change' and 'changed' groups on the anchor. The
247 notion of locating a threshold, lying along a continuum of perceived change, is supported by recent
248 research (38). As individuals will vary in their personal threshold, many methods use the mean of these
249 individual threshold locations or derive otherwise a threshold aggregated across individuals (e.g.,
250 receiver operating characteristic curves, logistic regression, discriminant analysis; (4,39)). Similarly,
251 the longitudinal item response model presented within this special section is designed to estimate the
252 location of this threshold (19). Therefore, from a theoretical standpoint we view anchor-based methods
253 such as receiver operating characteristic curves, logistic regression, discriminant analysis and
254 longitudinal item response theory models as useful techniques for identifying a threshold for *within-*
255 *individual* change to identify groups of responders and non-responders. However, regarding estimates
256 of mean score change within an 'improved' anchor group, we maintain that they do not target the
257 location of a threshold and are therefore theoretically biased estimators of within-individual change
258 thresholds (4). Instead, mean change within an 'improved' anchor group has been proposed as more
259 appropriate to guide thresholds for *within-group* changes over time (40,41). Similarly, calculating the
260 difference in mean change in scores between an 'improved' and 'stable' anchor group is not a
261 theoretically appropriate estimator of a within-individual change threshold (38), but instead has been
262 proposed as more suited to between-group differences in change over time (41,42). However,
263 simulations presented within this special section (28) suggest that deviations from normally distributed
264 score changes may pose a challenge to these theoretical ideals. Further planned simulations should help
265 to confirm this (43).

266 A second observation is that current methods for within-individual thresholds and their clinical
267 application use estimators relying on between-individual variability (4,7). For example, meaningful
268 change threshold estimation typically compares between-individual variation in an anchor measure with
269 between-individual variation in change between two assessment points. And distribution-based indices
270 are based on between-individual variance (e.g., standard deviation of a test score multiplied by a
271 constant representing the level of accepted uncertainty and another variable such as the reliability
272 coefficient). If researchers or clinicians are interested in understanding how a group of patients is

273 classified over the course of time (and not making a statement about individual patients), then using
274 measures that are based on between-individual variance is likely an appropriate approach (4). However,
275 if a statement about an individual patient is the goal, then we know that between-individual variability
276 is not always a good or justifiable proxy for within-individual variability (29,44–47). In such a situation,
277 the use of within-individual methods (e.g., EMA or related methods to explore intra-individual variation
278 (18)) might be more appropriate. In the call for papers, we encouraged authors to explicitly justify
279 whether thresholds were intended for between-group, within-group, or within-individual interpretations
280 and why it was appropriate to do so. This has led to calls for more nuance in interpretation (7); to
281 pragmatic responses that within-individual change methodology faces challenges in practical
282 applications ((5); but see (18,29) for contrasting examples); to detailed statements on how to interpret
283 a given index and when and where it is appropriate to use (4); as well as wider discussions and
284 explanations of the methods leading to such indices (19,31,35). We especially see the development of
285 appropriate within-individual methods for the identification of change as a key priority that also aligns
286 with current technological developments for practice.

287 A third point is that in many submissions the variability or uncertainty associated with either
288 the threshold or the change estimate is an important element in interpretation. Knowing the uncertainty
289 associated with a threshold estimate is important, but not always explained or provided. Regardless of
290 the type of variability used and whether a threshold based on meaning to patients or distributions is
291 sought, recognizing and making transparent that there is uncertainty associated with these thresholds is
292 a valuable reminder that none of the methods discussed in this special section offer absolute results.
293 Because the use of meaningful change methodology and distribution-based thresholds has been
294 ritualized to a degree, it is not always considered whether a particular method to determine thresholds
295 is the most appropriate one for a given context. Additionally, emerging mixed methods research relies
296 on classifying particular patients as "changed" for identification in case studies, with limited or no
297 allowance for measurement error, as well as assuming that the classification threshold applies to this
298 particular patient (48,49). Transparency about uncertainty in thresholds and classifications as well as
299 whether it is appropriate to apply a threshold for group or individual change is therefore a key

300 consideration for developing mixed-methods research agendas around how health outcome measures
301 are used by patients more broadly (50–53). We think that this intersection between epistemology,
302 psychometrics, and various fields of clinical practice contains one of the strongest development
303 opportunities for our understanding of (subjective) health outcome measurement, but substantial work
304 is needed to align theories and practices for a coordinated research effort in this area.

305 The call for papers was issued to invite discussion, development, as well as state-of-the-art
306 research and practice. We are grateful for the excellent range of submissions received and to all authors
307 and reviewers involved in selecting the published papers, which represent a two-year collective effort.
308 We hope that readers find these papers useful both in developing their own research, but also to help
309 the field to further extend its efforts around patient-centeredness. When we can all agree on what a
310 meaningful change is and how to measure it for a particular patient, measure, and population, then we
311 will have the opportunity to bring about meaningful change in clinical practice and at the social and
312 policy level.

313

314 **Acknowledgement**

315 This paper is available as a preprint at <https://psyarxiv.com/75a8c>

316 Preprint DOI: 10.31234/osf.io/75a8c

317 License: CC-By Attribution 4.0 International

318 **REFERENCES**

- 319 1. Reeve BB, Wyrwich KW, Wu AW, Velikova G, Terwee CB, Snyder CF, et al. ISOQOL
320 recommends minimum standards for patient-reported outcome measures used in patient-centered
321 outcomes and comparative effectiveness research. *Quality of Life Research*. 2013 Oct
322 1;22(8):1889–905.
- 323 2. Wyrwich KW, Norquist JM, Lenderking WR, Acaster S, the Industry Advisory Committee of
324 International Society for Quality of Life Research (ISOQOL). Methods for interpreting change
325 over time in patient-reported outcome measures. *Qual Life Res*. 2013 Apr;22(3):475–83.
- 326 3. King MT, Dueck AC, Revicki DA. Can Methods Developed for Interpreting Group-level Patient-
327 reported Outcome Data be Applied to Individual Patient Management? *Med Care*. 2019 May;57
328 Suppl 5 Suppl 1:S38–45.
- 329 4. Terwee CB, Peipert JD, Chapman R, Lai JS, Terluin B, Cella D, et al. Minimal important change
330 (MIC): a conceptual clarification and systematic review of MIC estimates of PROMIS measures.
331 *Qual Life Res*. 2021 Oct;30(10):2729–54.
- 332 5. Hays RD, Peipert JD. Between-group minimally important change versus individual treatment
333 responders. *Qual Life Res*. 2021 Oct;30(10):2765–72.
- 334 6. Musoro ZJ, Hamel JF, Ediebah DE, Cocks K, King MT, Groenvold M, et al. Establishing anchor-
335 based minimally important differences (MID) with the EORTC quality-of-life measures: a meta-
336 analysis protocol. *BMJ Open*. 2018 Jan;8(1):e019117.
- 337 7. Trigg A, Griffiths P. Triangulation of multiple meaningful change thresholds for patient-reported
338 outcome scores. *Qual Life Res*. 2021 Oct;30(10):2755–64.
- 339 8. Staunton H, Willgoss T, Nelsen L, Burbridge C, Sully K, Rofail D, et al. An overview of using
340 qualitative techniques to explore and define estimates of clinically important change on clinical
341 outcome assessments. *J Patient Rep Outcomes*. 2019 Dec;3(1):16.
- 342 9. Cocks K, King MT, Velikova G, de Castro G, Martyn St-James M, Fayers PM, et al. Evidence-
343 based guidelines for interpreting change scores for the European Organisation for the Research
344 and Treatment of Cancer Quality of Life Questionnaire Core 30. *European Journal of Cancer*.
345 2012 Jul;48(11):1713–21.
- 346 10. Cocks K, Buchanan J. How scoring limits the usability of minimal important differences (MIDs)
347 as responder definition (RD): an exemplary demonstration using EORTC QLQ-C30 subscales.
348 *Qual Life Res* [Internet]. 2022 Jul 9 [cited 2023 Feb 24]; Available from:
349 <https://link.springer.com/10.1007/s11136-022-03181-4>
- 350 11. Wang Y, Devji T, Qasim A, Hao Q, Wong V, Bhatt M, et al. A systematic survey identified
351 methodological issues in studies estimating anchor-based minimal important differences in
352 patient-reported outcomes. *Journal of Clinical Epidemiology*. 2022 Feb;142:144–51.
- 353 12. King MT. A point of minimal important difference (MID): a critique of terminology and
354 methods. *Expert Review of Pharmacoeconomics & Outcomes Research*. 2011 Apr;11(2):171–84.
- 355 13. Jaeschke R, Singer J, Guyatt GH. Measurement of health status. *Controlled Clinical Trials*. 1989
356 Dec;10(4):407–15.
- 357 14. SISAQOL-IMI | IMI Innovative Medicines Initiative [Internet]. [cited 2023 Mar 24]. Available
358 from: <https://www.imi.europa.eu/projects-results/project-factsheets/sisaqol-imi>

- 359 15. de Vet HC, Terwee CB, Ostelo RW, Beckerman H, Knol DL, Bouter LM. Minimal changes in
360 health status questionnaires: distinction between minimally detectable change and minimally
361 important change. *Health Qual Life Outcomes*. 2006 Dec;4(1):54.
- 362 16. U.S. Department of Health and Human Services FDA Center for Drug Evaluation and Research,
363 U.S. Department of Health and Human Services FDA Center for Biologics Evaluation and
364 Research & U.S. Department of Health and Human Services FDA Center for Devices and
365 Radiological Health. Guidance for industry: patient-reported outcome measures: use in medical
366 product development to support labeling claims: draft guidance. *Health and Quality of Life*
367 *Outcomes*. 2006;4(1):79.
- 368 17. Blampied NM. Reliable change and the reliable change index: still useful after all these years?
369 *tCBT*. 2022;15:e50.
- 370 18. Smit AC, Snippe E, Bringmann LF, Hoenders HJR, Wichers M. Transitions in depression: if,
371 how, and when depressive symptoms return during and after discontinuing antidepressants. *Qual*
372 *Life Res [Internet]*. 2022 Nov 23 [cited 2023 Feb 24]; Available from:
373 <https://link.springer.com/10.1007/s11136-022-03301-0>
- 374 19. Bjorner JB, Terluin B, Trigg A, Hu J, Brady KJS, Griffiths P. Establishing thresholds for
375 meaningful within-individual change using longitudinal item response theory. *Qual Life Res*
376 *[Internet]*. 2022 Jul 23 [cited 2023 Feb 24]; Available from:
377 <https://link.springer.com/10.1007/s11136-022-03172-5>
- 378 20. Li Y. Inferring meaningful change in quality of life with posterior predictive distribution: an
379 alternative to standard error of measurement. *Qual Life Res [Internet]*. 2022 Sep 9 [cited 2023
380 Feb 24]; Available from: <https://link.springer.com/10.1007/s11136-022-03239-3>
- 381 21. Wyrwich KW, Kitchen H, Knight S, Aldhouse NVJ, Macey J, Mesinkovska N, et al. Using
382 qualitative methods to establish the clinically meaningful threshold for treatment success in
383 alopecia areata. *Qual Life Res [Internet]*. 2022 Jul 12 [cited 2023 Feb 24]; Available from:
384 <https://link.springer.com/10.1007/s11136-022-03170-7>
- 385 22. Bartlett SJ, Bykerk VP, Schieir O, Valois MF, Pope JE, Boire G, et al. “From Where I Stand”:
386 using multiple anchors yields different benchmarks for meaningful improvement and worsening
387 in the rheumatoid arthritis flare questionnaire (RA-FQ). *Qual Life Res [Internet]*. 2022 Sep 8
388 [cited 2023 Feb 24]; Available from: <https://link.springer.com/10.1007/s11136-022-03227-7>
- 389 23. Jones SMW, Gaffney A, Unger JM. Common methods of determining meaningful change in
390 clinical practice: implications for precision patient-reported outcomes. *Qual Life Res [Internet]*.
391 2022 Sep 10 [cited 2023 Feb 24]; Available from: [https://link.springer.com/10.1007/s11136-022-](https://link.springer.com/10.1007/s11136-022-03246-4)
392 [03246-4](https://link.springer.com/10.1007/s11136-022-03246-4)
- 393 24. Wyrwich KW, Norman GR. The challenges inherent with anchor-based approaches to the
394 interpretation of important change in clinical outcome assessments. *Qual Life Res [Internet]*.
395 2022 Nov 18 [cited 2023 Feb 24]; Available from: [https://link.springer.com/10.1007/s11136-022-](https://link.springer.com/10.1007/s11136-022-03297-7)
396 [03297-7](https://link.springer.com/10.1007/s11136-022-03297-7)
- 397 25. Griffiths P, Sims J, Williams A, Williamson N, Cella D, Brohan E, et al. How strong should my
398 anchor be for estimating group and individual level meaningful change? A simulation study
399 assessing anchor correlation strength and the impact of sample size, distribution of change scores
400 and methodology on establishing a true meaningful change threshold. *Qual Life Res [Internet]*.
401 2022 Nov 19 [cited 2023 Feb 24]; Available from: [https://link.springer.com/10.1007/s11136-022-](https://link.springer.com/10.1007/s11136-022-03286-w)
402 [03286-w](https://link.springer.com/10.1007/s11136-022-03286-w)

- 403 26. Griffiths P, Terluin B, Trigg A, Schuller W, Bjorner JB. A confirmatory factor analysis approach
404 was found to accurately estimate the reliability of transition ratings. *Journal of Clinical*
405 *Epidemiology*. 2022 Jan;141:36–45.
- 406 27. Terluin B, Koopman JE, Hoogendam L, Griffiths P, Terwee CB, Bjorner JB. Estimating
407 meaningful thresholds for multi-item questionnaires using item response theory. *Qual Life Res*
408 [Internet]. 2023 Feb 13 [cited 2023 Feb 24]; Available from:
409 <https://link.springer.com/10.1007/s11136-023-03355-8>
- 410 28. Qin S, Nelson L, Williams N, Williams V, Bender R, McLeod L. Comparison of anchor-based
411 methods for estimating thresholds of meaningful within-patient change using simulated PROMIS
412 PF 20a data under various joint distribution characteristic conditions. *Qual Life Res* [Internet].
413 2022 Nov 13 [cited 2023 Feb 24]; Available from: [https://link.springer.com/10.1007/s11136-022-](https://link.springer.com/10.1007/s11136-022-03285-x)
414 [03285-x](https://link.springer.com/10.1007/s11136-022-03285-x)
- 415 29. Bringmann LF, van der Veen DC, Wichers M, Riese H, Stulp G. ESMvis: a tool for visualizing
416 individual Experience Sampling Method (ESM) data. *Qual Life Res* [Internet]. 2020/11/23 ed.
417 2020 Nov 22; Available from: <https://www.ncbi.nlm.nih.gov/pubmed/33222049>
- 418 30. Poon JL, Marshall C, Johnson C, Pegram HC, Hunter M, Kan H, et al. A qualitative study to
419 examine meaningful change in physical function associated with weight-loss. *Qual Life Res*
420 [Internet]. 2022 Jul 22 [cited 2023 Feb 24]; Available from:
421 <https://link.springer.com/10.1007/s11136-022-03191-2>
- 422 31. Peipert JD, Hays RD, Cella D. Likely change indexes improve estimates of individual change on
423 patient-reported outcomes. *Qual Life Res* [Internet]. 2022 Aug 3 [cited 2023 Feb 24]; Available
424 from: <https://link.springer.com/10.1007/s11136-022-03200-4>
- 425 32. Beckerman H, Roebroek ME, Lankhorst GJ, Becher JG, Bezemer PD, Verbeek ALM. Smallest
426 real difference, a link between reproducibility and responsiveness. *Quality of Life Research*.
427 2001;10(7):571–8.
- 428 33. Lee MK, Peipert JD, Cella D, Yost KJ, Eton DT, Novotny PJ, et al. Identifying meaningful
429 change on PROMIS short forms in cancer patients: a comparison of item response theory and
430 classic test theory frameworks. *Qual Life Res* [Internet]. 2022 Sep 24 [cited 2023 Feb 24];
431 Available from: <https://link.springer.com/10.1007/s11136-022-03255-3>
- 432 34. Ho EH, Verkuilen J, Fischer F. Measuring individual true change with PROMIS using IRT-based
433 plausible values. *Qual Life Res* [Internet]. 2022 Oct 25 [cited 2023 Feb 24]; Available from:
434 <https://link.springer.com/10.1007/s11136-022-03264-2>
- 435 35. Andrae DA, Foster B, Peipert JD. Comparison of raw and regression approaches to capturing
436 change on patient-reported outcome measures. *Qual Life Res* [Internet]. 2022 Sep 22 [cited 2023
437 Feb 24]; Available from: <https://link.springer.com/10.1007/s11136-022-03196-x>
- 438 36. Cronbach LJ, Furby L. How we should measure ‘change’: Or should we? *Psychological Bulletin*.
439 1970;74(1):68–80.
- 440 37. Lord FM, Novick MR. *Statistical theories of mental test scores*. Reading, MA: Addison-Wesley;
441 1968.
- 442 38. Vanier A, Sébille V, Blanchin M, Hardouin JB. The minimal perceived change: a formal model
443 of the responder definition according to the patient’s meaning of change for patient-reported
444 outcome data analysis and interpretation. *BMC Med Res Methodol*. 2021 Dec;21(1):128.

- 445 39. Gerlinger C, Schmelter T. Determining the non-inferiority margin for patient reported outcomes:
446 Determining the non-inferiority margin for patient reported outcomes. *Pharmaceut Statist*. 2011
447 Sep;10(5):410–3.
- 448 40. Musoro JZ, Coens C, Fiteni F, Katarzyna P, Cardoso F, Russell NS, et al. Minimally Important
449 Differences for Interpreting EORTC QLQ-C30 Scores in Patients With Advanced Breast Cancer.
450 *JNCI Cancer Spectrum*. 2019 Sep 1;3(3):pkz037.
- 451 41. Sabah SA, Alvand A, Beard DJ, Price AJ. Minimal important changes and differences were
452 estimated for Oxford hip and knee scores following primary and revision arthroplasty. *Journal of*
453 *Clinical Epidemiology*. 2022 Mar;143:159–68.
- 454 42. Bell ML, Dhillon HM, Bray VJ, Vardy JL. Important differences and meaningful changes for the
455 Functional Assessment of Cancer Therapy-Cognitive Function (FACT-Cog). *J Patient Rep*
456 *Outcomes*. 2018 Dec;2(1):48.
- 457 43. Vanier A, Leroy M, Hardouin JB. Toward a rigorous assessment of the statistical performances of
458 methods to estimate the Minimal Important Difference of Patient-Reported Outcomes: A protocol
459 for a large-scale simulation study. *Methods*. 2022 Aug;204:396–409.
- 460 44. Fisher AJ, Medaglia JD, Jeronimus BF. Lack of group-to-individual generalizability is a threat to
461 human subjects research. *Proceedings of the National Academy of Sciences*. 2018 Jul
462 3;115(27):E6106–15.
- 463 45. Harvill LM. An NCME Instructional Module on. Standard Error of Measurement. *Educational*
464 *Measure: Issues Practice*. 1991 Jun;10(2):33–41.
- 465 46. McAleavey AA. When (Not) to Rely on the Reliable Change Index [Internet]. *Open Science*
466 *Framework*; 2021 Nov [cited 2023 Feb 25]. Available from: <https://osf.io/3kthg>
- 467 47. Molenaar P, Campbell CG. The new person-specific paradigm in psychology. *Current Directions*
468 *in Psychological Science*. 2009;18(2):112–7.
- 469 48. De Smet M, Acke E, Cornelis S, Truijens F, Notaerts L, Reitske Meganck, et al. Understanding
470 ‘patient deterioration’ in psychotherapy from patients’ perspectives: A mixed methods multiple
471 case study. 2022 [cited 2023 Feb 25]; Available from:
472 <https://rgdoi.net/10.13140/RG.2.2.17796.60802>
- 473 49. Desmet M, Van Nieuwenhove K, De Smet M, Meganck R, Deeren B, Van Huele I, et al. What
474 too strict a method obscures about the validity of outcome measures. *Psychother Res*. 20210204th
475 ed. 2021 Sep;31(7):882–94.
- 476 50. McClimans L. Interpretability, validity, and the minimum important difference. *Theor Med*
477 *Bioeth*. 2011/07/28 ed. 2011 Dec;32(6):389–401.
- 478 51. McClimans LM. First person epidemiological measures: vehicles for patient centered care.
479 *Synthese*. 2021 May;198(S10):2521–37.
- 480 52. Truijens FL, Desmet M, De Coster E, Uyttenhove H, Deeren B, Meganck R. When quantitative
481 measures become a qualitative storybook: A phenomenological case analysis of validity and
482 performativity of questionnaire administration in psychotherapy research. *Qualitative Research in*
483 *Psychology*. 2019;19(1):244–87.
- 484 53. Truijens FL, Van Nieuwenhove K, De Smet MM, Desmet M, Meganck R. How questionnaires
485 shape experienced symptoms. A qualitative case comparison study of questionnaire

486 administration in psychotherapy research. *Qualitative Research in Psychology*. 2021;19(3):806–
487 30.

488

489

490

491

492

493

494 Table 1: Classification of meaningful change thresholds

Dimension	Options
1: Level of interpretation	A: Group B: Individual
2: Type of comparison	A: Difference (cross-sectional) B: Change over time* C: Difference in change over time
3: Magnitude	A: Minimal B: Larger than minimal (e.g., moderate or large)

495 *Note.* *May be further split according to improvement / worsening.

496

497

498

499

500

501 Table 2: Characteristics of included papers.

Paper	Meaningful change or Distribution-based	Level of interpretation	Type of comparison	Magnitude
Andrae et al. (35)	Distribution-based	Individual	Change over time	N/A
Bartlett et al. (22)	Meaningful change	Group	Change over time	Minimal and larger than minimal
Bjorner et al. (19)	Meaningful change	Individual	Change over time	Minimal
Cocks & Buchanan (10)	N/A	Individual	Change over time	N/A
Griffiths et al. (25)	Meaningful change	Group, Individual	Change over time	Minimal
Ho et al. (34)	Distribution-based	Individual, Group	Change over time	N/A
Jones et al. (23)	Meaningful change	Individual	Change over time	Not specified
Lee et al. (33)	Both	Individual	Change over time	Minimal
Li et al. (20)	Distribution-based	Individual	Change over time	N/A
Peipert et al. (31)	Distribution-based	Individual	Change over time	N/A
Poon et al. (30)	Meaningful change	Individual	Change over time (hypothetical)	Minimal
Qin et al. (28)	Meaningful change	Individual	Change over time	Not specified
Smit et al. (18)	Both	Individual	Change over time	Meaningful ^a
Wyrwich & Norman (24)	Meaningful change	General	General	General
Wyrwich et al. (21)	Meaningful change	Individual	Change over time (hypothetical)	Meaningful ^b

502 *Note.* ^a ‘meaningful’ was defined as the patient clearly noticing a change in daily life and/or
 503 experiencing discomfort as a result of the change.

504 ^b ‘meaningful’ was defined as the amount of change that patients considered a (hypothetical)
 505 treatment success.

506

507

508