

Original Article

*These authors contributed jointly to this work



Cite this article: Hobbs C, Lewis G, Dowrick C, Kounali D, Peters TJ, Lewis G (2021). Comparison between self-administered depression questionnaires and patients' own views of changes in their mood: a prospective cohort study in primary care. *Psychological Medicine* **51**, 853–860. <https://doi.org/10.1017/S0033291719003878>

Received: 4 February 2019
Revised: 9 October 2019
Accepted: 20 October 2019
First published online: 20 January 2020

Key words:
Depression; epidemiology; primary care.

Author for correspondence:
Catherine Hobbs, E-mail: c.hobbs@bath.ac.uk

Comparison between self-administered depression questionnaires and patients' own views of changes in their mood: a prospective cohort study in primary care

Catherine Hobbs^{1*} , Gemma Lewis^{2,*} , Christopher Dowrick³,
Daphne Kounali⁴, Tim J. Peters⁵ and Glyn Lewis²

¹Department of Psychology, University of Bath, Bath BA2 7AY, UK; ²Division of Psychiatry, Faculty of Brain Sciences, University College London, London W1T 7NF, UK; ³Institute of Psychology Health and Society, University of Liverpool, Waterhouse Building Block B, Liverpool L69 3BX, UK; ⁴Bristol Medical School, University of Bristol, Oakfield House, Oakfield Grove, Bristol BS8 2BN, UK and ⁵Bristol Medical School, University of Bristol, First Floor, Learning and Research, Southmead Hospital, Bristol BS10 5NB, UK

Abstract

Background. Self-administered questionnaires are widely used in primary care and other clinical settings to assess the severity of depressive symptoms and monitor treatment outcomes. Qualitative studies have found that changes in questionnaire scores might not fully capture patients' experience of changes in their mood but there are no quantitative studies of this issue. We examined the extent to which changes in scores from depression questionnaires disagreed with primary care patients' perceptions of changes in their mood and investigated factors influencing this relationship.

Methods. Prospective cohort study assessing patients on four occasions, 2 weeks apart. Patients ($N = 554$) were recruited from primary care surgeries in three UK sites (Bristol, Liverpool and York) and had reported depressive symptoms or low mood in the past year [68% female, mean age 48.3 (s.d. 12.6)]. Main outcome measures were changes in scores on patient health questionnaire (PHQ-9) and beck depression inventory (BDI-II) and the patients' own ratings of change.

Results. There was marked disagreement between clinically important changes in questionnaire scores and patient-rated change, with disagreement of 51% (95% CI 46–55%) on PHQ-9 and 55% (95% CI 51–60%) on BDI-II. Patients with more severe anxiety were less likely, and those with better mental and physical health-related quality of life were more likely, to report feeling better, having controlled for depression scores.

Conclusions. Our results illustrate the limitations of self-reported depression scales to assess clinical change. Clinicians should be cautious in interpreting changes in questionnaire scores without further clinical assessment.

Introduction

Self-administered screening questionnaires that assess the severity of depressive symptoms have been recommended in UK primary care and in North America and some parts of Europe (Kendrick et al., 2009; Thombs & Ziegelstein, 2014). These recommendations were made in response to concerns that depression is under-diagnosed and under-treated in primary care, with the aim of improving detection and monitoring treatment response. In 2006 the quality outcomes framework (QOF) in the UK encouraged the use of three questionnaires through monetary compensation to practices: the patient health questionnaire (PHQ-9), the beck depression inventory (BDI-II) and the hospital anxiety and depression scale (HADS). These questionnaires are no longer incentivized but remain widely used in UK primary care and continue to influence treatment decisions (Kendrick et al., 2009). The PHQ-9 along with other questionnaires is also used as a routine outcome measure in improving access to psychological therapies (IAPT) services in the UK (Clark et al., 2018).

Self-administered depression questionnaires have been compared to diagnostic assessments and their sensitivity and specificity are fairly good, at around 80% (Beck, Guth, Steer, & Ball, 1997; Moriarty, Gilbody, McMillan, & Manea, 2015). However, their use in clinical settings has been criticized (Dowrick et al., 2009; Toop, 2011). One concern is that changes in scores might not fully capture the patient's experience of improvement or deterioration in their mood. Such disagreement has important implications for treatment decisions and patient-centred care (Malpass et al., 2016; Robinson et al., 2017).

© The Author(s) 2020. This is an Open Access article, distributed under the terms of the Creative Commons Attribution-NonCommercial-ShareAlike licence (<http://creativecommons.org/licenses/by-nc-sa/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the same Creative Commons licence is included and the original work is properly cited. The written permission of Cambridge University Press must be obtained for commercial re-use.

Improvement is commonly rated by clinicians within research settings with the single-item measure, the clinical global impression-improvement scale (CGI-I). Research comparing the CGI-I with routine patient self-report measures in psychiatric settings has found good levels of agreement, although differences were also observed (Berk *et al.*, 2008; Zaider, Heimberg, Fresco, Schneier, & Liebowitz, 2003). Clinicians also routinely ask patients whether their condition has improved, deteriorated or stayed the same (Fischer *et al.*, 1999; Kamper, Maher, & Mackay, 2009). Patient-rated change is measured in research settings with a single-item question, which asks patients retrospectively about how their whole condition has changed compared to a previous occasion, rather than asking about individual symptoms (Fischer *et al.*, 1999; Kamper *et al.*, 2009).

We have conducted qualitative studies of people whose self-rated changes in mood differed from their responses to self-administered depression scales (Malpass *et al.*, 2016; Robinson *et al.*, 2017). Patients explained the disagreement as resulting from the presence of co-morbid conditions, negative and positive life events, changes in social support and changes in quality of life (Robinson *et al.*, 2017). This supports other qualitative findings that patients often state that scales such as the PHQ-9 do not fully capture their experience of illness (Malpass *et al.*, 2016). We are not aware of any similar qualitative or quantitative investigations of this question.

In this study, we used a cohort of patients recruited from primary care to investigate the extent to which responses to the PHQ-9 and BDI-II disagreed with patients' perceptions of changes in their mood, assessed using a patient-rated change scale. We also investigated factors that might influence patient reports of self-improvement having controlled their responses on the PHQ-9 and BDI-II.

Methods

Participants

Participants were recruited from general practice (GP) surgeries in three UK sites: Bristol, Liverpool and York. Computerized records were used to identify patients aged 18–70 who had reported low mood, depressive episodes, depressed mood, depressive symptoms or a major depressive episode in the past year, irrespective of any treatment. We excluded patients who: were diagnosed with bipolar disorder, psychosis or eating disorder; had alcohol or substance use problems; were unable to complete study questionnaires; or were 30 weeks or more pregnant. A total of 7721 patients were sent an information letter and 1470 (19%) replied. Of these, 821 were willing to be contacted, 23 (3%) of whom were ineligible. The remaining 798 were contacted to arrange an interview. Of these, 563 consented (38%) and 559 (38%) were interviewed (four could not be contacted). Data were collected at four time-points, 2 weeks apart (baseline and follow-up 1, 2 and 3). Patients and public representatives were involved in management and steering groups for the PANDA programme grant and gave input into the design, conduct and interpretation of the study.

Ethical approval

All participants provided written informed consent and ethical approval was obtained from NRES Committee South West – Central Bristol. The authors assert that all procedures

contributing to this work comply with the ethical standards of the relevant national and institutional committees on human experimentation and with the Helsinki Declaration of 1975 as revised in 2008.

Measures

Depressive symptoms: The PHQ-9 and BDI-II were completed at each time-point. The PHQ-9 is a 9-item self-administered measure of depressive symptoms in the past 2 weeks and scores range from 0 to 27 (Kroenke, Spitzer, & Williams, 2001). Internal consistency was high at each time-point (Cronbach's α 0.89–0.92). The BDI-II is a 21-item self-administered measure of the severity of depressive symptoms in the past 2 weeks (Beck *et al.*, 1997) and scores range from 0 to 63. Internal consistency was high at each time-point (Cronbach's α 0.93–0.95). Higher scores indicate more severe depressive symptoms.

Patient-rated change: We used a single-item question based on 'Global Rating Scales' that are routinely used in musculoskeletal and chronic pain research and have high reliability and validity (Fischer *et al.*, 1999; Kamper *et al.*, 2009). Participants were asked 'compared to when we last saw you 2 weeks ago how have your moods and feelings changed?' Response options were: 'I feel a lot better' (1), 'I feel slightly better' (2), 'I feel about the same' (3), 'I feel slightly worse' (4), 'I feel a lot worse' (5). We used 'moods and feelings' instead of 'depression' because many people might not consider themselves 'depressed' and this wording should encourage a more general response. Our qualitative studies found evidence that patients viewed this question as more open-ended and explorative, stating that it allowed them to 'sum up' their mental health and express themselves outside of the parameters of the questionnaires (Robinson *et al.*, 2017). The patient-rated change scale was completed twice at each time-point, at the beginning and end of the questionnaire. Test-retest reliability was good with kappa (quadratic weights) of 0.89. The scale, or similar, has been used in prior randomized controlled trials (Button *et al.*, 2015; Malpass *et al.*, 2016; Robinson *et al.*, 2017).

Anxiety: The generalized anxiety disorder assessment (GAD-7) (Spitzer, Kroenke, Williams, & Löwe, 2006) was completed at each time-point and is a 7-item self-administered measure of the severity of anxiety symptoms in the past 2 weeks, scores ranging from 0 to 21. Higher scores indicate more severe symptoms.

Physical and Mental Health-Related Quality of Life: The 12-item short-form health survey (SF-12) (Ware, Kosinski, & Keller, 1996) was administered at each time-point. Separate physical and mental health-related quality of life scores were derived (Ware *et al.*, 1996). Scores range from 0 to 100, higher scores indicating better quality of life.

Negative Life Events: At baseline (only), participants were asked, using a self-administered computerized questionnaire, whether they had experienced the following in the previous 6 months: (i) bereavement, (ii) separation or divorce, (iii) a serious illness or injury, (iv) victimization (mugging, burglary, serious assault), (v) being in trouble with the law, (vi) debt, (vii) a serious dispute with a family member or friend, or (viii) being made compulsorily redundant from work. Due to the low frequency, a binary variable was created (none or 1 or more).

Social Support: At baseline (only), participants completed eight questions as part of the self-administered computerized questionnaire relating to: (i) feeling loved, (ii) having others that

can be relied on, (iii) feeling accepted, (iv) feeling supported, (v) having others to talk to, (vi) having others that make them happy, (vii) having others that care what happens to them and (viii) having others that make them feel an important part of their lives. Each question used a three-point scale (1) not true, (2) partly true and (3) certainly true. Scores were summed and ranged from 1 to 24, higher scores indicating more social support.

Potential confounders: We adjusted for variables previously shown to be associated with depressive symptoms, and site. Demographic variables (age, sex, ethnicity, employment status, financial status and education level) were measured at baseline. Due to small numbers, ethnic minority status was a binary variable. Employment status was categorized as employed, unemployed not by choice and unemployed by choice. Financial status was three categories: low ('Finding it very difficult to make ends meet' and 'Finding it difficult to make ends meet'), medium ('Just about getting by') and high ('Living comfortably' and 'Doing alright'). Education level was seven categories, from no qualifications to a higher degree.

Statistical analyses

Identifying disagreement between questionnaire scores and patient-rated change

To calculate change scores, mean PHQ-9 and BDI-II scores at each follow-up time-point were subtracted from mean scores at the previous time-point (to correspond to the patient-rated change scale which asks about change over the last 2 weeks). Possible change scores ranged from -27 to $+27$ for PHQ-9 and -63 to $+63$ for BDI-II. Greater negative scores indicated improvement and greater positive scores indicated deterioration.

We used the minimal clinically important difference (MCID), the smallest change in symptoms meaningful to patients, to assess the extent of disagreement (Button et al., 2015). The MCID has been estimated in the PANDA cohort to be around a 20% reduction in PHQ-9 or BDI-II scores (manuscript in preparation). We used the MCID of a 20% reduction or increase in questionnaire scores to create the following categories: clinically important decreases (a decline in scores of 20% or more), no clinically important change (a decline or increase in scores smaller than 20%) and clinically important increases (an increase in scores greater than or equal to 20%) (Button et al., 2015). For each response option on the patient-rated change scale, we report the proportion of patients in each of the above MCID categories. As test-retest reliability for the patient-rated change scale was good ($\kappa = 0.89$) we used the rating from the beginning of the assessment. Disagreement did not vary substantially when the rating completed at the end of the assessment were used (online Supplementary Table S3).

We defined disagreement as (i) a clinically important change in PHQ-9/BDI-II scores and a rating of change response that indicated either no change or a change in the opposite direction (ii) no clinically important change in PHQ-9/BDI-II scores and a rating of change response that indicated a change in either direction. The proportion of patients showing some form of disagreement overall was calculated overall by dividing the total number of people showing disagreement with the total number of people. Proportion disagreement was also calculated within each patient-rated response category. Quadratic weighted and unweighted kappa values were used to test agreement between the patient rating of change responses and MCID categories. In

a prior manuscript, we had identified a MCID of 15% for the BDI-II (Button et al., 2015) so we conducted sensitivity analyses with this estimate.

Reliability of disagreement

We further examined the extent of disagreement by tabulating the proportion of participants scoring within each category of the patient-rated change scale with the equivalent proportion scoring a corresponding change on the PHQ-9/BDI-II (online Supplementary Analyses). For example, if 10% of patients reported feeling much better, this was tabulated against the top 10% of change scores on the PHQ-9/BDI-II and so on for the percentage who reported feeling slightly better, the same, slightly worse or worse. Quadratic weighted and unweighted kappa values were used to test agreement between these proportions.

Variables that influence disagreement

We used a binary outcome (feeling better *v.* same or worse) to reflect that neither feeling the same nor worse is a good clinical outcome. As the patient-rated change scale asks about the last 2 weeks, we could construct logistic models with the 2, 4 or 6 week follow-up as the outcome. We adjusted for binary clinically important change (20% change in scores or not) over the previous 2-weeks. This binary variable reduced collinearity between depression scores and other exposures (e.g. anxiety) and was consistent with our approach to clinically important change and disagreement.

For exposures measured at multiple time-points (anxiety, mental and physical-health-related quality of life,) we did a principal components analyses of the exposure at the current and preceding time-points. Principal components analysis (PCA) can be used to transform two correlated variables into orthogonal (uncorrelated) factors or 'principal components.' The first component is a function of the average score on each variable. The second component is uncorrelated with the first and is a function of the difference between two scores (Jolliffe & Cadima, 2016). Models were adjusted for confounders known to be associated with depressive symptoms (age, sex, ethnicity, education level, current use of antidepressants and marital, financial and employment status) and site. All analyses were conducted using STATA 14.

Role of the funding source

The funding source had no role in study design, data collection, data analysis, interpretation or writing of the report. The corresponding author had full access to all data used in the study, and final responsibility for the decision to submit for publication.

Results

Due to extensive missing data at baseline five patients were excluded, leaving 554 for analyses. At follow-ups 1, 2, and 3: 476 (86%); 443 (80%), and 430 (78%) provided data, respectively. Baseline sample characteristics are presented in Table 1. Patients were aged 18–71 (mean 48.30, s.d. 12.56), 68% female and 96% white.

Table 1. Sample characteristics at baseline

Demographic variable	Overall sample (n = 554)
Age, mean (s.d.)	48.30 (12.56)
Female, N (%)	377 (68)
White, N (%)	530 (96)
Married or partnership, N (%)	278 (50)
Employed, N (%)	296 (53)
Higher education, N (%)	161 (29)
ICD-10 depression diagnosis, N (%)	238 (45)
Taking antidepressants, N (%)	377 (69)
Site	
Bristol	197 (36)
Liverpool	188 (34)
York	169 (30)

Identifying disagreement between questionnaire scores and patient-rated change

Disagreement between questionnaire scores and the patient-rated change scale was similar across time-points, so data from baseline to follow-up 1 are presented for brevity. Results for further follow-ups are available in online Supplementary Table S1).

Depression change scores according to patient-rated change

Change in depression questionnaire scores was related to patients' responses on the rating scale. Patients who reported 'feeling a lot better' had the largest mean decrease in scores and patients who reported 'feeling a lot worse' the largest increase (Table 2, first row in PHQ9 and BDI-II sections).

Clinically important change in depression scores according to patient-rated change

When clinically important differences in depression scores were compared to patient ratings, there was evidence of disagreement. The proportion of patients showing each type of clinically important change in questionnaire scores (increase, no change, decrease), in comparison to their responses is presented in Table 2.

Disagreement was most common in patients who reported feeling worse on the patient-rated change scale. PHQ-9 scores showed no change or an improvement for 76% (95% CI 66–83%) of those who reported 'feeling slightly worse', and 81% (95% CI 54–94%) of those who reported 'feeling a lot worse' (Table 2, last row in PHQ-9 section). These results were very similar for the BDI-II (Table 2, last row in BDI-II section). Disagreement was also common in patients who reported feeling better. PHQ-9 scores remained the same or deteriorated in 65% (95% CI 55–74%) of those who reported 'feeling slightly better', and 53% (95% CI 37–67%) of those who reported 'feeling a lot better' (Table 2, last row in PHQ-9 section). Disagreement was lower for patients who reported feeling better on the BDI-II: 43% (95% CI 34–53%) for those reporting feeling slightly better and 28% (95% CI 16–43%) for those reporting feeling much better (Table 2, last row in BDI-II section). Overall, the proportion of people showing some form of disagreement was 51% (95% CI 46–55%) on the PHQ-9 and 55% (95% CI 51–60%) on the BDI-II.

This was similar at follow-up time points for the PHQ-9 [follow-up 1–2: 49% (95% CI 45–54%), follow-up 2–3: 51% (95% CI 46–56%)] and BDI-II [follow-up 1–2: 50%, 95% CI 46–55%, follow-up 2–3: 52% (95% CI 47–57%)]. When using a more stringent minimal clinically important difference of 15%, results were comparable (online Supplementary Table S2).

Quadratic weighted Kappa scores indicated agreement between patient ratings and the categories generated from the change scores ranging 81.2–83.6% for the PHQ-9 and 78.6–83.1% the BDI-II. Unweighted Kappa scores indicated low levels of agreement (3.9–7.6%) for PHQ-9 and BDI-II.

Reliability of disagreement

Results were similar when the proportion of patients scoring within each category of the patient-rated change scale were compared with the relative proportion of patients scoring within these ranges on the PHQ-9 and BDI-II (online Supplementary Table S4). High agreement was observed between the patient-rated change scale and PHQ-9/BDI-II, with weighted kappa values indicating agreement ranging 91.4–93.1% across time-points. Unweighted kappa values indicated poorer agreement (37.9–42.4%). We found no evidence that disagreement differed according to gender (results available on request).

Variables that influence disagreement

Results for the PHQ-9 are shown in Table 3 and for the BDI-II, Table 4. We found evidence that an increase in anxiety symptoms was associated with a decreased odds of reporting feeling better after controlling for changes in depressive symptoms. This was consistent across time-points, for PHQ-9 and BDI-II. For example at follow-up 1, a four-point increase in anxiety scores was associated with a 0.67 (95% CI 0.55–0.82) decrease in the odds of feeling better, having controlled for change in PHQ-9 scores.

We also found consistent evidence that improved mental and physical health-related quality of life was associated with increased odds of reporting feeling better after controlling for changes in depressive symptoms. For example at follow-up 1, an eight-point increase in mental health-related quality of life was associated with a 1.43 (95% CI 1.11–1.61) increase in the odds of feeling better. For physical health-related quality life this odds ratio was 1.28 (1.08–1.54). There was no evidence of an influence of negative life events or social support on the likelihood of reporting improvement (Tables 3 and 4). We found no evidence that any of these associations differed according to gender (available on request).

Discussion

Summary of findings

We found evidence that changes in scores on self-administered depression questionnaires often differ from patients' own views of changes in their mood. Over 50% of people evidenced some form of disagreement between their questionnaire scores and self-rated mood. Even though on average, there is fairly good agreement between change in depressive symptoms and self-rated changes in mood, our results suggest that applying these questionnaires to individual patients will be prone to error.

Patients with more severe anxiety symptoms were less likely, and those with better mental and physical health-related quality of life more likely, to report feeling better having controlled for

Table 2. Change in depression severity according to the patient-rated change scale, compared to clinically important changes in PHQ-9 and BDI-II scores

	Patient-rated change scale				
	Feeling a lot better	Feeling slightly better	Feeling about the same	Feeling slightly worse	Feeling a lot worse
PHQ-9					
Mean (s.d.) change	-3.4 (4.1)	-2.7 (3.9)	-0.26 (3.6)	1.3 (4.3)	1.6 (5.4)
CID decrease, <i>n</i> (%) ^a	19 (47%)	34 (35%)	29 (14%)	9 (9%)	2 (13%)
No CID change, <i>n</i> (%) ^a	20 (50%)	56 (58%)	149 (70%)	65 (66%)	11 (69%)
CID increase, <i>n</i> (%) ^a	1 (3%)	7 (7%)	36 (16%)	24 (25%)	3 (18%)
Disagreement, <i>n</i> (%) ^b	21 (53%)	63 (65%)	65 (30%)	74 (75%)	13 (82%)
BDI-II					
Mean (s.d.) change	-8.0 (8.9)	-5.6 (6.5)	-1.2 (5.8)	0.0 (5.7)	3.2 (7.1)
CID decrease, <i>n</i> (%) ^a	29 (72%)	55 (57%)	74 (34%)	21 (22%)	3 (18%)
No CID change, <i>n</i> (%) ^a	9 (23%)	33 (34%)	92 (42%)	48 (49%)	9 (53%)
CID increase, <i>n</i> (%) ^a	2 (5%)	9 (9%)	51 (24%)	28 (29%)	5 (29%)
Disagreement, <i>n</i> (%) ^b	11 (28%)	42 (43%)	125 (58%)	69 (71%)	12 (71%)

CID, clinically important difference based on the minimal CID (MCID).

Disagreement (differing indications of change in depressive symptoms) is shaded in grey (*n* = 465 PHQ-9, *n* = 468 BDI-II).

^aPercentages represent the proportions of patients showing differing CID changes (decrease, no change, increase) within each category of the global rating of change scale.

^bPercentages represent the proportions of patients showing disagreement within each category of the global rating of change scale.

their depression questionnaire scores. Our results support the idea that self-administered scales only capture a subset of the subjective experience that contributes to patient-rated change and suggests that relying solely upon responses to self-administered scales could be misleading in a large proportion of situations.

Strengths and limitations

We set broad and inclusive entry criteria to reflect the patients consulting for depression in primary care. The MCID allowed us to infer that differences were clinically important, though we acknowledge that the MCID is itself an average determined by reference to patients self-rated change. Our results indicate that such average MCIDs are difficult to apply in individual cases, even if they are valuable overall in planning and interpreting studies.

The depression questionnaires and patient-rated change scale will be subject to measurement error, which could be a potential source of disagreement. Multi-item scales with specific prompts might be more reliable (Kamper et al., 2009), but the reliability of the patient-rated change scale was good. There could be other reasons for disagreement. The patient-rated change scale asks retrospectively about change and recall might be poor (Herrmann, 1995). However, the recall period (2 weeks) was the same for the depression questionnaires and patient-rated change scale. 'Response shift' is the concept that answers will differ across time not because the condition has changed but because the opinion on what the condition means has changed (Schwartz & Sprangers, 1999). This might also lead to disagreement if it occurred. Finally, it is unclear which aspects of the patients' condition have informed the response to the patient-rated change scale. However, these points are largely concerned with explaining the differences between the two contrasting approaches to assessment rather than casting doubt on our conclusions.

There was a low response rate for the study and this might have affected the representativeness of our target population which was

patients seeking help in primary care. However, it seems unlikely that our method of recruitment and the low response rate would inflate the level of disagreement although we cannot rule out that possibility. Our sample was from the UK and predominantly white and this may limit generalizability. Finally, there was attrition though retention was good with 78% at the final follow-up.

These quantitative findings are partly consistent with our previous qualitative findings (Malpass et al., 2016; Robinson et al., 2017). Of course, the PHQ-9 and BDI-II only measure depression symptoms so it is unsurprising that anxiety should affect patient-rated change in mood and feelings independently. Given the co-occurrence of depression and anxiety it is important to recognize that, from the patients' perspective, changes in anxiety will also be important.

The PHQ-9 and BDI-II are recommended for assessment of depressive illness and treatment response in UK primary care and other clinical settings. Our results emphasize the importance of using these measures alongside clinical assessments that take in the perspective of the patient. Sole reliance upon information from self-administered questionnaires can potentially be misleading and ignores areas that patients' regard as important. Our evidence supports the widespread scepticism among physicians about using self-administered questionnaires in clinical practice (Dowrick et al., 2009). We provide quantitative evidence that the results of these questionnaires need to be interpreted along with other clinical assessments and should not be relied upon alone. Our findings support the concept of 'personal recovery', developed in mental health services but also relevant in primary care (Bejerholm & Roe, 2018; Burgess, Pirkis, Coombs, & Rosen, 2011). Personal recovery emphasizes the importance of a holistic focus on patients' broad experiences rather than a restricted focus on 'clinical recovery' or symptom change. This makes the patients' voice of central importance and there are efforts under-way to devise better measurements of patient-reported recovery.

Table 3. Association between exposure variables and the odds of reporting feeling better (*v.* the same or worse), adjusted for change on the PHQ-9

Exposure variable	Odds ratio for reporting feeling better (<i>v.</i> the same or worse), 95% confidence interval and <i>p</i> value (<i>n</i> = 375)		
	Baseline to follow-up 1	Follow-up 1–2	Follow-up 3–4
Anxiety symptoms ^a			
Unadjusted			
Feeling same or worse	ref	ref	ref
Feeling better	0.67 (0.55–0.82) <0.0001	0.65 (0.53–0.79) <0.0001	0.71 (0.59–0.86) <0.0001
Adjusted ^b			
Feeling same or worse	ref	ref	ref
Feeling better	0.66 (0.54–0.82) 0.016	0.61 (0.49–0.76) <0.0001	0.72 (0.60–0.97) 0.001
Mental health-related quality of life ^a			
Unadjusted			
Feeling same or worse	ref	ref	ref
Feeling better	1.34 (1.11–1.61) 0.002	1.33 (1.11–1.59) 0.002	1.38 (1.15–1.64) 0.000
Adjusted ^b			
Feeling same or worse	ref	ref	ref
Feeling better	1.32 (1.08–1.61) 0.006	1.38 (1.14–1.66) 0.001	1.40 (1.17–1.68) <0.000
Physical health-related quality of life ^a			
Unadjusted			
Feeling same or worse	ref	ref	ref
Feeling better	1.28 (1.07–1.54) 0.007	1.25 (1.06–1.48) 0.009	1.20 (1.01–1.42) 0.039
Adjusted ^b			
Feeling same or worse	ref	ref	ref
Feeling better	1.32 (1.08–1.60) 0.006	1.32 (1.10–1.58) 0.002	1.19 (.99–1.43) 0.057
Negative life events ^c			
Unadjusted			
Feeling same or worse	ref	ref	ref
Feeling better	0.98 (0.61–1.59) 0.94	1.13 (0.72–1.79) 0.59	1.17 (0.74–1.85) 0.50
Adjusted ^b			
Feeling same or worse	ref	ref	ref
Feeling better	0.99 (0.60–1.65) 0.98	1.11 (0.69–1.78) 0.76	1.15 (0.72–1.85) 0.56
Social support ^d			
Unadjusted odds Ratio (95% CI) <i>p</i> value			
Feeling same or worse	ref	ref	ref
Feeling better	1.07 (1.00–1.14) 0.067	1.01 (0.95–1.07) 0.71	1.02 (0.96–1.08) 0.56
Adjusted ^b			
Feeling same or worse	ref	ref	ref
Feeling better	1.07 (1.00–1.15) 0.045	1.02 (0.96–1.08) 0.59	1.01 (0.95–1.08) 0.76

^aFor exposures measured at every time-point (anxiety and quality of life), odds ratios represent the odds of reporting feeling better for each four-point increase in anxiety symptoms over time (on a factor score obtained using principal components analysis), adjusted for a binary indicator of meaningful change on the PHQ9.

^bAdjusted for age, sex, ethnicity, site, education level, current use of antidepressants and marital, financial and employment status.

^cNegative life events were measured at baseline only. The odds ratio represents the odds of feeling better in those who reported one life event or more compared to those who reported no life events, adjusted for a binary indicator of meaningful change on the PHQ9.

^dSocial support was measured at baseline only. The odds ratio represents the odds of reporting feeling better for each standard deviation increase in social support, adjusted for a binary indicator of meaningful change on the PHQ9.

Some patients view self-administered questionnaires positively and request them to monitor their recovery (Moore *et al.*, 2012). Questionnaires can, therefore, play a useful role in outcome assessment, in conjunction with the clinical assessment that takes account of more holistic changes in mood. They are also useful as a guide for service level outcome assessment (Clark *et al.*, 2018). In clinical trials, self-administered questionnaires are widely used for comparing groups and such randomized

comparisons should be unbiased. Our findings suggest, though, that additional questions should also be used to assess the outcome of treatments in research studies.

Future research could examine the generalizability of our findings to international settings and mental health services, and the relationship between patient-rated change and other mental health measures including the outcomes used in the NHS improving access to psychological therapy services (Clark *et al.*, 2018).

Table 4. Association between exposure variables and the odds of reporting feeling better (v. the same or worse), adjusted for change on the BDI-II

Exposure variable	Odds ratio for reporting feeling better (v. the same or worse), 95% confidence interval and <i>p</i> value (<i>n</i> = 375)		
	Baseline to follow-up 1	Follow-up 1–2	Follow-up 3–4
Anxiety symptoms ^a			
Unadjusted			
Feeling same or worse	ref	ref	ref
Feeling better	0.67 (0.56–0.81) <0.0001	0.67 (0.56–0.81) <0.0001	0.70 (0.59–0.84) <0.0001
Adjusted ^b			
Feeling same or worse	ref	ref	ref
Feeling better	0.65 (0.53–0.81) <0.0001	0.61 (0.49–0.76) <0.0001	0.71 (0.59–0.86) <0.0001
Mental health-related quality of life ^a			
Unadjusted			
Feeling same or worse	ref	ref	ref
Feeling better	1.37 (1.13–1.65) 0.001	1.33 (1.12–1.58) 0.001	1.38 (1.16–1.64) <0.0001
Adjusted ^c			
Feeling same or worse	ref	ref	ref
Feeling better	1.34 (1.10–1.63) 0.004	1.38 (1.14–1.66) 0.001	1.38 (1.16–1.64) <0.0001
Physical health-related quality of life ^a			
Unadjusted			
Feeling same or worse	ref	ref	ref
Feeling better	1.25 (1.04–1.49) 0.016	1.24 (1.05–1.46) 0.013	1.22 (1.03–1.45) 0.021
Adjusted ^c			
Feeling same or worse	ref	ref	ref
Feeling better	1.27 (1.05–1.54) 0.015	1.30 (1.08–1.55) 0.005	1.22 (1.02–1.47) 0.030
Negative life events ^d			
Unadjusted			
Feeling same or worse	ref	ref	ref
Feeling better	1.03 (0.64–1.66) 0.89	1.18 (0.75–1.85) 0.49	1.14 (0.71–1.81) 0.59
Adjusted ^c			
Feeling same or worse	ref	ref	ref
Feeling better	1.04 (0.63–1.72) 0.87	1.15 (0.72–1.85) 0.55	1.11 (0.68–1.79) 0.68
Social support ^b			
Unadjusted			
Feeling same or worse	ref	ref	ref
Feeling better	1.07 (1–1.14) 0.06	1.01 (0.95–1.07) 0.71	1.02 (0.96–1.09) 0.52
Adjusted ^c			
Feeling same or worse	ref	ref	ref
Feeling better	1.07 (1.00–1.15) 0.044	1.02 (0.96–1.08) 0.59	1.01 (0.95–1.08) 0.70

^aFor exposures measured at every time-point (anxiety and quality of life), odds ratios represent the odds of reporting feeling better for each four-point increase in anxiety symptoms over time (on a factor score obtained using principal components analysis), adjusted for a binary indicator of meaningful change on the PHQ9.

^bSocial support was measured at baseline only. The odds ratio represents the odds of reporting feeling better for each standard deviation increase in social support, adjusted for a binary indicator of meaningful change on the PHQ9.

^cAdjusted for age, sex, ethnicity, site, education level, current use of antidepressants and marital, financial and employment status.

^dNegative life events were measured at baseline only. The odds ratio represents the odds of feeling better in those who reported one life event or more compared to those who reported no life events, adjusted for a binary indicator of meaningful change on the PHQ9.

Future clinical trials could also use the patient-rated change in mood question as an outcome that might help to address the limitations of existing measures.

Supplementary material. The supplementary material for this article can be found at <https://doi.org/10.1017/S0033291719003878>.

Acknowledgements. We would like to thank the participants who took part in the study and the general practice surgeries for their help with recruitment.

We also thank the research team: researchers, research nurses and administrative staff.

Author contributions. Glyn Lewis as chief investigator, along with the other co-applicants of the PANDA programme, conceived and designed the study. Catherine Hobbs and Gemma Lewis analyzed the data and drafted the manuscript. All authors interpreted the data and criticized the manuscript for important intellectual content. All authors have read and approved the final version of the manuscript. This article is the work of the authors. All authors,

external and internal, had full access to all of the data (including statistical reports and tables) in the study and can take responsibility for the integrity of the data and the accuracy of the data analysis. We would like to thank Professor Tony Ades for his advice on statistical analyses.

Financial support. The study is funded by a National Institute of Health Research (NIHR). The PANDA is independent research commissioned by the National Institute for Health Research Programme Grant for Applied Research (RP-PG-0610-10048). The views expressed in this publication are those of the author(s) and not necessarily those of the sponsor, NHS, the National Institute for Health Research or the Department of Health and Social Care. The funder had no role in the study design, data collection, data analysis, interpretation of data or writing of the report. This work was partially supported by the UCLH NIHR Biomedical Research Centre.

Conflict of interest. None.

Ethical standards. The authors assert that all procedures contributing to this work comply with the ethical standards of the relevant national and institutional committees on human experimentation and with the Helsinki Declaration of 1975, as revised in 2008.

References

- Beck, A. T., Guth, D., Steer, R. A., & Ball, R. (1997). Screening for major depression disorders in medical inpatients with the Beck Depression Inventory for Primary Care. *Behaviour Research and Therapy*, 35, 785–791.
- Bejerholm, U., & Roe, D. (2018). Personal recovery within positive psychiatry. *Nordic Journal of Psychiatry*, 72, 420–430.
- Berk, M., Ng, F., Dodd, S., Callaly, T., Campbell, S., Bernardo, M., & Trauer, T. (2008). The validity of the CGI severity and improvement scales as measures of clinical effectiveness suitable for routine clinical use. *Journal of Evaluation in Clinical Practice*, 14, 979–983.
- Burgess, P., Pirkis, J., Coombs, T., & Rosen, A. (2011). Assessing the value of existing recovery measures for routine use in Australian mental health services. *Australian & New Zealand Journal of Psychiatry*, 45, 267–280.
- Button, K. S., Kounali, D., Thomas, L., Wiles, N. J., Peters, T. J., Welton, N. J., ... Lewis, G. (2015). Minimal clinically important difference on the Beck Depression Inventory-II according to the patient's perspective. *Psychological Medicine*, 45, 3269–3279.
- Clark, D. M., Canvin, L., Green, J., Layard, R., Pilling, S., & Janecka, M. (2018). Transparency about the outcomes of mental health services (IAPT approach): An analysis of public data. *Lancet (London, England)*, 391, 679–686.
- Dowrick, C., Leydon, G. M., McBride, A., Howe, A., Burgess, H., Clarke, P., ... Kendrick, T. (2009). Patients' and doctors' views on depression severity questionnaires incentivised in UK quality and outcomes framework: Qualitative study. *BMJ (Clinical research ed.)*, 338, b663.
- Fischer, D., Stewart, A. L., Bloch, D. A., Lorig, K., Laurent, D., & Holman, H. (1999). Capturing the patient's view of change as a clinical outcome measure. *JAMA*, 282, 1157.
- Herrmann, D. (1995). Reporting current, past, and changed health status. What we know about distortion. *Medical Care*, 33, AS89–AS94.
- Jolliffe, I. T., & Cadima, J. (2016). Principal component analysis: A review and recent developments. *Philosophical Transactions. Series A, Mathematical, Physical, and Engineering Sciences*, 374, 20150202.
- Kamper, S. J., Maher, C. G., & Mackay, G. (2009). Global rating of change scales: A review of strengths and weaknesses and considerations for design. *The Journal of Manual & Manipulative Therapy*, 17, 163–170.
- Kendrick, T., Dowrick, C., McBride, A., Howe, A., Clarke, P., Maisey, S., ... Smith, P. W. (2009). Management of depression in UK general practice in relation to scores on depression severity questionnaires: Analysis of medical record data. *BMJ (Clinical research ed.)*, 338, b750.
- Kroenke, K., Spitzer, R. L., & Williams, J. B. (2001). The PHQ-9: Validity of a brief depression severity measure. *Journal of General Internal Medicine*, 16, 606–613.
- Malpass, A., Dowrick, C., Gilbody, S., Robinson, J., Wiles, N., Duffy, L., & Lewis, G. (2016). Usefulness of PHQ-9 in primary care to determine meaningful symptoms of low mood: A qualitative study. *The British Journal of General Practice: The Journal of the Royal College of General Practitioners*, 66, e78–e84.
- Moore, M., Ali, S., Stuart, B., Leydon, G. M., Ovens, J., Goodall, C., & Kendrick, T. (2012). Depression management in primary care: An observational study of management changes related to PHQ-9 score for depression monitoring. *British Journal of General Practice*, 62, e451–e457.
- Moriarty, A. S., Gilbody, S., McMillan, D., & Manea, L. (2015). Screening and case finding for major depressive disorder using the Patient Health Questionnaire (PHQ-9): A meta-analysis. *General Hospital Psychiatry*, 37, 567–576.
- Robinson, J., Khan, N., Fusco, L., Malpass, A., Lewis, G., & Dowrick, C. (2017). Why are there discrepancies between depressed patients' Global Rating of Change and scores on the Patient Health Questionnaire depression module? A qualitative study of primary care in England. *BMJ Open*, 7, e014519.
- Schwartz, C. E., & Sprangers, M. A. (1999). Methodological approaches for assessing response shift in longitudinal health-related quality-of-life research. *Social Science & Medicine* (1982), 48, 1531–1548.
- Spitzer, R. L., Kroenke, K., Williams, J. B. W., & Löwe, B. (2006). A brief measure for assessing generalized anxiety disorder. *Archives of Internal Medicine*, 166, 1092.
- Thombs, B. D., & Ziegelstein, R. C. (2014). Does depression screening improve depression outcomes in primary care? *BMJ (Clinical research ed.)*, 348, g1253.
- Toop, L. (2011). The QOF, NICE, and depression: A clumsy mechanism that undermines clinical judgment. *The British Journal of General Practice: the Journal of the Royal College of General Practitioners*, 61, 432–433.
- Ware, J., Kosinski, M., & Keller, S. D. (1996). A 12-Item Short-Form Health Survey: Construction of scales and preliminary tests of reliability and validity. *Medical Care*, 34, 220–233.
- Zaider, T. I., Heimberg, R. G., Fresco, D. M., Schneier, F. R., & Liebowitz, M. R. (2003). Evaluation of the clinical global impression scale among individuals with social anxiety disorder. *Psychological Medicine*, 33, 611–622.