

**Adaptation of the PHQ-4 Short Screening for Depression and Anxiety to
increase its Sensitivity to Subclinical Variability**

Dominique Makowski¹, An Shu Te¹, & S.H. Annabel Chen^{1, 2, 3, 4}

¹ School of Social Sciences, Nanyang Technological University, Singapore

² LKC Medicine, Nanyang Technological University, Singapore

³ National Institute of Education, Singapore

⁴ Centre for Research and Development in Learning, Nanyang Technological University,
Singapore

The authors made the following contributions. Dominique Makowski:
Conceptualization, Data curation, Formal Analysis, Funding acquisition, Investigation,
Methodology, Project administration, Resources, Software, Supervision, Validation,
Visualization, Writing – original draft; An Shu Te: Project administration, Resources,
Investigation, Writing – original draft; S.H. Annabel Chen: Project administration,
Supervision, Writing – review & editing.

Correspondence concerning this article should be addressed to Dominique Makowski,
HSS 04-18, 48 Nanyang Avenue, Singapore. E-mail: dom.makowski@gmail.com

Abstract

Something something

Significance Statement. Developing reliable and sensitive instruments for mood disorders screening is critical in a global context marked by international crises (pandemics, wars), where more and more surveys are done online. In this study, we show that a small modification to the widely used PHQ-4 scale (adding the “Once or twice” response option) can increase its ability to capture the mild fluctuations prevalent in subclinical samples.

Keywords: visual illusions, illusion game, Pyllusion, personality, general factor

Word count: 4085

Adaptation of the PHQ-4 Short Screening for Depression and Anxiety to increase its Sensitivity to Subclinical Variability

Introduction

The Patient Health Questionnaire-4 (PHQ-4) is an ultra brief measurement of core signs of depression and anxiety (Kroenke et al., 2009). It consists of two items for depression (PHQ-2) and anxiety (GAD-2), each corresponding to DSM-IV Diagnostic Criterion A symptoms for major depressive disorder (MDD) and generalized anxiety disorder (GAD). It has been validated across many languages and samples Materu et al. (2020), becoming one of the most popular screening instrument (Maurer et al., 2018).

While the scale has been validated and used in the general population Hajek & König (2020), its initial purpose was to reliably discriminate and identify potential MDD/GAD patients. This diagnostic and discriminative goal materializes in the scale's design and the existence of categorical cut-offs. However, it might not be best suited to capture subclinical variability. In particular, the gap between the two lowest possible answers, "Not at all" and "Several days", is quite large and leaves out the possibility of more subtle occurrences. While this is not necessarily an issue in clinical and diagnostic contexts, it might lead to a sub-optimal discrimination on the lower end of the spectrum, important for instance in the context of variability quantifications. The goal of this study is to increase, with minimal changes to the scale, the sensitivity to very mild mood alterations.

Methods

Original Scale

The instructions "Over the last 2 weeks, how often have you been bothered by the following problems?" are followed with 4 items (A1 - Feeling nervous, anxious or on edge; A2 - Not being able to stop or control worrying; D1 - Little interest or pleasure in doing things; D2 - Feeling down, depressed, or hopeless). The original answer options are "Not at

all” (0), “Several days” (1), “More than half the days” (2), “Nearly every day” (3). The total score is computed by summing the responses of each facet.

Revision

In order to better capture potential mild mood inflections without altering the scale scoring or structure, we added a “Once or twice” option between “Not at all” and “Several days” (see Dobson & Mothersill, 1979 for the choice of the label).

Participants

The sample consists of 500 English-speaking participants (**stats**) who were administered the refined PHQ-4 online as part of another study, which data is available in open-access at <https://github.com/RealityBending/IllusionGameReliability>.

This study was approved by the NTU Institutional Review Board (NTU **NUMBER**). All participants provided their informed consent prior to participation and were incentivized after completing the study.

Results

The fully reproducible analysis script is accessible at <https://github.com/DominiqueMakowski/PHQ4R>.

Descriptive

The reliability of the anxiety (*Cronbach's* $\alpha = 0.892$; RMSEA = 0.032) and depression (*Cronbach's* $\alpha = 0.822$; RMSEA = 0.032) subscales was excellent. The proportion of response types stratified by item (see **Figure 1A**) shows that the new “Once or twice” option was the most prevalent response for almost all items (on average selected in **30.07%** of cases).

Item Response Theory

Item Response Theory (IRT) provides insights into how well items and responses capture the underlying latent trait θ . For each of the subscales, we fitted a unidimensional graded response model (GRM, Samejima, 1997). For anxiety, the two items captured 88.1% of the variance of the latent anxiety dimension ($\theta_{anxiety}$). The discrimination parameters suggested that the first item was less precise ($\alpha = 3.20$) than the second item ($\alpha = 12.53$) in its ability to discriminate between various levels of anxiety (i.e., each response on item 1 covers a larger range of $\theta_{anxiety}$, as can be seen in **Figure 1B**). The two depression items captured 80.6% of the variance of its latent trait ($\theta_{depression}$), and the opposite pattern was found: the first item had a higher precision ($\alpha = 15.69$) than the first ($\alpha = 2.19$). Importantly to our objective, it seems that the added “Once or twice” option does cover a selective portion of the latent space.

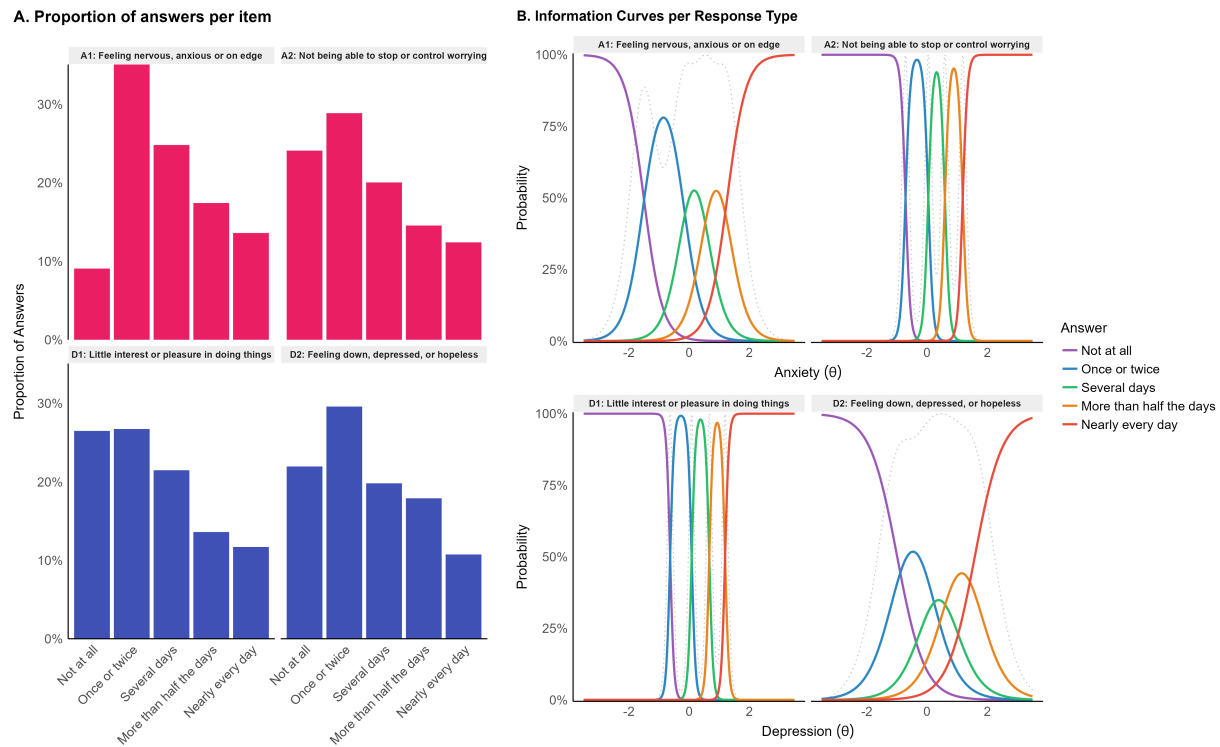


Figure 1. A. Proportion of answers of each type. B. Item Information Curves from IRT showing the coverage by each item and response of the latent dimension. Typically, an optimally informative item would display a large coverage over theta, with each response presenting a narrow coverage (high discrimination between different levels).

Scoring

We propose two types of scoring procedures. The first aims at minimally disrupting the original scale and making its scores comparable, enabling comparisons across studies and the usage of developed cut-offs and norms. We suggest to score the new option, “Once or twice”, as 0.5, occupying the middle space between 0 and 1 (preserving this way the total score range of 0 - 8 per dimension).

The second scoring method we propose takes into account the non-linear spacing between responses, as indicated by the peak of their measurement information. To develop it, we found for each dimension the lower and higher value of the latent trait θ for which the probability lower and higher probability of identification is 95%. We considered these points at 0 and 1, and normalized the location of the peaks of each response in between.

To illustrate, a person answering “Several days” and “Once or twice” to the first and second item of the anxiety subscale would have a score of 1.5 ($1 + 0.5$ following the basic scoring) and a normalized score of 0.50 ($((0.56 + 0.45)/2)$).

Discussion

The objective of this study was to test the introduction of a “Once or twice” response option to the PHQ-4 to make it more sensitive to milder fluctuations. The fact that this new response option was the most prevalent is in itself evidence for its usefulness, and IRT analysis further revealed how this response captures with precision a unique portion of anxiety and depression. Our results suggest that adding this option response increases the scale’s potential to discriminate average mood levels (which are superior to zero) from lower-end extremes (the true zero).

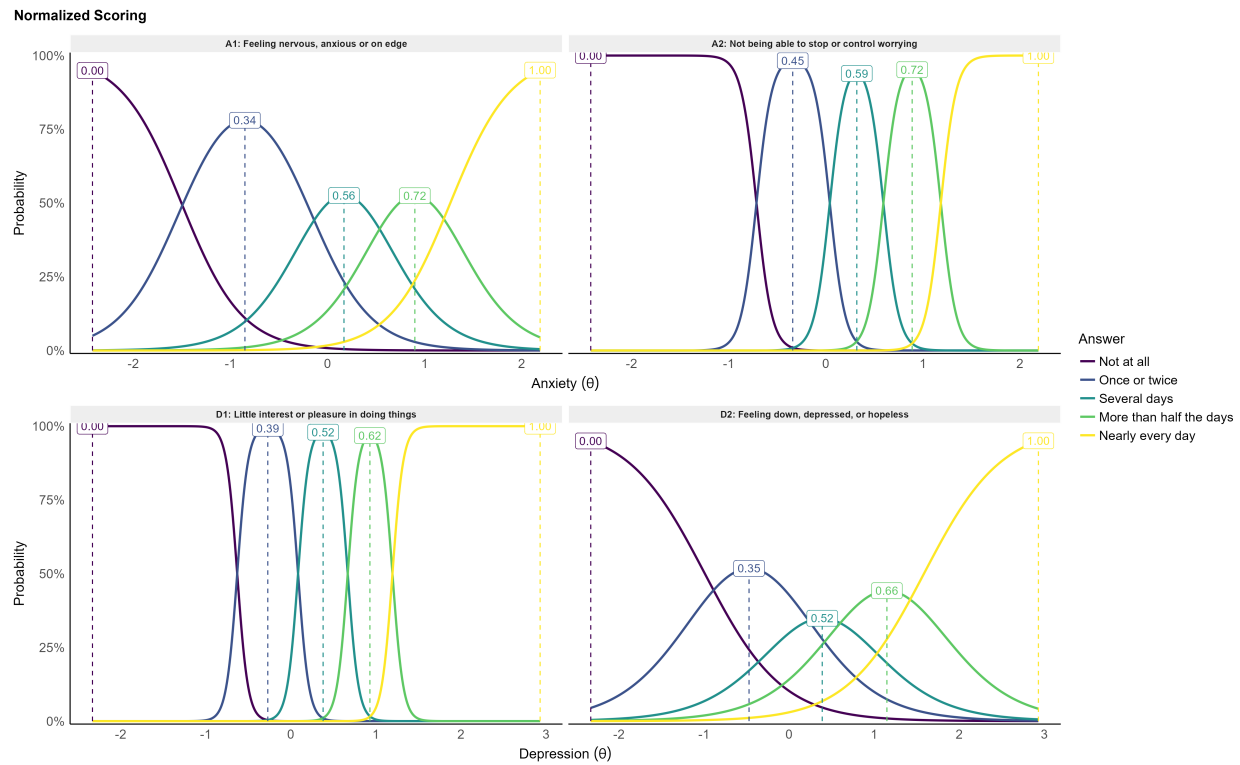


Figure 2. Normalized scoring of the items taking into account the non-linear spacing between responses (based on the peak of their measurement information).

Data Availability

The dataset analysed during the current study are available in the GitHub repository
<https://github.com/DominiqueMakowski/PHQ4R>

Funding

This work was supported by the Presidential Postdoctoral Fellowship Grant
 (NTU-PPF-2020-10014) from Nanyang Technological University (awarded to DM).

Table 1

Refined Patient Health Questionnaire-4 (PHQ-4R). The instructions are "Over the last 2 weeks, how often have you been bothered by the following problems?". The 0.5 scoring is introduced to preserve the compatibility with the original version, and enables to use its norms and cut-offs. The normalized scoring takes into account the non-linear spacing between item responses.

Facet	Item	Response	Basic Scoring	Normalized Scoring
Anxiety	Feeling nervous, anxious or on edge	<i>Not at all</i>	0	0
		<i>Once or twice</i>	0.5	0.5
		<i>Several days</i>	1	0.7
		<i>More than half the days</i>	2	0.8
		<i>Nearly every day</i>	3	1
	Not being able to stop or control worrying	<i>Not at all</i>	0	0
		<i>Once or twice</i>	0.5	0.5
		<i>Several days</i>	1	0.7
		<i>More than half the days</i>	2	0.8
		<i>Nearly every day</i>	3	1
Depression	Little interest or pleasure in doing things	<i>Not at all</i>	0	0
		<i>Once or twice</i>	0.5	0.5
		<i>Several days</i>	1	0.7
		<i>More than half the days</i>	2	0.8
		<i>Nearly every day</i>	3	1
	Feeling down, depressed, or hopeless	<i>Not at all</i>	0	0
		<i>Once or twice</i>	0.5	0.5
		<i>Several days</i>	1	0.7
		<i>More than half the days</i>	2	0.8
		<i>Nearly every day</i>	3	1

References

- Christodoulaki, A., Baralou, V., Konstantakopoulos, G., & Touloumi, G. (2022). Validation of the patient health questionnaire-4 (PHQ-4) to screen for depression and anxiety in the greek general population. *Journal of Psychosomatic Research*, 160, 110970.
- Dobson, K. S., & Mothersill, K. J. (1979). Equidistant categorical labels for construction of likert-type scales. *Perceptual and Motor Skills*, 49(2), 575–580.
- Hajek, A., & König, H.-H. (2020). Prevalence and correlates of individuals screening positive for depression and anxiety on the phq-4 in the german general population:

Findings from the nationally representative german socio-economic panel (GSOEP).

International Journal of Environmental Research and Public Health, 17(21), 7865.

Kroenke, K., Spitzer, R. L., Williams, J. B., & Löwe, B. (2009). An ultra-brief screening scale for anxiety and depression: The PHQ-4. *Psychosomatics*, 50(6), 613–621.

Löwe, B., Wahl, I., Rose, M., Spitzer, C., Glaesmer, H., Wingenfeld, K., Schneider, A., & Brähler, E. (2010). A 4-item measure of depression and anxiety: Validation and standardization of the patient health questionnaire-4 (PHQ-4) in the general population. *Journal of Affective Disorders*, 122(1-2), 86–95.

Materu, J., Kuringe, E., Nyato, D., Galishi, A., Mwanamsangu, A., Katebalila, M., Shao, A., Chungalucha, J., Nnko, S., & Wambura, M. (2020). The psychometric properties of PHQ-4 anxiety and depression screening scale among out of school adolescent girls and young women in tanzania: A cross-sectional study. *BMC Psychiatry*, 20(1), 1–8.

Maurer, D. M., Raymond, T. J., & Davis, B. N. (2018). Depression: Screening and diagnosis. *American Family Physician*, 98(8), 508–515.

Mendoza, N. B., Frondozo, C. E., Dizon, J. I. W. T., & Buenconsejo, J. U. (2022). The factor structure and measurement invariance of the PHQ-4 and the prevalence of depression and anxiety in a southeast asian context amid the COVID-19 pandemic. *Current Psychology*, 1–10.

Samejima, F. (1997). Graded response model. In *Handbook of modern item response theory* (pp. 85–100). Springer.