



HHS Public Access

Author manuscript

Psychol Med. Author manuscript; available in PMC 2021 June 01.

Published in final edited form as:

Psychol Med. 2020 June ; 50(8): 1368–1380. doi:10.1017/S0033291719001314.

Corresponding authors: Andrea Benedetti, PhD; Centre for Outcomes Research & Evaluation, Research Institute of the McGill University Health Centre, 5252 Boulevard de Maisonneuve, Montréal, QC, H4A 3S5, Canada; Tel (514) 934-1934 ext. 32161; andrea.benedetti@mcgill.ca, Brett D. Thombs, PhD; Jewish General Hospital; 4333 Cote Ste Catherine Road; Montreal, Quebec H3T 1E4; Tel (514) 340-8222 ext. 25112; brett.thombs@mcgill.ca.

Full Addresses:

Yin Wu

Lady Davis Institute for Medical Research, Jewish General Hospital, 4333 Chemin de la Côte-Sainte-Catherine, Montréal, QC, H3T 1E4, Canada

Brooke Levis

Lady Davis Institute for Medical Research, Jewish General Hospital, 4333 Chemin de la Côte-Sainte-Catherine, Montréal, QC, H3T 1E4, Canada

Kira E. Riehm

Lady Davis Institute for Medical Research, Jewish General Hospital, 4333 Chemin de la Côte-Sainte-Catherine, Montréal, QC, H3T 1E4, Canada

Nazanin Saadat

Lady Davis Institute for Medical Research, Jewish General Hospital, 4333 Chemin de la Côte-Sainte-Catherine, Montréal, QC, H3T 1E4, Canada

Alexander W. Levis

Lady Davis Institute for Medical Research, Jewish General Hospital, 4333 Chemin de la Côte-Sainte-Catherine, Montréal, QC, H3T 1E4, Canada

Marleine Azar

Lady Davis Institute for Medical Research, Jewish General Hospital, 4333 Chemin de la Côte-Sainte-Catherine, Montréal, QC, H3T 1E4, Canada

Danielle B. Rice

Lady Davis Institute for Medical Research, Jewish General Hospital, 4333 Chemin de la Côte-Sainte-Catherine, Montréal, QC, H3T 1E4, Canada

Jill Boruff

Macdonald-Stewart Library Building, 809 Sherbrooke Street West, Montreal, Quebec H3A 0C1, Canada

Pim Cuijpers

Department of Clinical, Neuro and Developmental Psychology, Faculty of Behavioural and Movement Sciences, Vrije Universiteit Amsterdam, Van der Boechorststraat 1, 1081 BT Amsterdam, The Netherlands

Simon Gilbody

Mental Health and Addiction Research Group, Department of Health Sciences and Hull York Medical School, University of York, Heslington YO10 5DD, United Kingdom

John P.A. Ioannidis

Stanford University, 1265 Welch Road, MSOB X306, Stanford, CA, 94305, USA

Lorie A. Kloda

Concordia University, 1455, boul. de Maisonneuve Ouest, LB-331, Montréal, QC, H3G 1M8, Canada

Dean McMillan

Mental Health and Addiction Research Group, Department of Health Sciences and Hull York Medical School, University of York, Heslington YO10 5DD, United Kingdom

Scott B. Patten

Department of Community Health Sciences, 3rd Floor, TRW Building, University of Calgary, 3280 Hospital Drive NW, Calgary, AB, T2N 4Z6, Canada

Ian Shrier

Centre for Clinical Epidemiology, Lady Davis Institute for Medical Research, Jewish General Hospital, 3755 Cote Ste-Catherine Rd, Montréal, QC, H3T 1E2, Canada

Roy C. Ziegelstein

Johns Hopkins University School of Medicine, Miller Research Building, 733 N. Broadway, Suite 115, Baltimore, MD, 21205, USA

Dickens H. Akena

Department of Psychiatry, Makerere University College of Health Sciences, P.O.Box 7062 Kampala, Uganda

Bruce Arroll

Department of General Practice and Primary Health Care, University of Auckland, Private Bag 92019, Auckland 1142, New Zealand

Liat Ayalon

Louis and Gabi Weisfeld School of Social Work, Ramat Gan, Bar Ilan University, 52900, Israel

Hamid R. Baradaran

Endocrine Research Center, Institute of Endocrinology and Metabolism, Iran University of Medical Sciences, Tehran 15937-16615, Iran

Murray Baron

Jewish General Hospital, Suite A 725, 3755 Cote St Catherine Rd, Montréal, QC, H3T 1E2, Canada

Charles H. Bombardier

Division of Clinical and Neuropsychology, Department of Rehabilitation Medicine, University of Washington, Box 359612, Harborview Medical Center, 325 9th Avenue, Seattle, WA, 98104, USA

Peter Butterworth

Research School of Population Health, Florey Building (54), Australian National University, Canberra, ACT, 2600, Australia

Gregory Carter

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Locked Bag #7, Hunter Region Mail Centre, NSW 2310, Australia
 Locked Mail Bag 9718, Wellington, New Zealand

Equivalency of the diagnostic accuracy of the PHQ-8 and PHQ-9:

University of São Paulo, Av. Bandeirantes, 3900, 14048-900-Ribeirão Preto, SP, Brazil
 Juliana C. N. Chan
 Department of Medicine and Therapeutics, The Chinese University of Hong Kong, 9/F Lui Che Woo Clinical Sciences Building,
 Prince of Wales Hospital, Shatin, Hong Kong
 Rushina Cholera
 UNC Department of Pediatric, 260 MacNider Building, CB# 7220, 321 S. Columbia Street, UNC School of Medicine, Chapel Hill,
 NC 27599-7220, USA
 Yeates Conwell
 University of Rochester Medical Center, 300 Crittenden Blvd., Rochester, NY, 14642, USA
 Janneke M. de Man-van Ginkel
 University Medical Center Utrecht, Internal mail no Str. 6.131, P.O. Box 85500, 3508 GA, Utrecht, The Netherlands
 Jesse R. Fann
 Department of Psychiatry & Behavioral Sciences, University of Washington, Box 356560, Seattle, WA 98195
 Felix H. Fischer
 Medizinische Klinik mit Schwerpunkt Psychosomatik, Charité - Universitätsmedizin Berlin, Charitéplatz 1, 10098 Berlin, Germany
 Daniel Fung
 Institute of Mental Health, 10 Buangkok View, 539747, Singapore
 Bizu Gelaye
 Department of Epidemiology, 677 Huntington Ave, Room 505F, Boston, MA, 02115, USA
 Felicity Goodyear-Smith
 Department of General Practice and Primary Health Care, University of Auckland, PB 92019, Auckland, 1142, New Zealand
 Catherine G. Greeno
 2204 Cathedral of Learning, University of Pittsburgh, 4200 Fifth Ave, Pittsburgh, PA, 15260, USA
 Brian J. Hall
 Department of Psychology, Faculty of Social Sciences, Humanities and Social Sciences Building E21-3040, University of Macau, E21
 Avenida da Universidade, Taipa, Macau, China
 Patricia A. Harrison
 City of Minneapolis Health Department, 250 S. Fourth St., Room 510, Minneapolis, MN 55415, USA
 Martin Härter
 University Medical Center Hamburg-Eppendorf, Department of Medical Psychology, Martinistraße 52, 20246 Hamburg, Germany
 Ulrich Hegerl
 University of Leipzig, Department of Psychiatry and Psychotherapy, Semmelweisstrasse 10, 04103 Leipzig, Germany
 Leanne Hides
 School of Psychology, University of Queensland, St Lucia, Brisbane, Queensland, 4072, Australia
 Stevan E. Hobfoll
 STAR-Stress, Anxiety & Resilience Consultants, 1000 N Lake Shore Plaza 9B, Chicago, IL, 60611
 Marie Hudson
 Jewish General Hospital and Lady Davis Research Institute, 3755 Côte Ste-Catherine Rd, Room A725, Montréal, QC, H3T 1E2,
 Canada
 Thomas Hyphantis
 Department of Psychiatry, Faculty of Medicine, School of Health Sciences, University of Ioannina, Ioannina 451 10, Greece
 Masatoshi Inagaki
 Department of Psychiatry, Faculty of Medicine, Shimane University, 89-1 Enya-cho, Izumo, Shimane, Japan
 Nathalie Jetté
 Department of Neurology, Icahn School of Medicine at Mount Sinai, 1 Gustave L. Levy Pl, New York, NY 10029, USA
 Mohammad E. Khamseh
 Endocrine Research Center, Institute of Endocrinology and Metabolism, Iran University of Medical Sciences, Tehran 15937-16615,
 Iran
 Kim M. Kiely
 Neuroscience Research Australia, Margarete Ainsworth Building, Barker Street, NSW 2031, Sydney, Australia
 Yunxin Kwan
 Tan Tock Seng Hospital, 11 Jalan Tan Tock Seng, 308433, Singapore
 Femke Lamers
 Amsterdam UMC, Vrije Universiteit, Department Psychiatry, Oldenaller 1, 1081 HJ Amsterdam, The Netherlands
 Shen-Ing Liu
 Department of Psychiatry, Mackay Memorial Hospital, No. 92, Section 2, Chung-Shan North Rd, Taipei, Taiwan
 Manote Lotrakul
 Department of Psychiatry, Faculty of Medicine, Ramathibodi Hospital, Mahidol University, Bangkok 10400, Thailand
 Sonia R. Loureiro
 Pós-graduação em Saúde Mental, Depto. de Neurociências e Ciências do Comportamento da Faculdade de Medicina de Ribeirão
 Preto, USP, Rua Tenente Catão Roxo, 2650, CEP 14051-140, Ribeirão Preto, SP, Brazil
 Bernd Löwe
 Universitätsklinikum Hamburg-Eppendorf, Institut und Poliklinik für Psychosomatische Medizin und Psychotherapie, Martinistr. 52,
 Gebäude O25, 20246 Hamburg, Germany
 Anthony McGuire
 Department of Nursing, St. Joseph's College, 278 Whites Bridge Rd., Standish, ME, 04084, USA

A systematic review and individual participant data meta-

Sherma Mond-Siddik

Cancer Resource & Education Centre/ Department of Psychiatry, Faculty of Medicine & Health Sciences, Universiti Putra Malaysia, 43400 UPM Serdang, Selangor, Malaysia

Tiago N. Munhoz

Depto Medicina Social, Programa Pós-graduação Epidemiologia, Universidade Federal de Pelotas, Rua Marechal Deodoro 1160, 3º piso, 96020-220 - Pelotas, RS, Brasil

Kumiko Muramatsu

Department of Clinical Psychology, Graduate School of Niigata Seiryō University, 1-5939, Suidocho, Chuo-ku, Niigata 951-8121, Japan.

Flávia L. Osório

Department of Neurosciences and Behaviour, Medical School of Ribeirão Preto, University of São Paulo, Avenida dos Bandeirantes 3900, CEP 14048-900, Brazil

Vikram Patel

Department of Global Health and Social Medicine, Harvard Medical School, Boston, USA 02119, USA

Brian W. Pence

Department of Epidemiology, UNC-Chapel Hill, McGavran-Greenberg 2103C, CB#7435, 135 Dauer Dr, Chapel Hill NC 27599-7435, USA

Philippe Persoons

Katholieke Universiteit Leuven, Department of Neurosciences, Research Group Psychiatry, University Psychiatric Center KU Leuven, Herestraat 49, 3000 Leuven, Belgium

Angelo Picardi

Italian National Institute of Health, Centre for Behavioural Sciences and Mental Health, Viale Regina Elena 299, 00161 Rome, Italy

Katrin Reuter

Outpatient Service for Psychotherapy and Psychooncology, Stadtstrasse 11, 79104 Freiburg, Germany

Alasdair G. Rooney

Division of Psychiatry, University of Edinburgh, Royal Edinburgh Hospital

Edinburgh, EH10 5HF, Scotland

Iná S. Santos

Depto Medicina Social, Programa Pós-graduação Epidemiologia, Universidade Federal de Pelotas, Rua Marechal Deodoro 1160, 3º piso

96020-220 - Pelotas, RS, Brasil

Juwita Shaaban

School of Medical Science, Health Campus Universiti Sains Malaysia, 16150 Kubang Kerian, Kelantan, Malaysia

Abbey Sidebottom

Allina Health, 800 E 28th Street, MR 15521, Minneapolis, MN 55407-3799, USA

Adam Simning

University of Rochester Medical Center, School of Medicine and Dentistry, 601 Elmwood Ave, Box PSYCH, Rochester, NY 14642, USA

Lesley Stafford

Centre for Women's Mental Health, The Royal Women's Hospital, Locked Bag 300, Parkville Victoria 3052, Australia

Sharon Sung

Programme in Health Services & Systems Research, Duke-NUS Medical School, 20 College Road, Level 6, 169856, Singapore

Pei Lin Lynnette Tan

Tan Tock Seng Hospital, 11 Jalan Tan Tock Seng, 308433, Singapore

Alyna Turner

IMPACT SRC, Deakin University, HERB Building Level 3, PO Box 281 Geelong 3220, Australia

Henk C. van Weert

Department of General Practice, Amsterdam Institute for General Practice and Public Health, Amsterdam University Medical Centers, location AMC, Meibergdreef 9, 1105 AZ Amsterdam, The Netherlands

Jennifer White

Department of Physiotherapy, School of Primary and Allied Health Care, Monash University, Building G, Level 3, McMahons Road, Frankston Victoria 3199

Australia

Mary A. Whooley

Department of Veterans Affairs Medical Center, 4150 Clement Street (111A1), San Francisco, CA 94121, USA

Kirsty Winkley

Faculty of Nursing, Midwifery & Palliative Care, King's College London, Strand, London WC2R 2LS, UK

Mitsuhiko Yamada

National Center of Neurology and Psychiatry, 4-1-1 Ogawahigashi, Kodaira, Tokyo 187-8553, Japan

Andrea Benedetti

Centre for Outcomes Research & Evaluation, Research Institute of the McGill University Health Centre, 5252 Boulevard de Maisonneuve, Montréal, QC, H4A 3S5, Canada

Brett D. Thombs

Jewish General Hospital; 4333 Cote Ste Catherine Road; Montreal, Quebec H3T 1E4, Canada

Author Contributions.

YW, BLevis, JB, PC, SG, JPAI, LAK, DM, SBP, IS, RCZ, AB, and BDT were responsible for the study conception and design. JB and LAK designed and conducted database searches to identify eligible studies. DHA, BA, LA, HRB, MB, CHB, PB, GC, MHC, JCNC, RC, YC, JMG, JRF, FHF, DF, BG, FGS, CGG, BJH, PAH, MHärter, UH, LH, SEH, MHudson, TH, MI, NJ, MEK, KMK, YK, FL,

analysis

A full list of authors and affiliations appears at the end of the article.

Abstract

Background: Item 9 of the Patient Health Questionnaire-9 (PHQ-9) queries about thoughts of death and self-harm, but not suicidality. Although it is sometimes used to assess suicide risk, most positive responses are not associated with suicidality. The PHQ-8, which omits Item 9, is thus increasingly used in research. We assessed equivalency of total score correlations and the diagnostic accuracy to detect major depression of the PHQ-8 and PHQ-9.

Methods: We conducted an individual patient data meta-analysis. We fit bivariate random-effects models to assess diagnostic accuracy.

Results: 16,742 participants (2,097 major depression cases) from 54 studies were included. The correlation between PHQ-8 and PHQ-9 scores was 0.996 (95% confidence interval 0.996 to 0.996). The standard cutoff score of 10 for the PHQ-9 maximized sensitivity + specificity for the PHQ-8 among studies that used a semi-structured diagnostic interview reference standard (N = 27). At cutoff 10, the PHQ-8 was less sensitive by 0.02 (-0.06 to 0.00) and more specific by 0.01 (0.00 to 0.01) among those studies (N = 27), with similar results for studies that used other types of interviews (N = 27). For all 54 primary studies combined, across all cutoffs, the PHQ-8 was less sensitive than the PHQ-9 by 0.00 to 0.05 (0.03 at cutoff 10), and specificity was within 0.01 for all cutoffs (0.00 to 0.01).

Conclusions: PHQ-8 and PHQ-9 total scores were similar. Sensitivity may be minimally reduced with the PHQ-8, but specificity is similar.

INTRODUCTION

The 9-item Patient Health Questionnaire (PHQ-9) (Kroenke, Spitzer & Williams, 2001) is a self-report questionnaire that is commonly used for identifying people who may have depression based on matching symptoms to diagnostic criteria or, more commonly, on a

SL, ML, SRL, BLöwe, AM, SM, TNM, KM, FLO, VP, BWP, PP, AP, KR, AGR, ISS, JS, ASidebottom, ASimming, LS, SS, PLLT, AT, HCvW, JW, MAW, KW, MY, and BDT were responsible for collection of primary data included in this study. BLevis, KER, NS, AWL, MA, DBR, and BDT contributed to data extraction and coding for the meta-analysis. YW, BLevis, AB, and BDT contributed to the data analysis and interpretation. YW, BLevis, AB, and BDT contributed to drafting the manuscript. All authors provided a critical review and approved the final manuscript. AB and BDT are guarantors.

Conflict of Interest.

Drs. Jetté and Patten declare that they received a grant, outside the submitted work, from the University of Calgary Hotchkiss Brain Institute, which was jointly funded by the Institute and Pfizer. Pfizer was the original sponsor of the development of the PHQ-9, which is now in the public domain. Dr. Chan is a steering committee member or consultant of Astra Zeneca, Bayer, Lilly, MSD and Pfizer. She has received sponsorships and honorarium for giving lectures and providing consultancy and her affiliated institution has received research grants from these companies. Dr. Hegerl declares that within the last three years, he was an advisory board member for Lundbeck and Servier; a consultant for Bayer Pharma; a speaker for Roche Pharma and Servier; and received personal fees from Janssen, all outside the submitted work. Dr. Inagaki declares that he has received a grant from Novartis Pharma, and personal fees from Meiji, Mochida, Takeda, Novartis, Yoshitomi, Pfizer, Eisai, Otsuka, MSD, Technomics, and Sumitomo Dainippon, all outside of the submitted work. All authors declare no other relationships or activities that could appear to have influenced the submitted work. No funder had any role in the design and conduct of the study; collection, management, analysis, and interpretation of the data; preparation, review, or approval of the manuscript; and decision to submit the manuscript for publication.

standard cutoff of a score of 10 or greater (Moriarty et al. 2015; Levis et al. 2019). It is also used as a continuous measure to assess depressive symptom severity in research and clinical care (Kroenke, Spitzer & Williams, 2001). The nine items of the PHQ-9 are designed to capture the nine Diagnostic and Statistical Manual of Mental Disorders (DSM) symptom criteria for a major depressive episode.¹ Response options on the items range from “not at all” (0 points) to “nearly every day” (3 points). Per the DSM-5, the ninth criterion for major depression requires “Recurrent thoughts of death (not just fear of dying), recurrent suicidal ideation without a specific plan, or a suicide attempt or a specific plan for committing suicide” (American Psychiatric Association 2013). Item 9 of the PHQ-9 taps into this criterion but also assesses self-harm, which is not part of the DSM criterion, or passive thoughts of death within the last two weeks: “...how often have you been bothered by... thoughts that you would be better off dead or of hurting yourself in some way?” It does not query specifically about suicidality, and positive responses may be due to thoughts about death or to thoughts about self-harm.

Item 9 is sometimes used as an indicator of suicide risk, and it may be useful as a component of modelling approaches for stratifying suicide risk among participants in psychiatric settings (Simon et al. 2016; Simon et al. 2013). However, responses to the item may not accurately reflect whether or not suicide risk is present, particularly among patients with serious medical conditions for whom thoughts of death may not reflect suicidal ideation, and it appears to perform poorly in identifying individuals at risk in these settings. Four studies in non-psychiatric settings have compared positive responses on Item 9 to responses to questions that explicitly assess suicidal thoughts or intentionality. In these studies, which included US military veterans in primary care (Corson, Gerrity & Dobscha, 2004), patients with coronary artery disease (Razykov et al. 2012; Suarez et al. 2015), and cancer patients (Walker et al. 2011), 7% to 21% of all study participants had positive responses on Item 9, but of those, only 18% to 35% had thoughts of suicide based on questions designed specifically to address suicide risk, and only 3% to 20% had a plan (Corson, Gerrity & Dobscha, 2004; Razykov et al. 2012; Suarez et al. 2015; Walker et al. 2011). Thus, concerns have been raised that using Item 9 to identify individuals at risk would result in a high rate of false indications, compared to questions designed specifically to assess suicidal thoughts or intentionality (Razykov et al. 2012; Suarez et al. 2015; Walker et al. 2011).

The PHQ-8 omits Item 9 from the PHQ-9. Many research studies use the PHQ-8 as a depression screening tool or to assess depressive symptom severity in order to avoid the high risk of inaccurate indications of suicide risk based on Item 9 (Corson, Gerrity & Dobscha, 2004; Razykov et al. 2012; Kroenke et al. 2009; Wells et al. 2013; Barrera et al. 2017). This is a particularly important consideration in studies that are not focused on depression or psychiatric disorders, but would need to allocate substantial resources to follow-up on responses to Item 9 of the PHQ-9. Similarly, many large epidemiological studies that include assessments of depressive symptoms are not able to provide adequate assessment and intervention with telephone or internet surveys (Kroenke et al. 2009).

Although differences in performance between the PHQ-8 and PHQ-9 might be expected to be minimal, to the best of our knowledge, only one study has attempted to verify this by comparing diagnostic accuracy between the PHQ-8 and PHQ-9 (Razykov et al. 2012). That

study evaluated the diagnostic testing accuracy of the PHQ-8 versus the PHQ-9 and the correlation between PHQ-8 and PHQ-9 scores in a sample of 1,022 coronary artery disease outpatients (233 major depression cases). Differences between sensitivity and specificity for the PHQ-8 (50%, 91%) and PHQ-9 (54%, 90%) based on a cutoff score of 10 or greater were minimal. In addition, PHQ-8 and PHQ-9 scores were highly correlated ($r = 0.997$) (Razykov et al. 2012). One additional study reported correlations between continuous PHQ-8 and PHQ-9 scores (Corson, Gerrity & Dobscha, 2004). That study, which included over 1000 patients from a US Department of Veterans Affairs primary care setting, reported a correlation of $r = 0.998$ (Corson, Gerrity & Dobscha, 2004).

We have synthesized a large database of individual participant data (IPD) from primary studies on the PHQ-9 (Levis et al. 2019; Levis et al. 2018). In the present study we included studies from that database that provided individual item scores (not just total PHQ-9 scores), which allowed for calculation of PHQ-8 scores. The objectives of the present study were (1) to evaluate the equivalency of the correlation between PHQ-8 and PHQ-9 scores for assessing depressive symptom severity; and (2) to assess the equivalency of the diagnostic accuracy of PHQ-8 and PHQ-9 across relevant cutoffs for screening to detect major depression.

METHOD

Data Source

The present study used a subset of participants from an IPD database of PHQ-9 (Levis et al. 2019; Levis et al. 2018). The main PHQ-9 IPD meta-analysis (IPDMA) was registered in PROSPERO (CRD42014010673), and a protocol was published (Thombs et al. 2014). Analyses of the diagnostic accuracy of the PHQ-8 were conducted according to protocol with two exceptions: (1) we stratified results by reference standard categories and (2) we added an examination of equivalency with the PHQ-9. Results from the main IPDMA of the PHQ-9 are available elsewhere (Levis et al. 2019).

Search Strategy

A medical librarian searched Medline, Medline In-Process & Other Non-Indexed Citations, PsycINFO, and Web of Science from January 1, 2000 through February 7, 2015, using a peer-reviewed search strategy (Canadian Agency for Drugs and Technologies in Health 2016) (SupplementaryMethods1). We limited our search to these databases based on research showing that adding other databases when the Medline search is highly sensitive does not identify additional eligible studies (Rice et al. 2016). The search was limited to the year 2000 forward because the PHQ-9 was originally published in 2001 (Kroenke, Spitzer & Williams, 2001). In addition to the database search, we reviewed reference lists of relevant reviews and queried contributing authors about non-published studies. Search results were uploaded into RefWorks (RefWorks-COS, Bethesda, MD, USA). After de-duplication, unique citations were uploaded into DistillerSR (Evidence Partners, Ottawa, Canada), which was used to store and track search results, conduct screening for eligibility, document correspondence with primary study authors, and extract study characteristics.

Identification of Eligible Studies

Datasets from articles in any language were eligible for inclusion if they included diagnostic classification among participants aged 18 or older for current Major Depressive Disorder (MDD) or Major Depressive Episode (MDE) based on a validated semi-structured or fully structured interview conducted within two weeks of PHQ-9 administration, since diagnostic criteria for major depression are for symptoms in the last two weeks. Datasets where not all participants were at least 18 years of age were included if the primary data allowed us to select participants who were at least 18 years of age. Datasets where not all participants were administered the PHQ-9 within two weeks of the diagnostic interview were included if the primary data allowed us to select participants who were administered both instruments within two weeks. Data from studies where the PHQ-9 was administered exclusively to individuals with known psychiatric diagnoses or symptoms or who were seeking psychiatric care were excluded, because screening is not indicated for patients already seeking care or managed in psychiatric settings. For defining major depression cases, we considered MDD or MDE based on the DSM or the International Classification of Diseases (ICD). If more than one was reported, we prioritized MDE over MDD and DSM over ICD. Across all studies, there were only 23 discordant diagnoses that depended on classification prioritization (0.1% of participants). For the present study, we only included primary studies that provided individual PHQ-9 item scores and not just PHQ-9 total scores, because only those datasets allowed us to generate PHQ-8 scores and compare the PHQ-8 with the PHQ-9.

Two investigators independently reviewed titles and abstracts for eligibility. If either reviewer deemed a study potentially eligible, full-text article review was done by two investigators, independently. Disagreement between reviewers after full-text review was resolved by consensus, consulting a third investigator when necessary. Translators were consulted to evaluate titles, abstracts and full-text articles for languages other than those for which team members were fluent.

Data Contribution and Synthesis

Authors of eligible datasets were invited to contribute de-identified primary data. Primary study country, clinical setting, language, and diagnostic interview administered were extracted from published reports by two investigators independently, with disagreements resolved by consensus. Countries were categorized as “very high”, “high”, or “low-medium” development level based on the United Nation’s human development index (Whiting et al. 2011). Recruitment settings were categorized as “non-medical”, “primary care”, “inpatient specialty care”, or “outpatient specialty care.” Participant-level data included age, sex, major depression status, current diagnosis or treatment for a mental health problem, and PHQ-9 item scores. In two primary studies, multiple recruitment settings were included, thus recruitment setting was coded at the participant-level. When primary study datasets included appropriate statistical weighting to reflect sampling procedures, we used the provided weights. For studies where sampling procedures merited weighting, but the original study did not weight, we constructed appropriate weights using inverse selection probabilities. Weighting occurred, for instance, when all participants with positive screens, but only a

random subset of participants with negative screens, were administered a diagnostic interview.

Individual participant data were converted to a standard format and entered into a single dataset that also included study-level data. We compared published participant characteristics and diagnostic accuracy results with those obtained using the raw datasets. When primary data and original publications were discrepant, we identified and corrected errors when possible and resolved any outstanding discrepancies in consultation with the original investigators.

Statistical Analyses

To evaluate the equivalence of the PHQ-8 and PHQ-9 scores for assessing depressive symptom severity, a Pearson correlation with a 95% confidence interval (CI) was calculated between the total scores of PHQ-8 (which excluded Item 9) and PHQ-9.

To estimate the diagnostic accuracy of the PHQ-8 and compare with the PHQ-9, we analyzed primary studies separately by the type of diagnostic interview that was used as the reference standard, as we did in the previously published main PHQ-9 meta-analysis (Levis et al. 2019). This was done because of differences in the performance of the different types of interviews. Semi-structured interviews involve clinical judgement and are designed to be administered by clinically trained professionals; fully structured interviews are completely scripted and designed for lay administration, but the resulting increased standardization and reliability across interviewers may lead to increased misclassification (Brugha, Bebbington & Jenkins, 1999; Nosen & Woody 2008). The Mini International Neuropsychiatric Interview (MINI), which is a fully structured interview, was developed to be administered in a fraction of the time necessary for other fully structured interviews and was described by its developers as designed to be over-inclusive (Robins et al. 1988; Sheehan et al. 1997). In a previous study, we found that semi-structured and fully structured diagnostic interviews are not interchangeable reference standards for major depression and that fully structured interviews may diagnose depression at higher rates than semi-structured interviews at low symptom levels and at lower rates at high symptom levels (Levis et al. 2018). We also found that the MINI classifies approximately twice as many participants as cases compared to the most commonly used fully-structured interview, the Composite International Diagnostic Interview (CIDI) (Levis et al. 2018). In the main PHQ-9 meta-analysis, the diagnostic accuracy of the PHQ-9 differed substantially depending on the reference standard used for the comparison (Levis et al. 2019). Thus, for the present study, we analyzed primary studies separately based on whether they used a semi-structured interview, a fully structured interview (non-MINI), or the MINI.

For each reference standard and for the PHQ-8 and PHQ-9 cutoffs 5–15, separately, bivariate random-effects models were fitted using an adaptive Gauss Hermite quadrature with 1 quadrature point (Riley et al. 2008). This 2-stage meta-analytic approach models sensitivity and specificity at the same time, taking the inherent correlation between them and the precision of estimates within studies into account. A random-effects model was used as we assumed true values of sensitivity and specificity would likely to vary across primary studies.

In order to examine the equivalence between PHQ-8 and PHQ-9 across reference standards, for each analysis, we used the results of the random-effects meta-analyses at each cutoff to construct separate empirical receiver operating characteristic (ROC) curves based on the pooled estimates. Equivalence tests between PHQ-8 and PHQ-9 sensitivity and specificity were conducted at each cutoff. This allowed us to test whether the sensitivity and specificity of the PHQ-8 was similar to that of the PHQ-9, up to a pre-specified maximum clinically acceptable difference, that is, an equivalence margin (Walker & Nowacki 2011). In the present study, an equivalence margin of $\delta = 0.05$ was used, which is the same margin that was used in a previous study that used the same IPD database (Ishihara et al. 2018). CIs for the differences between PHQ-8 and PHQ-9 sensitivity and specificity at each cutoff were constructed via a cluster bootstrap approach (van der Leeden, Busing & Meijer, 1997; van der Leeden, Meijer & Busing, 2008), with resampling at the study and subject level. For each comparison, we ran 1000 iterations of the bootstrap. For each bootstrap iteration, the bivariate random-effects model was fitted to the PHQ-8 and PHQ-9 data, and the pooled sensitivities and specificities were computed separately, as described above, for each cutoff score. Equivalence tests were done by comparing the CIs around the pooled sensitivity and specificity differences to the equivalence margin of $\delta = 0.05$. If the entire CI was included within the interval of ± 0.05 , then we rejected the hypothesis that there were differences large enough to be important and concluded that equivalence was present. If the entire CI was outside of the interval, then we failed to reject the hypothesis that the PHQ-8 and PHQ-9 were not equivalent. When the CIs crossed the ± 0.05 threshold, findings were deemed equivocal, and the equivalence was indeterminate.

Although we previously found that sensitivity and specificity of the PHQ-9 differs by type of reference standard, we did not believe that the differences in sensitivity and specificity between the PHQ-8 and PHQ-9 would vary depending on the reference standard. This is because for each included study the PHQ-8 and PHQ-9 were compared to the same reference standard. Thus, we reported pooled sensitivity and specificity only stratified by reference standards, but we investigated equivalence both stratified by reference standards and pooled across all studies. To investigate heterogeneity across studies, by reference standard and overall, we generated forest plots for the differences in sensitivity and specificity estimates between PHQ-8 and PHQ-9 for the standard cutoff 10 for each study. We also quantified heterogeneity at cutoff 10, by reporting the estimated variances of the random effects for the differences in PHQ-8 and PHQ-9 sensitivity and specificity (τ^2) (Fagerland, Lydersen & Laake, 2014; Higgins & Thompson 2002).

All analyses were run in R (R version R 3.5.0 and R Studio version 1.1.423) using the lme4 package.

RESULTS

Search Results and Inclusion of Primary Data

For the main IPDMA, of 5,248 unique titles and abstracts identified from the database search, 5,039 were excluded after title and abstract review and 113 after full-text (SupplementaryList1), leaving 96 eligible articles with data from 69 unique participant samples (SupplementaryFigure1). Of the 69 unique samples, 55 contributed data (80%). In

addition, authors of included studies contributed data from three unpublished studies, for a total of 58 PHQ-9 datasets contributed to our IPDMA. Four studies without PHQ-9 individual item scores were excluded from the present study (see SupplementaryTable1b). Thus, 16,742 participants (2,097 major depression cases) from 54 studies were analyzed (78% of 21,572 participants from the 69 eligible published studies and 3 eligible unpublished studies). Included study characteristics are shown in SupplementaryTable1a. Characteristics of eligible studies that did not provide data, including the 4 studies excluded because they only provided PHQ-9 total scores, are shown in SupplementaryTable1b.

There were 27 included primary studies that used semi-structured interviews to assess major depression (6,362 participants), 13 that used fully structured interviews other than the MINI (7,596 participants), and 14 that used the MINI (2,784 participants). The Structured Clinical Interview for DSM Disorders (SCID) was the most commonly used semi-structured interview (24 studies, 4,378 participants), and the CIDI was the most commonly used fully structured interview (10 studies, 6,291 participants). The average study sample size and number of major depression cases was 236 and 29 for studies that used a semi-structured interview; 584 and 61 for studies that used a fully structured interview; and 199 and 37 for studies that used the MINI.

Participant characteristics are shown in Table 1.

PHQ-8 and PHQ-9 Scores

Among all participants in all studies, the mean (standard deviation [SD]) PHQ-8 score (range = 0–24) was 5.3 (5.2), and the mean (SD) PHQ-9 score (range 0–27) was 5.4 (5.4). Overall, 11.8% of participants had a non-zero score on Item 9 (score of 1–3). As shown in Table 2, this included 1.9% among participants with PHQ-8 scores 0–4 and increased to 64.7% among those with scores 20–24. The correlation (95% CI) between PHQ-8 and PHQ-9 scores was 0.996 (0.996, 0.996). The correlation of the score of Item 9 with PHQ-8 scores was 0.480 (0.469, 0.492).

Diagnostic Accuracy of the PHQ-8 and PHQ-9

ROC curves comparing sensitivity and specificity estimates for cutoffs 5–15 between the PHQ-8 and PHQ-9 for the three reference standard categories, separately, are shown in Figure 1. The curves for the PHQ-8 and PHQ-9 were highly overlapping for each reference standard, and the area under the curve for the PHQ-8 and PHQ-9 were similar for semi-structured interviews (0.930 versus 0.933), fully structured interviews (excluding the MINI; 0.852 versus 0.856), and the MINI (0.894 versus 0.899).

Comparisons of sensitivity and specificity estimates between PHQ-8 and PHQ-9 cutoffs 5–15 across the three reference standard categories are shown in Table 3. Cutoff 10 maximized combined sensitivity and specificity for PHQ-8 (sensitivity [95% CI] = 0.86 [0.80, 0.90], specificity [95% CI] = 0.86 [0.83, 0.89]) and PHQ-9 (sensitivity [95% CI] = 0.88 [0.82, 0.92], specificity [95% CI] = 0.86 [0.82, 0.88]) among studies using a semi-structured interview as the reference standard; cutoff 8 for PHQ-8 (sensitivity [95% CI] = 0.77 [0.66, 0.85], specificity [95% CI] = 0.78 [0.71, 0.84]) and PHQ-9 (sensitivity [95% CI] = 0.79 [0.68, 0.86], specificity [95% CI] = 0.77 [0.70, 0.83]) among studies using a fully structured

interview; and cutoff 8 for PHQ-8 (sensitivity [95% CI] = 0.83 [0.75, 0.89], specificity [95% CI] = 0.80 [0.75, 0.84]) and PHQ-9 (sensitivity [95% CI] = 0.86 [0.77, 0.91], specificity [95% CI] = 0.79 [0.74, 0.83]) among studies using the MINI.

In comparisons stratified by reference standard, for sensitivity, results of equivalence tests showed that for semi-structured diagnostic interviews, estimates were equivalent from cutoffs 5 through 9 and indeterminate from cutoffs 10 through 15; for fully structured interviews (excluding the MINI), they were equivalent on cutoffs 5 and 7 and indeterminate at cutoffs 6 and 8 through 15; and for the MINI, they were equivalent from cutoffs 5 through 7 and indeterminate from cutoffs 8 through 15. Estimates of specificity were equivalent in all analyses, regardless of reference standards and cutoffs. See Table 3.

Overall, including all 54 primary studies, as shown in Table 4, across cutoffs, sensitivity was between 0.00 and 0.05 percentage points lower for the PHQ-8 compared to the PHQ-9. At cutoff 10, the difference (95% CI) was -0.03 (-0.06 , -0.02). For specificity, the PHQ-8 and PHQ-9 were within 0.01 for all cutoffs. For sensitivity, estimates were equivalent for cutoffs 5 to 8 and indeterminate for cutoffs 9 to 15. For specificity, estimates were equivalent for all cutoffs.

A forest plot of the difference of sensitivity and specificity estimates for cutoff 10 between PHQ-8 and PHQ-9 for all studies is shown in Figure 2. At the commonly used cutoff of 10 or greater, there was low heterogeneity in the differences across the 54 studies with estimated inter-study heterogeneity (τ^2) <0.01 for sensitivity and <0.01 for specificity. Forest plots of the differences of sensitivity and specificity estimates for cutoff 10 between PHQ-8 and PHQ-9 among studies by reference standard category are shown in SupplementaryFigure2.

DISCUSSION

In the present study, we assessed the correlation of continuous PHQ-8 and PHQ-9 scores for assessing depression severity in research and clinical practice, and we compared the diagnostic accuracy of the PHQ-8 and PHQ-9 across all cutoffs for detecting major depression. There were two main findings. First, the correlation of continuous PHQ-8 and PHQ-9 scores was high (0.996). Second, to screen for major depression, the PHQ-8 at different possible cutoffs including the standard cutoff of 10 or greater, was similarly accurate compared to the PHQ-9 overall and across all three types of reference standards. The cutoffs that maximized combined sensitivity and specificity were the same for the PHQ-8 and PHQ-9 across reference standard categories.

Overall, for all 54 primary studies combined, across all cutoffs, the PHQ-8 was slightly less sensitive than the PHQ-9 by 0.00 to 0.05 (0.03 at cutoff 10). For specificity, the PHQ-8 and PHQ-9 were within 0.01 for all cutoffs. Although the CIs for the difference in sensitivity for cutoff 10 did not fit the study definition of equivalency, the reduction in sensitivity if the PHQ-8 is used is small, and specificity is equivalent.

Previous studies have shown that Item 9 of the PHQ-9 does not accurately assess suicide risk and identifies far more patients or study participants as at risk than would be identified with

items designed to assess suicide risk (Corson, Gerrity & Dobscha, 2004; Razykov et al. 2012; Walker et al. 2011). Thus, unintended consequences of using Item 9 could include substantial additional costs for research, as well as possible harms or inconvenience to patients. Research ethics boards sometimes require follow-up for all patients with positive responses to Item 9. Using the PHQ-8, which is minimally different from PHQ-9 in terms of diagnostic accuracy characteristics, would reduce unintended consequences of false signals of suicide risk without substantive changes to continuous measurement properties or diagnostic accuracy for major depression.

It is possible that use of the PHQ-8 could result in not identifying a small subset of people with suicidal thoughts, although, if the case, based on our findings, this number would be small. Furthermore, there is no evidence that using questionnaires to screen for suicide in general medical settings, above and beyond screening for depression, would reduce risk of suicide (Allaby 2010; Crawford et al. 2011; Siu & the US Preventive Services Task Force 2016). Tools are available to screen patients or stratify by risk for suicidality. However, a review concluded that they are not accurate enough at this point for use in practice and that alternative methods are more appropriate (Carter & Spittal 2018). Indeed, in mental health settings or when there is reason to suspect possible suicidality, standards of care indicate that engagement with patients is needed to assess suicide risk and determine the best management plan, as appropriate (Carter & Spittal 2018).

To our knowledge, this is the first meta-analysis and also the first study using a large IPD database to compare diagnostic accuracy characteristics of PHQ-8 and PHQ-9. Strengths of this study included the large overall sample size, the ability to compare results for PHQ-8 and PHQ-9 from all cutoffs from all studies (rather than just published cutoff results), and the ability to assess diagnostic accuracy separately in studies that used semi- and fully structured diagnostic interviews as the reference standard. There are also limitations to consider. First, for the full IPDMA data, we were unable to include primary data from 14 of 69 published eligible datasets (20% datasets; 17% of eligible participants), and we restricted our analyses to those with complete data for all individual PHQ-9 item scores (95% of available data). Nonetheless, this sample was much larger than the few previous primary studies that have compared the PHQ-8 and PHQ-9. Second, we categorized studies based on the diagnostic interview that was used, but adaptations to interviews are sometimes made and, thus, all studies may have not used the diagnostic interviews in the way that they were originally designed. However, when we analyzed data from all studies, regardless of reference standard, heterogeneity was minimal, suggesting that findings can be applied across reference standards.

CONCLUSIONS

In summary, although the PHQ-9 was designed to reflect the 9 symptoms included in DSM criteria for major depression, the item assessing suicide risk also assesses self-harm. This study used a large IPD dataset and found that the PHQ-8 performs similarly to the PHQ-9 in terms of the correlation of continuous scores and the diagnostic accuracy across all cutoffs for detecting major depression. Removing Item 9 and using the PHQ-8 instead of the PHQ-9 has minimal influence on performance of the measure and will likely reduce the number of

false positive signals from people who endorse this item but would not be considered to be at risk for suicide based on measures intended to assess suicide risk.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Authors

Yin Wu, PhD, Brooke Levis, MSc, Kira E. Riehm, MSc, Nazanin Saadat, MSc, Alexander W. Levis, MSc, Marleine Azar, MSc, Danielle B. Rice, MSc, Jill Boruff, MLIS, Pim Cuijpers, PhD, Simon Gilbody, PhD, John P.A. Ioannidis, MD, Lorie A. Kloda, PhD, Dean McMillan, PhD, Scott B. Patten, MD, Ian Shrier, MD, Roy C. Ziegelstein, MD, Dickens H. Akena, PhD, Bruce Arroll, MBChB, Liat Ayalon, PhD, Hamid R. Baradaran, MD, Murray Baron, MD, Charles H. Bombardier, PhD, Peter Butterworth, PhD, Gregory Carter, FRANZCP, Marcos H. Chagas, MD, Juliana C. N. Chan, MD, Rushina Cholera, MD, Yeates Conwell, MD, Janneke M. de Man-van Ginkel, PhD, Jesse R. Fann, MD, Felix H. Fischer, PhD, Daniel Fung, MD, Bizu Gelaye, PhD, Felicity Goodyear-Smith, MD, Catherine G. Greeno, PhD, Brian J. Hall, PhD, Patricia A. Harrison, PhD, Martin Härter, MD, PhD Dipl Psych, Ulrich Hegerl, MD, Leanne Hides, PhD(Clin), Stevan E. Hobfoll, PhD, Marie Hudson, MD, Thomas Hyphantis, MD, PhD, Masatoshi Inagaki, MD, Nathalie Jetté, MD, Mohammad E. Khamseh, MD, Kim M. Kiely, PhD, Yunxin Kwan, MMed (Psychiatry), Femke Lamers, PhD, Shen-Ing Liu, MD, Manote Lotrakul, MD, Sonia R. Loureiro, PhD, Bernd Löwe, MD, Anthony McGuire, PhD, Sherina Mohd-Sidik, PhD, Tiago N. Munhoz, PhD, Kumiko Muramatsu, MD, Flávia L. Osório, PhD, Vikram Patel, MD, Brian W. Pence, PhD, Philippe Persoons, MD, Angelo Picardi, MD, Katrin Reuter, PhD, Alasdair G. Rooney, MD, Iná S. Santos, MD, Juwita Shaaban, MMed (Fam. Med), Abbey Sidebottom, PhD, Adam Simning, MD, PhD, Lesley Stafford, PhD, Sharon Sung, PhD, Pei Lin Lynnette Tan, MMed (Psychiatry), Alyna Turner, PhD, Henk C. van Weert, MD, Jennifer White, PhD, Mary A. Whooley, MD, Kirsty Winkley, PhD, Mitsuhiro Yamada, MD, Andrea Benedetti, PhD, Brett D. Thombs, PhD.

Affiliations

Lady Davis Institute for Medical Research, Jewish General Hospital, Montréal, Québec, Canada (Wu, Levis B, Riehm, Saadat, Levis A, Azar, Rice, Shrier, Baron, Hudson, Thombs); Department of Psychiatry, McGill University, Montréal, Québec, Canada (Wu, Thombs); Department of Epidemiology, Biostatistics and Occupational Health, McGill University, Montréal, Québec, Canada (Wu, Levis B, Levis A, Azar, Shrier, Benedetti, Thombs); Department of Medicine, McGill University, Montréal, Québec, Canada (Baron, Hudson, Benedetti, Thombs); Department of Psychology, McGill University, Montréal, Québec, Canada (Rice, Thombs); Schulich Library of Physical Sciences, Life Sciences, and Engineering, McGill University, Montreal, Quebec, Canada (Boruff); Department of Clinical, Neuro and Developmental Psychology, Amsterdam Public Health Research Institute, Vrije Universiteit, Amsterdam, the Netherlands (Cuijpers); Hull York Medical School and the

Department of Health Sciences, University of York, Heslington, York, UK (Gilbody, McMillan); Department of Medicine, Department of Health Research and Policy, Department of Biomedical Data Science, Department of Statistics, Stanford University, Stanford, California, USA (Ioannidis); Library, Concordia University, Montréal, Québec, Canada (Kloda); Department of Community Health Sciences, University of Calgary, Calgary, Alberta, Canada (Patten, Jetté); Hotchkiss Brain Institute and O'Brien Institute for Public Health, University of Calgary, Calgary, Alberta, Canada (Patten, Jetté); Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, Maryland, USA (Ziegelstein); Department of Psychiatry, Makerere University College of Health Sciences, Kampala, Uganda (Akena); Department of General Practice and Primary Health Care, University of Auckland, New Zealand (Arroll, Goodyear-Smith); Louis and Gabi Weisfeld School of Social Work, Bar Ilan University, Ramat Gan, Israel (Ayalon); Endocrine Research Center, Institute of Endocrinology and Metabolism, Iran University of Medical Sciences, Tehran, Iran (Baradaran, Khamseh); Ageing Clinical & Experimental Research Team, Institute of Applied Health Sciences, School of Medicine, Medical Sciences and Nutrition, University of Aberdeen, Aberdeen, Scotland, UK (Baradaran); Department of Rehabilitation Medicine, University of Washington, Seattle, Washington, USA (Bombardier); Centre for Research on Ageing, Health and Wellbeing, Research School of Population Health, The Australian National University, Canberra, Australia (Butterworth); Melbourne Institute of Applied Economic and Social Research, University of Melbourne, Melbourne, Australia (Butterworth); Centre for Brain and Mental Health Research, University of Newcastle, New South Wales, Australia (Carter); Department of Neurosciences and Behavior, Ribeirão Preto Medical School, University of São Paulo, Ribeirão Preto, Brazil (Chagas, Osório, Loureiro); Department of Medicine and Therapeutics, Prince of Wales Hospital, The Chinese University of Hong Kong, Hong Kong Special Administrative Region, China (Chan); Asia Diabetes Foundation, Prince of Wales Hospital, Hong Kong Special Administrative Region, China (Chan); Hong Kong Institute of Diabetes and Obesity, Hong Kong Special Administrative Region, China (Chan); Department of Pediatrics, University of North Carolina at Chapel Hill School of Medicine, Chapel Hill, North Carolina, USA (Cholera); Department of Psychiatry, University of Rochester Medical Center, New York, USA (Conwell, Simning); Julius Center for Health Sciences and Primary Care, University Medical Center Utrecht, Utrecht, the Netherlands (de Man-van Ginkel); Department of Psychiatry and Behavioral Sciences, University of Washington, Seattle, Washington, USA (Fann); Department of Psychosomatic Medicine, Center for Internal Medicine and Dermatology, Charité - Universitätsmedizin Berlin, Germany (Fischer); Department of Child & Adolescent Psychiatry, Institute of Mental Health, Singapore (Fung, Sung); Yong Loo Lin School of Medicine, National University of Singapore, Singapore (Fung); Lee Kong Chian School of Medicine, Nanyang Technological University, Singapore (Fung); Programme in Health Services & Systems Research, Duke-NUS Medical School, Singapore (Fung, Liu, Sung); Department of Epidemiology, Harvard T. H. Chan School of Public Health, Boston, Massachusetts, USA (Gelaye); School of Social Work, University of Pittsburgh, Pittsburgh,

Pennsylvania, USA (Greeno); Global and Community Mental Health Research Group, Department of Psychology, Faculty of Social Sciences, University of Macau, Macau Special Administrative Region, China (Hall); Department of Health, Behavior, and Society, Johns Hopkins Bloomberg School of Public Health, Baltimore, Maryland, USA (Hall); City of Minneapolis Health Department, Minneapolis, Minnesota, USA (Harrison); Department of Medical Psychology, University Medical Center Hamburg-Eppendorf, Hamburg, Germany (Härter); Department of Psychiatry and Psychotherapy, University Hospital Leipzig, Leipzig, Germany (Hegerl); School of Psychology, University of Queensland, Brisbane, Queensland, Australia (Hides); STAR-Stress, Anxiety & Resilience Consultants, Chicago, Illinois, USA (Hobfoll); Department of Psychiatry, University of Ioannina, Ioannina, Greece (Hyphantis); Department of Psychiatry, Faculty of Medicine, Shimane University, Shimane, Japan (Inagaki); Department of Neurology, Icahn School of Medicine at Mount Sinai, New York, New York, USA (Jetté); School of Psychology, University of New South Wales, Sydney, Australia (Kiely); Neuroscience Research Australia, Sydney, Australia (Kiely); Department of Psychological Medicine, Tan Tock Seng Hospital, Singapore (Kwan, Tan); Department of Psychiatry, Amsterdam UMC, Vrije Universiteit, Amsterdam Public Health Research Institute, Amsterdam, the Netherlands (Lamers); Department of Psychiatry, Mackay Memorial Hospital, Taipei, Taiwan (Liu); Department of Medical Research, Mackay Memorial Hospital, Taipei, Taiwan (Liu); Department of Medicine, Mackay Medical College, Taipei, Taiwan (Liu); Department of Psychiatry, Faculty of Medicine, Ramathibodi Hospital, Mahidol University, Bangkok, Thailand (Lotrakul); Department of Psychosomatic Medicine and Psychotherapy, University Medical Center Hamburg-Eppendorf, Hamburg, Germany (Löwe); Department of Nursing, St. Joseph's College, Standish, Maine, USA (McGuire); Cancer Resource & Education Centre, and Department of Psychiatry, Faculty of Medicine and Health Sciences, Universiti Putra Malaysia, Serdang, Selangor, Malaysia (Mohd-Sidik); Post-graduate Program in Epidemiology, Federal University of Pelotas, Pelotas, RS, Brazil (Munhoz, Santos); Department of Clinical Psychology, Graduate School of Niigata Seiryō University, Niigata, Japan (Muramatsu); National Institute of Science and Technology, Translational Medicine, Ribeirão Preto, Brazil (Osório); Department of Global Health and Social Medicine, Harvard Medical School, Boston, Massachusetts, USA (Patel); Department of Global Health and Population, Harvard T.H. Chan School of Public Health, Boston, Massachusetts, USA (Patel); Department of Epidemiology, Gillings School of Global Public Health, The University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA (Pence); Department of Adult Psychiatry, University Hospitals Leuven, Leuven, Belgium (Persoons); Department of Neurosciences, Katholieke Universiteit Leuven, Leuven, Belgium (Persoons); Centre for Behavioural Sciences and Mental Health, Italian National Institute of Health, Rome, Italy (Picardi); Department of Psychiatry and Psychotherapy, University Medical Center Freiburg, Freiburg, Germany (Reuter); Division of Psychiatry, Royal Edinburgh Hospital, University of Edinburgh, Edinburgh, Scotland, UK (Rooney); Department of Family Medicine, School of Medical Sciences, Universiti Sains Malaysia, Kelantan, Malaysia (Shaaban); Allina Health, Minneapolis, Minnesota, USA (Sidebottom);

Centre for Women's Mental Health, Royal Women's Hospital, Parkville, Australia (Stafford); Melbourne School of Psychological Sciences, University of Melbourne, Australia (Stafford); School of Medicine and Public Health, University of Newcastle, New South Wales, Newcastle, Australia (Turner); IMPACT Strategic Research Centre, School of Medicine, Deakin University, Geelong, Victoria, Australia (Turner); Department of General Practice, Amsterdam Institute for General Practice and Public Health, Amsterdam University Medical Centers, location AMC (van Weert); Monash University, Melbourne, Australia (White); Department of Epidemiology and Biostatistics, University of California San Francisco, San Francisco, California, USA (Whooley); Department of Medicine, Veterans Affairs Medical Center, San Francisco, California, USA (Whooley); Department of Medicine, University of California San Francisco, San Francisco, California, USA (Whooley); Florence Nightingale Faculty of Nursing, Midwifery & Palliative Care, King's College London, London, UK (Winkley), Department of Neuropsychopharmacology, National Institute of Mental Health, National Center of Neurology and Psychiatry, Ogawa-Higashi, Kodaira, Tokyo, Japan (Yamada); Respiratory Epidemiology and Clinical Research Unit, McGill University Health Centre, Montréal, Québec, Canada (Benedetti); Department of Educational and Counselling Psychology, McGill University, Montréal, Québec, Canada (Thombs).

Acknowledgements

This study was funded by the Canadian Institutes of Health Research (KRS-134297; PCG-155468). Dr. Wu was supported by an Utting Postdoctoral Fellowship from the Jewish General Hospital, Montreal, Quebec. Drs. Benedetti and Thombs were supported by Fonds de recherche du Québec - Santé (FRQS) researcher salary awards. Ms. Levis was supported by a CIHR Frederick Banting and Charles Best Canada Graduate Scholarship doctoral award. Ms. Riehm and Ms. Saadat were supported by CIHR Frederick Banting and Charles Best Canadian Graduate Scholarships – Master's Awards. Mr. Levis and Ms. Azar were supported by FRQS Masters Training Awards. Ms. Rice was supported by a Vanier Canada Graduate Scholarship. Collection of data for the study by Arroll et al. was supported by a project grant from the Health Research Council of New Zealand. Data collection for the study by Ayalon et al. was supported from a grant from Lundbeck International. The primary study by Khamseh et al. was supported by a grant (M-288) from Tehran University of Medical Sciences. The primary study by Bombardier et al. was supported by the Department of Education, National Institute on Disability and Rehabilitation Research, Spinal Cord Injury Model Systems: University of Washington (grant no. H133N060033), Baylor College of Medicine (grant no. H133N060003), and University of Michigan (grant no. H133N060032). Dr. Butterworth was supported by Australian Research Council Future Fellowship FT130101444. Dr. Cholera was supported by a United States National Institute of Mental Health (NIMH) grant (5F30MH096664), and the United States National Institutes of Health (NIH) Office of the Director, Fogarty International Center, Office of AIDS Research, National Cancer Center, National Heart, Blood, and Lung Institute, and the NIH Office of Research for Women's Health through the Fogarty Global Health Fellows Program Consortium (1R25TW00934001) and the American Recovery and Reinvestment Act. Dr. Conwell received support from NIMH (R24MH071604) and the Centers for Disease Control and Prevention (R49 CE002093). The primary studies by Amoozegar and by Fiest et al. were funded by the Alberta Health Services, the University of Calgary Faculty of Medicine, and the Hotchkiss Brain Institute. The primary study by Fischer et al. was funded by the German Federal Ministry of Education and Research (01GY1150). Dr. Fischler was supported by a grant from the Belgian Ministry of Public Health and Social Affairs and a restricted grant from Pfizer Belgium. Data for the primary study by Gelaye et al. was supported by grant from the NIH (T37 MD001449). Collection of data for the primary study by Gjerdingen et al. was supported by grants from the NIMH (R34 MH072925, K02 MH65919, P30 DK50456). The primary study by Eack et al. was funded by the NIMH (R24 MH56858). Collection of data provided by Drs. Härter and Reuter was supported by the Federal Ministry of Education and Research (grants No. 01 GD 9802/4 and 01 GD 0101) and by the Federation of German Pension Insurance Institute. Collection of data for the primary study by Hobfoll et al. was made possible in part from grants from NIMH (RO1 MH073687) and the Ohio Board of Regents. Dr. Hall received support from a grant awarded by the Research and Development Administration Office, University of Macau (MYRG2015-00109-FSS). The primary study by Hides et al. was funded by the Perpetual Trustees, Flora and Frank Leith Charitable Trust, Jack Brockhoff Foundation, Grosvenor Settlement, Sunshine Foundation and Danks Trust. The primary study by Henkel et al. was funded by the German Ministry of Research and Education. Data for the study by Razykov et al. was

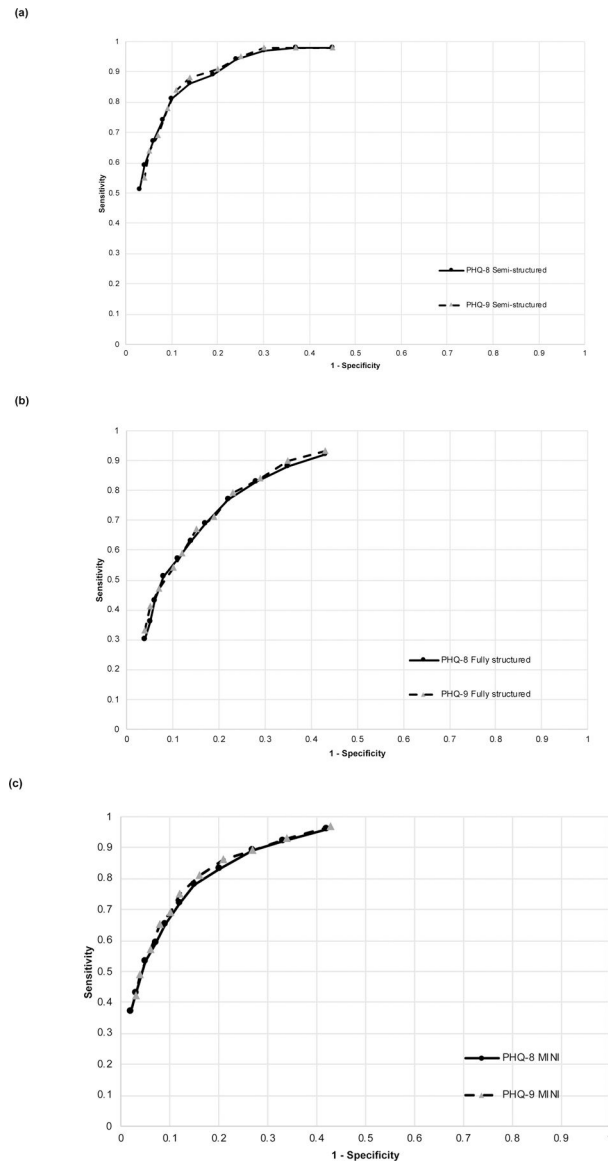
collected by the Canadian Scleroderma Research Group, which was funded by the CIHR (FRN 83518), the Scleroderma Society of Canada, the Scleroderma Society of Ontario, the Scleroderma Society of Saskatchewan, Sclérodermie Québec, the Cure Scleroderma Foundation, Inova Diagnostics Inc., Euroimmun, FRQS, the Canadian Arthritis Network, and the Lady Davis Institute of Medical Research of the Jewish General Hospital, Montreal, QC. Dr. Hudson was supported by a FRQS Senior Investigator Award. Collection of data for the primary study by Hyphantis et al. was supported by grant from the National Strategic Reference Framework, European Union, and the Greek Ministry of Education, Lifelong Learning and Religious Affairs (ARISTEIA-ABREVIATE, 1259). The primary study by Inagaki et al. was supported by the Ministry of Health, Labour and Welfare, Japan. Dr. Jetté was supported by a Canada Research Chair in Neurological Health Services Research. Collection of data for the primary study by Kiely et al. was supported by National Health and Medical Research Council (grant number 1002160) and Safe Work Australia. Dr. Kiely was supported by funding from an Australian National Health and Medical Research Council fellowship (grant number 1088313). The primary study by Lamers et al. was funded by the Netherlands Organisation for Health Research and Development (grant number 945-03-047). The primary study by Liu et al. was funded by a grant from the National Health Research Institute, Republic of China (NHRI-EX97-9706PI). The primary study by Lotrakul et al. was supported by the Faculty of Medicine, Ramathibodi Hospital, Mahidol University, Bangkok, Thailand (grant number 49086). Dr. Bernd Löwe received research grants from Pfizer, Germany, and from the medical faculty of the University of Heidelberg, Germany (project 12½000) for the study by Gräfe et al. The primary study by Mohd-Sidik et al. was funded under the Research University Grant Scheme from Universiti Putra Malaysia, Malaysia and the Postgraduate Research Student Support Accounts of the University of Auckland, New Zealand. The primary study by Santos et al. was funded by the National Program for Centers of Excellence (PRONEX/FAPERGS/CNPq, Brazil). The primary study by Muramatsu et al. was supported by an educational grant from Pfizer US Pharmaceutical Inc. Collection of primary data for the study by Dr. Pence was provided by NIMH (R34MH084673). The primary studies by Osório et al. were funded by Reitoria de Pesquisa da Universidade de São Paulo (grant number 09.1.01689.17.7) and Banco Santander (grant number 10.1.01232.17.9). Dr. Osório was supported by Productivity Grants (PQ-CNPq-2 -number 30132½016-7). The primary study by Picardi et al. was supported by funds for current research from the Italian Ministry of Health. Dr. Persoons was supported by a grant from the Belgian Ministry of Public Health and Social Affairs and a restricted grant from Pfizer Belgium. Dr. Shaaban was supported by funding from Universiti Sains Malaysia. The primary study by Rooney et al. was funded by the United Kingdom National Health Service Lothian Neuro-Oncology Endowment Fund. The primary study by Sidebottom et al. was funded by a grant from the United States Department of Health and Human Services, Health Resources and Services Administration (grant number R40MC07840). Simning et al.'s research was supported in part by grants from the NIH (T32 GM07356), Agency for Healthcare Research and Quality (R36 HS018246), NIMH (R24 MH071604), and the National Center for Research Resources (TL1 RR024135). Dr. Stafford received PhD scholarship funding from the University of Melbourne. Collection of data for the studies by Turner et al were funded by a bequest from Jennie Thomas through the Hunter Medical Research Institute. Collection of data for the primary study by Williams et al. was supported by a NIMH grant to Dr. Marsh (RO1-MH069666). The primary study by Thombs et al. was done with data from the Heart and Soul Study (PI Mary Whooley). The Heart and Soul Study was funded by the Department of Veterans Epidemiology Merit Review Program, the Department of Veterans Affairs Health Services Research and Development service, the National Heart Lung and Blood Institute (R01 HL079235), the American Federation for Aging Research, the Robert Wood Johnson Foundation, and the Ischemia Research and Education Foundation. Dr. Thombs was supported by an Investigator Award from the Arthritis Society. The primary study by Twist et al. was funded by the UK National Institute for Health Research under its Programme Grants for Applied Research Programme (grant reference number RP-PG-0606-1142). The study by Wittkamp et al. was funded by The Netherlands Organization for Health Research and Development (ZonMw) Mental Health Program (nos. 100.003.005 and 100.002.021) and the Academic Medical Center/University of Amsterdam. Collection of data for the primary study by Zhang et al. was supported by the European Foundation for Study of Diabetes, the Chinese Diabetes Society, Lilly Foundation, Asia Diabetes Foundation and Liao Wun Yuk Diabetes Memorial Fund. No other authors reported funding for primary studies or for their work on the present study.

REFERENCES

- Allaby M (2010). Screening for depression: A report for the UK National Screening Committee (Revised report) UK National Screening Committee: London, United Kingdom.
- American Psychiatric Association. (2013). Diagnostic and statistical manual of mental disorders (5th ed.). Arlington, VA: American Psychiatric Publishing.
- Barrera TL, Cummings JP, Armento M, Cully JA, Bush Amspoker A, Wilson NL, Mallen MJ, Shrestha S, Kunik ME, Stanley MA (2017). Telephone-delivered cognitive-behavioral therapy for older, rural veterans with depression and anxiety in home-based primary care. *Clinical gerontologist* 40, 114–123. [PubMed: 28452676]
- Brugha TS, Bebbington PE, Jenkins R (1999). A difference that matters: comparisons of structured and semi-structured psychiatric diagnostic interviews in the general population. *Psychological Medicine* 29, 1013–1020. [PubMed: 10576294]

- Canadian Agency for Drugs and Technologies in Health. (2016). PRESS – Peer review of electronic search strategies: 2015 Guideline explanation and elaboration (PRESS E&E) CADTH: Ottawa.
- Carter G, Spittal MJ (2018). Suicide risk assessment: Risk stratification is not accurate enough to be clinically useful and alternative approaches are needed. *Crisis* 39, 229. [PubMed: 29972324]
- Corson K, Gerrity MS, Dobscha SK (2004). Screening for depression and suicidality in a VA primary care setting: 2 items are better than 1 item. *The American Journal of Managed Care* 10, 839–845. [PubMed: 15609737]
- Crawford MJ, Thana L, Methuen C, Ghosh P, Stanley SV, Ross J, Gordon F, Blair G, Bajaj P (2011). Impact of screening for risk of suicide: randomised controlled trial. *British Journal of Psychiatry* 198, 379–384. [PubMed: 21525521]
- Fagerland MW, Lydersen S, Laake P (2014). Recommended tests and confidence intervals for paired binomial proportions. *Statistics in Medicine* 33, 2850–2875. [PubMed: 24648355]
- Higgins JP, Thompson SG (2002). Quantifying heterogeneity in a meta-analysis. *Statistics in Medicine* 21, 1539–1558. [PubMed: 12111919]
- Ishihara M, Harel D, Levis B, Levis AW, Riehm KE, Saadat N, Azar M, Rice DB, Sanchez TA, Chiovitti MJ, Cuijpers P (2018). Shortening self-report mental health symptom measures through optimal test assembly methods: Development and validation of the Patient Health Questionnaire-Depression-4. *Depression and anxiety* Advance online publication. doi: 10.1002/da.22841.
- Kroenke K, Spitzer RL, Williams JBW (2001). The PHQ-9: Validity of a brief depression severity measure. *Journal of General Internal Medicine* 16, 606–613. [PubMed: 11556941]
- Kroenke K, Strine TW, Spitzer RL, Williams JB, Berry JT, Mokdad AH (2009). The PHQ-8 as a measure of current depression in the general population. *Journal of Affective Disorders* 114, 163–173. [PubMed: 18752852]
- Levis B, Benedetti A, Riehm KE, Saadat N, Levis AW, Azar M, Rice DB, Chiovitti MJ, Sanchez TA, Cuijpers P, Gilbody S (2018). Probability of major depression diagnostic classification using semi-structured versus fully structured diagnostic interviews. *British Journal of Psychiatry* 212, 377–385. [PubMed: 29717691]
- Levis B, Benedetti A, Thombs BD, on behalf of the DEPRESSion Screening Data (DEPRESSD) Collaboration (2019). Accuracy of Patient Health Questionnaire-9 (PHQ-9) for screening to detect major depression: individual participant data meta-analysis. *BMJ* 2019;365:11476. [PubMed: 30967483]
- Moriarty AS, 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. [PubMed: 26195347]
- Nosen E, Woody SR (2008). Chapter 8: Diagnostic assessment in research. In *Handbook of Research Methods in Abnormal and Clinical Psychology* (ed. McKay D). Sage: Thousand Oaks.
- Razykov I, Ziegelstein R, Whooley M, Thombs BD (2012). The PHQ-9 versus the PHQ-8 – Is item 9 useful for assessing suicide risk in coronary artery disease patients? Data from the heart and Soul study. *Journal of Psychosomatic Research* 73, 163–168. [PubMed: 22850254]
- Rice DB, Kloda LA, Levis B, Qi B, Kingsland E, Thombs BD (2016). Are MEDLINE searches sufficient for systematic reviews and meta-analyses of the diagnostic accuracy of depression screening tools? A review of meta-analyses. *Journal of Psychosomatic Research* 87, 7–13. [PubMed: 27411746]
- Riley RD, Dodd SR, Craig JV, Thompson JR, Williamson PR (2008). Meta-analysis of diagnostic test studies using individual patient data and aggregate data. *Statistics in Medicine* 27, 6111–6136. [PubMed: 18816508]
- Robins LN, Wing J, Wittchen HU, Helzer JE, Babor TF, Burke J, Farmer A, Jablenski A, Pickens R, Regier DA, Sartorius N (1988). The Composite International Diagnostic Interview: an epidemiologic instrument suitable for use in conjunction with different diagnostic systems and in different cultures. *Archives of General Psychiatry* 45, 1069–1077. [PubMed: 2848472]
- Sheehan DV, Lecrubier Y, Sheehan KH, Janavs J, Weiller E, Keskiner A, Schinka J, Knapp E, Sheehan MF, Dunbar GC (1997). The validity of the Mini International Neuropsychiatric Interview (MINI) according to the SCID-P and its reliability. *European Psychiatry* 12, 232–241.

- Simon GE, Coleman KJ, Rossom RC, Beck A, Oliver M, Johnson E, Whiteside U, Operskalski B, Penfold RB, Shortreed SM, Rutter C (2016). Risk of suicide attempt and suicide death following completion of the Patient Health Questionnaire depression module in community practice. *The Journal of Clinical Psychiatry* 77, 221–227. [PubMed: 26930521]
- Simon GE, Rutter CM, Peterson D, Oliver M, Whiteside U, Operskalski B, Ludman EJ (2013). Does response on the PHQ-9 Depression Questionnaire predict subsequent suicide attempt or suicide death?. *Psychiatric Services* 64, 1195–1202. [PubMed: 24036589]
- Suarez L, Beach SR, Moore SV, Mastromauro CA, Januzzi JL, Celano CM, Chang TE, Huffman JC (2015). Use of the Patient Health Questionnaire-9 and a detailed suicide evaluation in determining imminent suicidality in distressed patients with cardiac disease. *Psychosomatics* 56, 181–189. [PubMed: 25660436]
- Siu AL, and the US Preventive Services Task Force (USPSTF) (2016). Screening for Depression in Adults: US Preventive Services Task Force Recommendation Statement. *JAMA: The Journal of the American Medical Association* 315, 380–387. [PubMed: 26813211]
- Thombs BD, Benedetti A, Kloda LA, Levis B, Nicolau I, Cuijpers P, Gilbody S, Ioannidis JP, McMillan D, Patten SB, Shrier I (2014). The diagnostic accuracy of the Patient Health Questionnaire-2 (PHQ-2), Patient Health Questionnaire-8 (PHQ-8), and Patient Health Questionnaire-9 (PHQ-9) for detecting major depression: protocol for a systematic review and individual patient data meta-analyses. *Systematic Reviews* 3, 124. [PubMed: 25348422]
- Walker J, Hansen CH, Butcher I, Sharma N, Wall L, Murray G, Sharpe M (2011). Thoughts of death and suicide reported by cancer patients who endorsed the “suicidal thoughts” item of the PHQ-9 during routine screening for depression. *Psychosomatics* 52, 424–427. [PubMed: 21907060]
- Walker E, Nowacki AS (2011). Understanding equivalence and noninferiority testing. *Journal of General Internal Medicine* 26, 192–196. [PubMed: 20857339]
- Wells TS, Horton JL, LeardMann CA, Jacobson IG, Boyko EJ (2013). A comparison of the PRIME-MD PHQ-9 and PHQ-8 in a large military prospective study, the Millennium Cohort Study. *Journal of Affective Disorders* 148, 77–83. [PubMed: 23246365]
- Whiting PF, Rutjes AW, Westwood ME, Mallett S, Deeks JJ, Reitsma JB, Leeflang MM, Sterne JA, Bossuyt PM (2011). QUADAS-2: a revised tool for the quality assessment of diagnostic accuracy studies. *Annals of Internal Medicine* 155, 529–536. [PubMed: 22007046]
- van der Leeden R, Busing FMTA, Meijer E. (1997). Bootstrap methods for two-level models. Technical report PRM 97–04 Leiden University, Department of Psychology: Leiden, The Netherlands.
- van der Leeden R, Meijer E, Busing FMTA (2008). Chapter 11: Resampling multilevel models. In *Handbook of Research Methods in Abnormal and Clinical Psychology* (McKay D). In *Handbook of Multilevel Analysis* (ed. Leeuw J, Meijer E), pp. 401–433. Springer: New York.

**Fig 1.**

(a) ROC curves for PHQ-8 and PHQ-9 among studies that used a semi-structured reference standard. (b) ROC curves for PHQ-8 and PHQ-9 among studies that used a fully structured reference standard (MINI excluded). (c) ROC curves for PHQ-8 and PHQ-9 among studies that used the MINI reference standard

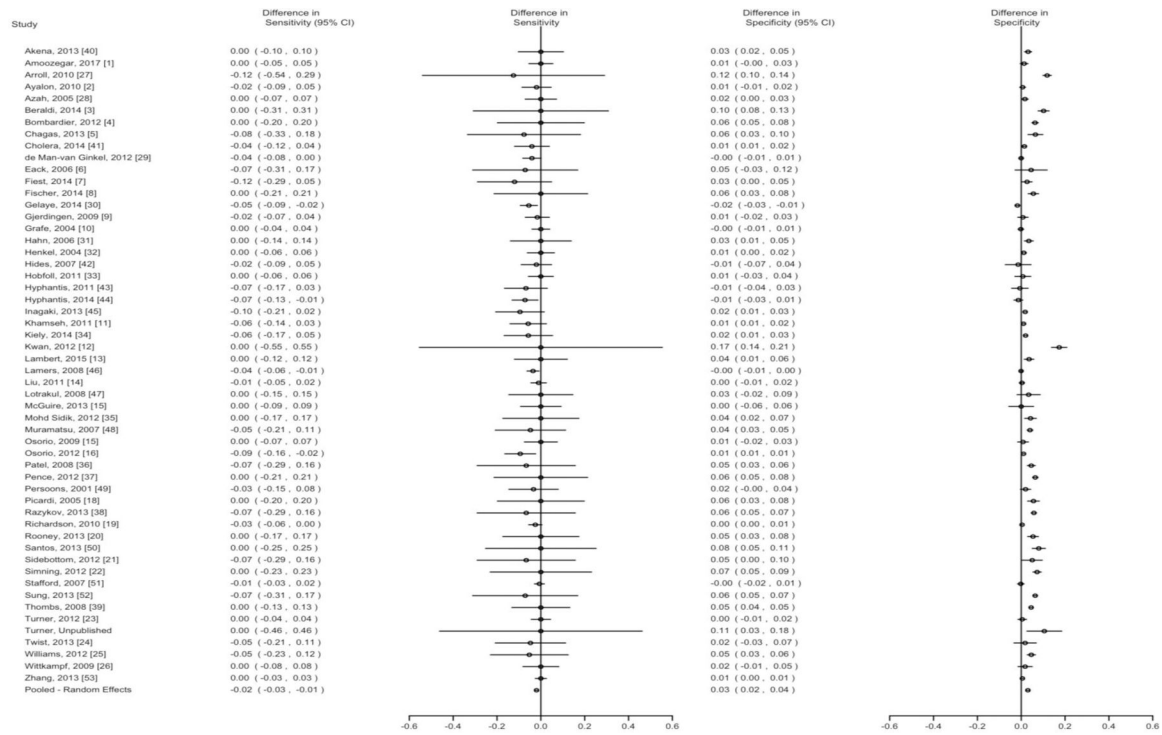


Fig 2. Forest plots of the difference in sensitivity and specificity estimates at cutoff 10 between PHQ-8 and PHQ-9 among all studies^a (N Studies = 54^b; N Participants = 16,742; N major depression = 2,097)^c
^a τ^2 for the difference of sensitivity and specificity were both <0.001.
^b The reference numbers refer to Supplementary Material References.
^c Amoozegar, 2017 [1] and Lambert, 2015 [13] were unpublished at the time of the electronic database search then subsequently published.

Table 1.

Participant characteristics by subgroup

Participant Subgroup	N Participants	N (%) Major Depression
All participants	16,742	2,097 (13)
Type of diagnostic interview		
Semi-structured diagnostic interview	6,362	790 (12)
Fully structured diagnostic interview	7,596	790 (10)
Mini International Neuropsychiatric Interview	2,784	517 (19)
Age^a		
< 60	11,144	1,402 (13)
60	5,552	692 (12)
Sex^a		
Women	9,552	1,259 (13)
Men	7,180	835 (12)
Care setting		
Non-medical care	1,832	252 (14)
Primary care	7,846	760 (10)
Inpatient specialty care	1,245	136 (11)
Outpatient specialty care	5,819	949 (16)
Country human development index		
Very high	13,297	1,577 (12)
High	1,337	276 (21)
Low-medium	2,108	244 (12)

^aDue to missing participant data, total participant numbers for these variables were <16,742.

Table 2.

Characteristics of participants who rated Item 9 as present for several days, more than half the days, or nearly every day (i.e., scores 1–3) in last two weeks by total PHQ-8 score

Total PHQ-8 Score	N of participants ^a	% with non-zero Item 9	Item 9 Mean (SD)
0–4	11,034	1.9%	0.02 (0.16)
5–9	5,071	13.2%	0.16 (0.45)
10–14	2,231	31.3%	0.44 (0.74)
15–19	1,044	48.3%	0.78 (0.97)
20–24	380	64.7%	1.39 (1.25)

^aNumbers of participants add up to >16,742 as they were weighted by sampling weights.

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 3a.

Comparison of sensitivity and specificity estimates between PHQ-8 and PHQ-9 among studies that used a semi-structured reference standard

Cutoff	PHQ-8 ^a				PHQ-9				PHQ-8 – PHQ-9 ^b			
	Sensitivity	95% CI	Specificity	95% CI	Sensitivity	95% CI	Specificity	95% CI	Sensitivity	95% CI	Specificity	95% CI
5	0.98	(0.95, 0.99)	0.55	(0.50, 0.60)	0.98	(0.95, 0.99)	0.55	(0.50, 0.60)	0.00	(-0.01, 0.00)	0.00	(0.00, 0.01)
6	0.98	(0.95, 0.99)	0.63	(0.58, 0.68)	0.98	(0.95, 0.99)	0.63	(0.58, 0.67)	0.00	(-0.00, 0.00)	0.00	(0.00, 0.01)
7	0.97	(0.93, 0.99)	0.70	(0.66, 0.74)	0.98	(0.93, 0.99)	0.70	(0.65, 0.74)	-0.01	(-0.02, 0.00)	0.00	(0.00, 0.01)
8	0.94	(0.89, 0.96)	0.76	(0.72, 0.79)	0.95	(0.90, 0.97)	0.75	(0.71, 0.79)	-0.01	(-0.03, 0.00)	0.01	(0.00, 0.01)
9	0.89	(0.84, 0.92)	0.81	(0.78, 0.84)	0.91	(0.87, 0.95)	0.80	(0.77, 0.83)	-0.02	(-0.06, -0.00)	0.01	(0.00, 0.01)
10 ^b	0.86	(0.80, 0.90)	0.86	(0.83, 0.89)	0.88	(0.82, 0.92)	0.86	(0.82, 0.88)	-0.02	(-0.06, -0.00)	0.00	(0.00, 0.02)
11	0.81	(0.75, 0.87)	0.90	(0.87, 0.92)	0.84	(0.84, 0.84)	0.89	(0.89, 0.89)	-0.03	(-0.06, -0.00)	0.01	(0.00, 0.02)
12	0.74	(0.68, 0.79)	0.92	(0.89, 0.93)	0.78	(0.71, 0.83)	0.91	(0.89, 0.93)	-0.04	(-0.09, -0.01)	0.01	(0.00, 0.01)
13	0.67	(0.60, 0.73)	0.94	(0.92, 0.95)	0.69	(0.63, 0.75)	0.93	(0.91, 0.95)	-0.02	(-0.07, -0.00)	0.01	(0.00, 0.01)
14	0.59	(0.53, 0.65)	0.96	(0.94, 0.97)	0.64	(0.57, 0.70)	0.95	(0.93, 0.96)	-0.05	(-0.09, -0.01)	0.01	(0.00, 0.01)
15	0.51	(0.44, 0.57)	0.97	(0.95, 0.98)	0.55	(0.48, 0.62)	0.96	(0.94, 0.97)	-0.04	(-0.09, -0.02)	0.01	(0.00, 0.01)

^aN Studies = 27; N Participants = 6,362; N major depression = 790

^bFor PHQ-8 cutoff 10, among studies that used semi-structured interviews as the reference standard, the default optimizer in gimer failed to converge, thus bobyqa was used instead.

CI: confidence interval

Table 3b. Comparison of sensitivity and specificity estimates between PHQ-8 and PHQ-9 among studies that used a fully structured reference standard (MINI excluded)

Cutoff	PHQ-8 ^a				PHQ-9				PHQ-8 – PHQ-9			
	Sensitivity	95% CI	Specificity	95% CI	Sensitivity	95% CI	Specificity	95% CI	Sensitivity	95% CI	Specificity	95% CI
5	0.92	(0.85, 0.96)	0.57	(0.49, 0.66)	0.93	(0.85, 0.96)	0.57	(0.48, 0.65)	-0.01	(-0.03, 0.00)	0.00	(0.00, 0.02)
6	0.88	(0.79, 0.93)	0.65	(0.57, 0.73)	0.90	(0.81, 0.95)	0.65	(0.56, 0.72)	-0.02	(-0.07, 0.00)	0.00	(0.00, 0.02)
7	0.83	(0.73, 0.90)	0.72	(0.64, 0.79)	0.84	(0.73, 0.90)	0.71	(0.63, 0.78)	-0.01	(-0.01, 0.00)	0.01	(0.00, 0.02)
8	0.77	(0.66, 0.85)	0.78	(0.71, 0.84)	0.79	(0.68, 0.86)	0.77	(0.70, 0.83)	-0.02	(-0.07, -0.00)	0.01	(0.00, 0.01)
9	0.69	(0.59, 0.77)	0.83	(0.76, 0.87)	0.71	(0.62, 0.80)	0.81	(0.75, 0.86)	-0.02	(-0.07, -0.00)	0.02	(0.01, 0.02)
10	0.63	(0.52, 0.72)	0.86	(0.81, 0.90)	0.67	(0.57, 0.76)	0.85	(0.80, 0.90)	-0.04	(-0.09, -0.01)	0.01	(0.00, 0.02)
11	0.57	(0.45, 0.67)	0.89	(0.85, 0.93)	0.59	(0.49, 0.69)	0.88	(0.84, 0.92)	-0.02	(-0.07, -0.01)	0.01	(0.00, 0.02)
12	0.51	(0.38, 0.64)	0.92	(0.88, 0.94)	0.54	(0.43, 0.65)	0.90	(0.86, 0.93)	-0.03	(-0.16, -0.01)	0.02	(0.01, 0.02)
13	0.43	(0.32, 0.55)	0.94	(0.91, 0.96)	0.47	(0.36, 0.58)	0.93	(0.89, 0.95)	-0.04	(-0.12, -0.01)	0.01	(0.00, 0.01)
14	0.36	(0.26, 0.47)	0.95	(0.93, 0.97)	0.41	(0.31, 0.53)	0.95	(0.92, 0.96)	-0.05	(-0.14, -0.01)	0.00	(0.00, 0.01)
15	0.30	(0.22, 0.39)	0.96	(0.95, 0.98)	0.33	(0.24, 0.42)	0.96	(0.94, 0.97)	-0.03	(-0.07, -0.00)	0.00	(0.00, 0.00)

^a N Studies = 13; N Participants = 7,596; N major depression = 790

CI: confidence interval

Table 3c.

Comparison of sensitivity and specificity estimates between PHQ-8 and PHQ-9 among studies that used the MINI reference standard

Cutoff	PHQ-8 ^a				PHQ-9				PHQ-8 – PHQ-9			
	Sensitivity	95% CI	Specificity	95% CI	Sensitivity	95% CI	Specificity	95% CI	Sensitivity	95% CI	Specificity	95% CI
5	0.96	(0.93, 0.98)	0.58	(0.50, 0.65)	0.97	(0.93, 0.98)	0.57	(0.49, 0.65)	-0.01	(-0.01, 0.00)	0.01	(0.00, 0.02)
6	0.92	(0.85, 0.96)	0.67	(0.59, 0.74)	0.93	(0.86, 0.97)	0.66	(0.59, 0.73)	-0.01	(-0.03, 0.00)	0.01	(0.00, 0.02)
7	0.89	(0.81, 0.94)	0.73	(0.67, 0.79)	0.89	(0.81, 0.94)	0.73	(0.66, 0.79)	0.00	(-0.03, 0.00)	0.01	(0.00, 0.01)
8	0.83	(0.75, 0.89)	0.80	(0.75, 0.84)	0.86	(0.77, 0.91)	0.79	(0.74, 0.83)	-0.03	(-0.06, 0.00)	0.01	(0.00, 0.02)
9	0.78	(0.69, 0.85)	0.85	(0.81, 0.89)	0.81	(0.71, 0.88)	0.84	(0.80, 0.88)	-0.03	(-0.07, -0.00)	0.01	(0.00, 0.02)
10	0.72	(0.63, 0.79)	0.88	(0.84, 0.91)	0.75	(0.66, 0.82)	0.88	(0.84, 0.91)	-0.03	(-0.08, -0.01)	0.01	(0.00, 0.02)
11	0.65	(0.57, 0.73)	0.91	(0.88, 0.94)	0.69	(0.61, 0.77)	0.90	(0.87, 0.93)	-0.04	(-0.08, -0.01)	0.01	(0.00, 0.02)
12	0.59	(0.51, 0.66)	0.93	(0.91, 0.95)	0.65	(0.56, 0.73)	0.92	(0.90, 0.94)	-0.06	(-0.11, -0.02)	0.01	(0.00, 0.02)
13	0.53	(0.44, 0.62)	0.95	(0.93, 0.97)	0.57	(0.49, 0.66)	0.94	(0.92, 0.96)	-0.04	(-0.09, -0.01)	0.01	(0.00, 0.02)
14	0.43	(0.35, 0.51)	0.97	(0.95, 0.98)	0.49	(0.49, 0.49)	0.96	(0.96, 0.96)	-0.06	(-0.11, -0.02)	0.01	(0.00, 0.02)
15	0.37	(0.29, 0.45)	0.98	(0.96, 0.99)	0.42	(0.42, 0.42)	0.97	(0.97, 0.97)	-0.05	(-0.10, -0.02)	0.01	(0.00, 0.01)

^aN Studies = 14; N Participants = 2,784; N major depression = 517

CI: confidence interval; MINI: Mini International Neuropsychiatric Interview

Table 4.

Comparison of sensitivity and specificity estimates between PHQ-8 and PHQ-9 across cutoffs 5–15 for all studies

PHQ-8 – PHQ-9				
Cutoff	Sensitivity	95% CI	Specificity	95% CI
5	–0.01	(–0.01, 0.00)	0.00	(0.00, 0.01)
6	0.00	(–0.01, 0.00)	0.01	(0.00, 0.01)
7	–0.01	(–0.02, 0.00)	0.00	(0.00, 0.01)
8	–0.01	(–0.04, –0.01)	0.00	(0.01, 0.01)
9	–0.03	(–0.06, –0.01)	0.01	(0.01, 0.01)
10	–0.03	(–0.06, –0.02)	0.01	(0.00, 0.01)
11	–0.03	(–0.06, –0.01)	0.01	(0.01, 0.01)
12	–0.05	(–0.08, –0.03)	0.01	(0.00, 0.01)
13	–0.04	(–0.06, –0.02)	0.00	(0.00, 0.01)
14	–0.05	(–0.08, –0.03)	0.01	(0.00, 0.01)
15	–0.04	(–0.07, –0.03)	0.01	(0.00, 0.01)

^aN Studies = 54; N Participants = 16,742; N major depression = 2,097

CI: confidence interval