

ORIGINAL ARTICLE

Patient Health Questionnaire-9 scores do not accurately estimate depression prevalence: individual participant data meta-analysis

Brooke Levis^{a,b}, Andrea Benedetti^{b,c,d}, John P.A. Ioannidis^{e,f,g,h}, Ying Sun^a, Zelalem Negeri^{a,b}, Chen He^a, Yin Wu^{a,b,i}, Ankur Krishnan^a, Parash Mani Bhandari^{a,b}, Dipika Neupane^{a,b}, Mahrukh Imran^a, Danielle B. Rice^{a,j}, Kira E. Riehm^{a,k}, Nazanin Saadat^a, Marleine Azar^{a,b}, Jill Boruff^l, Pim Cuijpers^m, Simon Gilbodyⁿ, Lorie A. Kloda^o, Dean McMillan^k, Scott B. Patten^{p,q}, Ian Shrier^{a,b,r}, Roy C. Ziegelstein^s, Sultan H. Alamri^t, Dagmar Amtmann^u, Liat Ayalon^v, Hamid R. Baradaran^{w,x}, Anna Beraldi^y, Charles N. Bernstein^{z,aa}, Arvin Bhana^{bb,cc}, Charles H. Bombardier^u, Gregory Carter^{dd}, Marcos H. Chagas^{ee}, Dixon Chibanda^{ff}, Kerrie Clover^{dd}, Yeates Conwell^{gg}, Crisanto Diez-Quevedo^{hh,ii}, Jesse R. Fann^{jj}, Felix H. Fischer^{i,kk}, Leila Gholizadeh^{ll}, Lorna J. Gibson^{mm}, Eric P. Greenⁿⁿ, Catherine G. Greeno^{oo}, Brian J. Hall^{pp,qq}, Emily E. Haroz^{rr}, Khalida Ismail^{ss}, Nathalie Jetté^{p,q,tt}, Mohammad E. Khamseh^w, Yunxin Kwan^{uu}, Maria Asunción Lara^{vv}, Shen-Ing Liu^{ww,xx,yy,zz}, Sonia R. Loureiro^{ee}, Bernd Löwe^{aaa}, Ruth Ann Marrie^{bbb}, Laura Marsh^{ccc}, Anthony McGuire^{ddd}, Kumiko Muramatsu^{eee}, Laura Navarrete^{fff}, Flávia L. Osório^{ee,ggg}, Inge Petersen^{hhh}, Angelo Picardiⁱⁱⁱ, Stephanie L. Pugh^{jjj,kkk}, Terence J. Quinn^{lll}, Alasdair G. Rooney^{mmm}, Eileen H. Shinnⁿⁿⁿ, Abbey Sidebottom^{ooo}, Lena Spangenberg^{ppp}, Pei Lin Lynnette Tan^{uuu}, Martin Taylor-Rowan^{qqq}, Alyna Turner^{rrr,sss}, Henk C. van Weert^{ttt}, Paul A. Vöhringer^{uuu,vvv,www}, Lynne I. Wagner^{xxx,yyy}, Jennifer White^{zzz}, Kirsty Winkley^{aaaa}, Brett D. Thombs^{a,b,d,i,j,bbbb,cccc,*}

^aLady Davis Institute for Medical Research, Jewish General Hospital, Montréal, Québec, Canada

^bDepartment of Epidemiology, Biostatistics and Occupational Health, McGill University, Montréal, Québec, Canada

^cRespiratory Epidemiology and Clinical Research Unit, McGill University Health Centre, Montréal, Québec, Canada

^dDepartment of Medicine, McGill University, Montréal, Québec, Canada

^eStanford Prevention Research Center, Department of Medicine, Stanford University School of Medicine, Stanford, CA, USA

^fDepartment of Health Research and Policy, Stanford University School of Medicine, Stanford, CA, USA

^gDepartment of Biomedical Data Science, Stanford University School of Medicine, Stanford, CA, USA

^hDepartment of Statistics, Stanford University School of Humanities and Sciences, Stanford, CA, USA

ⁱDepartment of Psychiatry, McGill University, Montréal, Québec, Canada

^jDepartment of Psychology, McGill University, Montréal, Québec, Canada

^kDepartment of Mental Health, Bloomberg School of Public Health, Johns Hopkins University, Baltimore, MD, USA

^lSchulich Library of Physical Sciences, Life Sciences, and Engineering, McGill University, Montréal, Québec, Canada

^mDepartment of Clinical, Neuro and Developmental Psychology, Amsterdam Public Health Research Institute, Vrije Universiteit Amsterdam, Amsterdam, the Netherlands

ⁿHull York Medical School and the Department of Health Sciences, University of York, Heslington, NY, UK

^oLibrary, Concordia University, Montréal, Québec, Canada

^pDepartment of Community Health Sciences, University of Calgary, Calgary, Alberta, Canada

^qHotchkiss Brain Institute and O'Brien Institute for Public Health, University of Calgary, Calgary, Alberta, Canada

^rDepartment of Family Medicine, McGill University, Montréal, Québec, Canada

^sDepartment of Medicine, Johns Hopkins University School of Medicine, Baltimore, MD, USA

^tFaculty of Medicine, King Abdulaziz University, Jeddah, Makkah, Saudi Arabia

^uDepartment of Rehabilitation Medicine, University of Washington, Seattle, WA, USA

^vLouis and Gabi Weisfeld School of Social Work, Bar Ilan University, Ramat Gan, Israel

^wEndocrine Research Center, Institute of Endocrinology and Metabolism, Iran University of Medical Sciences, Tehran, Iran

* Corresponding author. Lady Davis Institute for Medical Research, Jewish General Hospital, Montréal, Québec, Canada. Tel.: +1 (514) 340-8222x25112.

E-mail address: brett.thombs@mcgill.ca (B.D. Thombs).

- ^xAgeing Clinical & Experimental Research Team, Institute of Applied Health Sciences, School of Medicine, Medical Sciences and Nutrition, University of Aberdeen, Aberdeen, Scotland, UK
- ^yKbo-Lech-Mangfall-Klinik Garmisch-Partenkirchen, Klinik für Psychiatrie, Psychotherapie & Psychosomatik, Lehrkrankenhaus der Technischen Universität München, Munich, Germany
- ^zUniversity of Manitoba IBD Clinical and Research Centre, Winnipeg, Manitoba, Canada
- ^{aa}Department of Internal Medicine, Max Rady College of Medicine, Rady Faculty of Health Sciences, University of Manitoba, Winnipeg, Manitoba, Canada
- ^{bb}Centre for Rural Health, School of Nursing and Public Health, College of Health Sciences, University of KwaZulu-Natal, Durban, KwaZulu-Natal, South Africa
- ^{cc}Health Systems Research Unit, South African Medical Research Council, Cape Town, South Africa
- ^{dd}Centre for Brain and Mental Health Research, University of Newcastle, New South Wales, Australia
- ^{ee}Department of Neurosciences and Behavior, Ribeirão Preto Medical School, University of São Paulo, Ribeirão Preto, Brazil
- ^{ff}Department of Community Medicine, University of Zimbabwe, Harare, Zimbabwe
- ^{gg}Department of Psychiatry, University of Rochester Medical Center, Rochester, NY, USA
- ^{hh}Servei de Psiquiatria, Hospital Germans Trias i Pujol, Badalona, Spain
- ⁱⁱDepartament de Psiquiatria i Medicina Legal, Universitat Autònoma de Barcelona, Badalona, Spain
- ^{jj}Department of Psychiatry and Behavioral Sciences, University of Washington, Seattle, WA, USA
- ^{kk}Department of Psychosomatic Medicine, Center for Internal Medicine and Dermatology, Charité - Universitätsmedizin Berlin, Berlin, Germany
- ^{ll}Faculty of Health, University of Technology Sydney, Sydney, Australia
- ^{mm}Tropical Epidemiology Group, Faculty of Epidemiology and Population Health, London School of Hygiene and Tropical Medicine, London, UK
- ⁿⁿDuke Global Health Institute, Duke University, Durham, NC, USA
- ^{oo}School of Social Work, University of Pittsburgh, Pittsburgh, PA, USA
- ^{pp}Department of Psychology, Faculty of Social Sciences, Global and Community Mental Health Research Group, University of Macau, Macau Special Administrative Region, China
- ^{qq}Department of Health, Behavior, and Society, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA
- ^{rr}Center for American Indian Health, Department of International Health, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA
- ^{ss}Department of Psychological Medicine, Institute of Psychiatry, Psychology and Neurosciences, King's College London Weston Education Centre, London, UK
- ^{tt}Department of Neurology, Icahn School of Medicine at Mount Sinai, New York, NY, USA
- ^{uu}Department of Psychological Medicine, Tan Tock Seng Hospital, Singapore
- ^{vv}Instituto Nacional de Psiquiatria Ramón de la Fuente Muñiz, San Lorenzo Huipulco, Tlalpan, México D. F. Mexico
- ^{ww}Programme in Health Services & Systems Research, Duke-NUS Medical School, Singapore
- ^{xx}Department of Psychiatry, Mackay Memorial Hospital, Taipei, Taiwan
- ^{yy}Department of Medical Research, Mackay Memorial Hospital, Taipei, Taiwan
- ^{zz}Department of Medicine, Mackay Medical College, Taipei, Taiwan
- ^{aaa}Department of Psychosomatic Medicine and Psychotherapy, University Medical Center Hamburg-Eppendorf, Hamburg, Germany
- ^{bbb}Departments of Medicine and Community Health Sciences, Max Rady College of Medicine, Rady Faculty of Health Sciences, University of Manitoba, Winnipeg, Manitoba, Canada
- ^{ccc}Baylor College of Medicine, Houston and Michael E. DeBakey Veterans Affairs Medical Center, Houston, TX, USA
- ^{ddd}Department of Nursing, St. Joseph's College, Standish, ME, USA
- ^{eee}Department of Clinical Psychology, Graduate School of Niigata Seiryō University, Niigata, Japan
- ^{fff}Department of Epidemiology and Psychosocial Research, Instituto Nacional de Psiquiatria Ramón de la Fuente Muñiz, Ciudad de México, Mexico
- ^{ggg}National Institute of Science and Technology, Translational Medicine, Ribeirão Preto, Brazil
- ^{hhh}Centre for Rural Health, School of Nursing and Public Health, University of KwaZulu-Natal, Durban, South Africa
- ⁱⁱⁱCentre for Behavioural Sciences and Mental Health, Italian National Institute of Health, Rome, Italy
- ^{jjj}NRG Oncology Statistics and Data Management Center, Philadelphia, PA, USA
- ^{kkk}American College of Radiology, Philadelphia, PA, USA
- ^{lll}Institute of Cardiovascular & Medical Sciences, University of Glasgow, Glasgow, Scotland, UK
- ^{mmm}Division of Psychiatry, Royal Edinburgh Hospital, University of Edinburgh, Edinburgh, Scotland, UK
- ⁿⁿⁿDepartment of Behavioral Science, University of Texas M. D. Anderson Cancer Center, Houston, TX, USA
- ^{ooo}Allina Health, Minneapolis, MN, USA
- ^{ppp}Department of Medical Psychology and Medical Sociology, University of Leipzig, Leipzig, Germany
- ^{qqq}Institute of Cardiovascular and Medical Science, University of Glasgow, Glasgow, Scotland, UK
- ^{rrr}School of Medicine and Public Health, University of Newcastle, Newcastle, New South Wales, Australia
- ^{sss}Deakin University, IMPACT Strategic Research Centre, School of Medicine, Barwon Health, Geelong, Victoria, Australia
- ^{ttt}Department of General Practice, Amsterdam Institute for General Practice and Public Health, Amsterdam University Medical Centers, Location AMC, Amsterdam, the Netherlands
- ^{uuu}Department of Psychiatry and Mental Health, Clinical Hospital, Universidad de Chile, Santiago, Chile
- ^{vvv}Millennium Institute for Depression and Personality Research (MIDAP), Ministry of Economy, Macul, Santiago, Chile
- ^{www}Psychiatry Department, Tufts Medical Center, Tufts University, Boston, MA, USA
- ^{xxx}Department of Social Sciences and Health Policy, Wake Forest School of Medicine, Wake Forest University, Winston-Salem, NC, USA
- ^{yyy}Wake Forest Baptist Comprehensive Cancer Center, Winston-Salem, NC, USA
- ^{zzz}Department of Physiotherapy, School of Primary and Allied Health Care, Monash University, Melbourne, Australia
- ^{aaaa}Florence Nightingale Faculty of Nursing, Midwifery & Palliative Care, King's College London, London, UK
- ^{bbbb}Department of Educational and Counselling Psychology, McGill University, Montréal, Québec, Canada
- ^{cccc}Biomedical Ethics Unit, McGill University, Montréal, Québec, Canada

Abstract

Objectives: Depression symptom questionnaires are not for diagnostic classification. Patient Health Questionnaire-9 (PHQ-9) scores ≥ 10 are nonetheless often used to estimate depression prevalence. We compared PHQ-9 ≥ 10 prevalence to Structured Clinical Interview for Diagnostic and Statistical Manual of Mental Disorders (SCID) major depression prevalence and assessed whether an alternative PHQ-9 cutoff could more accurately estimate prevalence.

Study Design and Setting: Individual participant data meta-analysis of datasets comparing PHQ-9 scores to SCID major depression status.

Results: A total of 9,242 participants (1,389 SCID major depression cases) from 44 primary studies were included. Pooled PHQ-9 ≥ 10 prevalence was 24.6% (95% confidence interval [CI]: 20.8%, 28.9%); pooled SCID major depression prevalence was 12.1% (95% CI: 9.6%, 15.2%); and pooled difference was 11.9% (95% CI: 9.3%, 14.6%). The mean study-level PHQ-9 ≥ 10 to SCID-based prevalence ratio was 2.5 times. PHQ-9 ≥ 14 and the PHQ-9 diagnostic algorithm provided prevalence closest to SCID major depression prevalence, but study-level prevalence differed from SCID-based prevalence by an average absolute difference of 4.8% for PHQ-9 ≥ 14 (95% prediction interval: -13.6%, 14.5%) and 5.6% for the PHQ-9 diagnostic algorithm (95% prediction interval: -16.4%, 15.0%).

Conclusion: PHQ-9 ≥ 10 substantially overestimates depression prevalence. There is too much heterogeneity to correct statistically in individual studies. © 2020 Elsevier Inc. All rights reserved.

Keywords: Depression prevalence; PHQ-9; SCID; Individual participant data meta-analysis

1. Introduction

Disease prevalence estimates have important implications for interpreting medical research, understanding disease burden, and making decisions about health care resource utilization [1]. In mental health research, major depression classification requires using validated diagnostic interviews [2,3]. Administering diagnostic interviews in large enough samples to estimate prevalence, however, is resource intensive. Thus, researchers sometimes use self-report depression symptom questionnaires, or screening tools, instead, and label the percentage of participants scoring above a screening cutoff as depression prevalence [4,5]. A 2018 study identified 19 primary studies listed in PubMed in a 3-month period whose titles indicated that they assessed the prevalence of depression or depressive disorders and found that 89% were based on screening questionnaires only [4].

Some self-report questionnaires include the same symptoms evaluated in validated diagnostic interviews. None, however, include all components of diagnostic interviews, such as assessment of functional impairment or investigation of nonpsychiatric medical conditions that can cause similar symptoms [4]. Using depression symptom questionnaires and cutoffs intended for screening to assess depression prevalence may overestimate prevalence. This is because screening attempts to identify previously unrecognized cases; cutoffs are set to cast a wide net and identify many more patients who may have depression than meet diagnostic criteria.

A recent review examined meta-analyses of depression prevalence published in 2008–2017 [5]. Of 81 prevalence estimates reported in abstracts of 69 meta-analyses, 10% were based on diagnostic interviews, 44% were based on screening or rating tools, and 46% combined results from diagnostic interviews and screening or rating tools. The mean reported prevalence was 31% among meta-analyses

based on screening or rating tools compared with 17% with diagnostic interviews [5]. The degree to which screening tools exaggerate prevalence, however, depends on the screening tool and cutoff used [4,5].

We do not know of any studies that have evaluated the degree to which specific screening tool and cutoff combinations overestimate depression prevalence [4,5]. The Patient Health Questionnaire-9 (PHQ-9) [6–8] is the most commonly used depression screening tool in primary care [9]. Its nine items align with the Diagnostic and Statistical Manual of Mental Disorders (DSM) criteria for major depressive episode (MDE) [10–12]. The standard cutoff, ≥ 10 , is well established for screening to detect major depression and maximized combined sensitivity and specificity in a recent individual participant data meta-analysis (IPDMA) [6–8,13]. PHQ-9 ≥ 10 has been used to estimate depression prevalence in primary research studies and via synthesis in meta-analyses, including in very high-impact journals [14–16]. It is also sometimes used to diagnose depression and make treatment decisions for individual patients [6,17–19].

Our objective was to use an IPDMA approach to (1) compare PHQ-9 ≥ 10 prevalence to major depression prevalence based on a well-validated semistructured diagnostic interview, the Structured Clinical Interview for DSM (SCID) [20]; and (2) use a prevalence matching approach [4,21] to determine if a PHQ-9 cutoff could be set to match SCID-based prevalence with sufficiently low heterogeneity to accurately estimate prevalence in individual studies.

2. Methods

This study used a subset of data accrued for an IPDMA of the accuracy of the PHQ-9 for screening to detect major depression [13]. Detailed methods were registered in

What is new?

Key findings

- Based on 9,242 participants from 44 primary studies, the pooled Patient Health Questionnaire-9 (PHQ-9) ≥ 10 prevalence (25%) was double the pooled Structured Clinical Interview for DSM (SCID) major depression prevalence (12%); the pooled difference from individual studies was 12%.
- PHQ-9 ≥ 14 prevalence and PHQ-9 diagnostic algorithm prevalence most closely matched SCID major depression prevalence, but study-level PHQ-9 ≥ 14 and PHQ-9 diagnostic algorithm prevalence differed from SCID major depression prevalence with 95% prediction intervals of -14% to 15% and -16% to 15% , respectively.

What this adds to what was known?

- Although PHQ-9 ≥ 10 is often used to estimate depression prevalence, it overestimates major depression prevalence substantially.
- There is too much heterogeneity to correct statistically in individual studies.

What is the implication and what should change now?

- Estimates of depression prevalence should be based on validated diagnostic interviews designed for determining case status; users should evaluate published reports of depression prevalence to ensure that they are based on methods intended to classify major depression.

PROSPERO (CRD42014010673), and a protocol was published [22]. This analysis was not part of the original IPDMA protocol.

2.1. Study selection

In the main IPDMA, datasets from articles in any language were eligible for inclusion if (1) they included PHQ-9 scores; (2) they included diagnostic classifications for current MDE or major depressive disorder (MDD) based on DSM [10–12] or International Classification of Diseases [23] criteria, using a validated semistructured or fully structured interview; (3) the PHQ-9 and diagnostic interview were administered within 2 weeks of each other; (4) participants were aged ≥ 18 years and not recruited from youth or school-based settings; and (5) participants were not recruited from psychiatric settings or because they were identified as having depressive symptoms. Datasets where

not all participants were eligible were included if primary data allowed the selection of eligible participants.

For the present study, we included primary studies that based diagnoses on the SCID [20]. The SCID is a semi-structured diagnostic interview intended to be conducted by an experienced diagnostician; it requires clinical judgment and allows rephrasing questions and probes. The reason for including only SCID studies is that in analyses using large IPDMA databases [24–26], we found that, compared with semistructured interviews, fully structured interviews, which are designed for administration by lay interviewers, identify more participants with low-level symptoms as depressed but fewer participants with high-level symptoms. These results were consistent with the idea that semistructured interviews most closely replicate clinical interviews done by trained professionals, whereas fully structured interviews are less resource-intensive options that can be administered by research staff without diagnostic skills but may misclassify major depression in many participants. In our PHQ-9 IPDMA database, 44 of 47 studies that used semistructured interviews used the SCID. Thus, to reduce heterogeneity, we only included these 44 studies in main analyses.

In sensitivity analyses, we also included the three studies that used other semistructured interviews. We considered also incorporating published results from eligible studies that did not contribute data to the IPDMA. However, only 3 of 14 such studies [27–29] (970 participants, including 77 major depression cases) reported sufficient information to compare PHQ-9 ≥ 10 and SCID-based prevalence, and these studies did not report information necessary to be included in all prevalence matching analyses.

2.2. Data sources and searches

A medical librarian searched MEDLINE, MEDLINE In-Process & Other Non-Indexed Citations via Ovid, PsycINFO, and Web of Science from January 1, 2000, to May 9, 2018, using a peer-reviewed [30] search strategy (Supplementary Material: Appendix Methods). We also reviewed reference lists of relevant reviews and queried contributing authors about nonpublished studies.

Two investigators independently reviewed titles and abstracts for eligibility. If either deemed a study potentially eligible, the full text was reviewed by two investigators, independently, with disagreements resolved by consensus, consulting a third investigator when necessary.

2.3. Data contribution and synthesis

Authors of eligible datasets were invited to contribute deidentified primary data, including PHQ-9 scores and major depression classification status. We emailed corresponding authors of eligible studies at least three times, as necessary. If no response, we emailed coauthors and attempted phone contact.

Before integrating individual datasets into our synthesized dataset, we compared published participant characteristics and diagnostic accuracy results with results from raw datasets and resolved discrepancies with the original investigators. When datasets included statistical weights to reflect sampling procedures, we used provided weights. For studies where sampling procedures merited weighting, but the original study did not weigh, we constructed weights using inverse selection probabilities.

2.4. Data analysis

2.4.1. Comparison of PHQ-9 ≥ 10 prevalence and SCID major depression prevalence

For each primary study, we estimated the percentage of participants who scored ≥ 10 on the PHQ-9, the percentage of participants classified as having major depression based on the SCID, the difference of these percentages, and the ratio. Then, across studies, we pooled prevalence for PHQ-9 ≥ 10 , prevalence for the SCID, and differences in prevalence.

2.4.2. Prevalence matching

To identify which PHQ-9 scoring approach best matched SCID-based prevalence, we estimated pooled differences in prevalence for each possible PHQ-9 cutoff and the PHQ-9 diagnostic algorithm compared with SCID. The scoring approach with the smallest pooled difference was chosen to be the “prevalence match scoring approach.” Then, for each included study, we estimated the difference and ratio in prevalence for the prevalence match scoring approach vs. SCID. We determined the mean and median absolute difference and range of differences across all studies. To illustrate the range of difference values that would be expected if a new study were to compare prevalence based on the prevalence match scoring approach to prevalence based on SCID, we estimated 95% prediction intervals for the differences. For the diagnostic algorithm, which requires five or more items with scores of ≥ 2 points, with at least one being depressed mood or anhedonia [8], three studies [31–33] (524 participants) and 88 additional participants from other studies (612 participants total, 7%) were excluded, as they did not provide PHQ-9 item scores, which are necessary to determine diagnostic algorithm criteria. In sensitivity analyses, we evaluated if results differed if the 612 participants were excluded from all analyses rather than just those involving the diagnostic algorithm.

All meta-analyses incorporated sampling weights and were conducted in R (R version 3.4.1; R Studio version 1.0.143) using the lme4 package. To estimate pooled prevalence values, generalized linear mixed-effects models with a logit link function were fit using the glmer function. To estimate pooled difference values, linear mixed-effects models were fit using the lmer function. To account for the correlation between subjects within the same primary study, random intercepts were fit for each primary study.

To quantify heterogeneity, we reported the estimated between-study variance (τ^2) for each analysis.

In post-hoc analyses, we investigated whether differences in prevalence for the PHQ-9 prevalence match scoring approach and SCID were associated with study and participant characteristics. To do this, we fit additional linear mixed-effects models for pooled prevalence difference, including age, sex, country human development index (“very high,” “high,” or “low-medium,” based on the United Nation’s 2018 Human Development Index) and recruitment setting category (primary care, nonmedical care, inpatient specialty care, or outpatient specialty care) as fixed-effect covariates. For these analyses, we excluded 56 participants (<1%) missing age or sex data.

3. Results

3.1. Search results and inclusion of primary study datasets

Of 9,674 unique titles and abstracts identified from the database search for the main IPDMA, 9,198 were excluded after title and abstract review and 297 were excluded after full-text review, leaving 179 eligible articles with data from 123 unique participant samples, of which 95 (77.2%) contributed datasets. Authors of included studies contributed data from five unpublished studies, for a total of 100 datasets. Of these, for the present study’s main analyses, we excluded 56 studies that classified major depression using a diagnostic interview other than the SCID (Fig. 1). Thus, the main analyses of the present study included 9,242 participants (1,389 major depression cases) from 44 primary studies [31–72]. Among the 28 eligible primary studies that did not provide datasets for the main IPDMA, 14 used the SCID (4,408 participants). Thus, the main analyses included 75.9% of eligible studies that used the SCID (44/58) and 67.7% of eligible participants (9,242/13,650). Table 1 shows the characteristics of each included study.

In sensitivity analyses, we included data from three additional studies (1,992 participants, including 139 major depression cases) that provided individual participant data but administered a semistructured interview other than the SCID (Table 1) [73–75].

3.2. Comparison of PHQ-9 ≥ 10 prevalence and SCID major depression prevalence

The percentage of participants with PHQ-9 ≥ 10 in each of the 44 SCID studies ranged from 5.3% to 64.8%; pooled prevalence was 24.6% (95% confidence interval [CI]: 20.8%, 28.9%; τ^2 : 0.505). The percentage of participants with SCID major depression ranged from 0.6% to 56.4%; pooled prevalence was 12.1% (95% CI: 9.6%, 15.2%; τ^2 : 0.703).

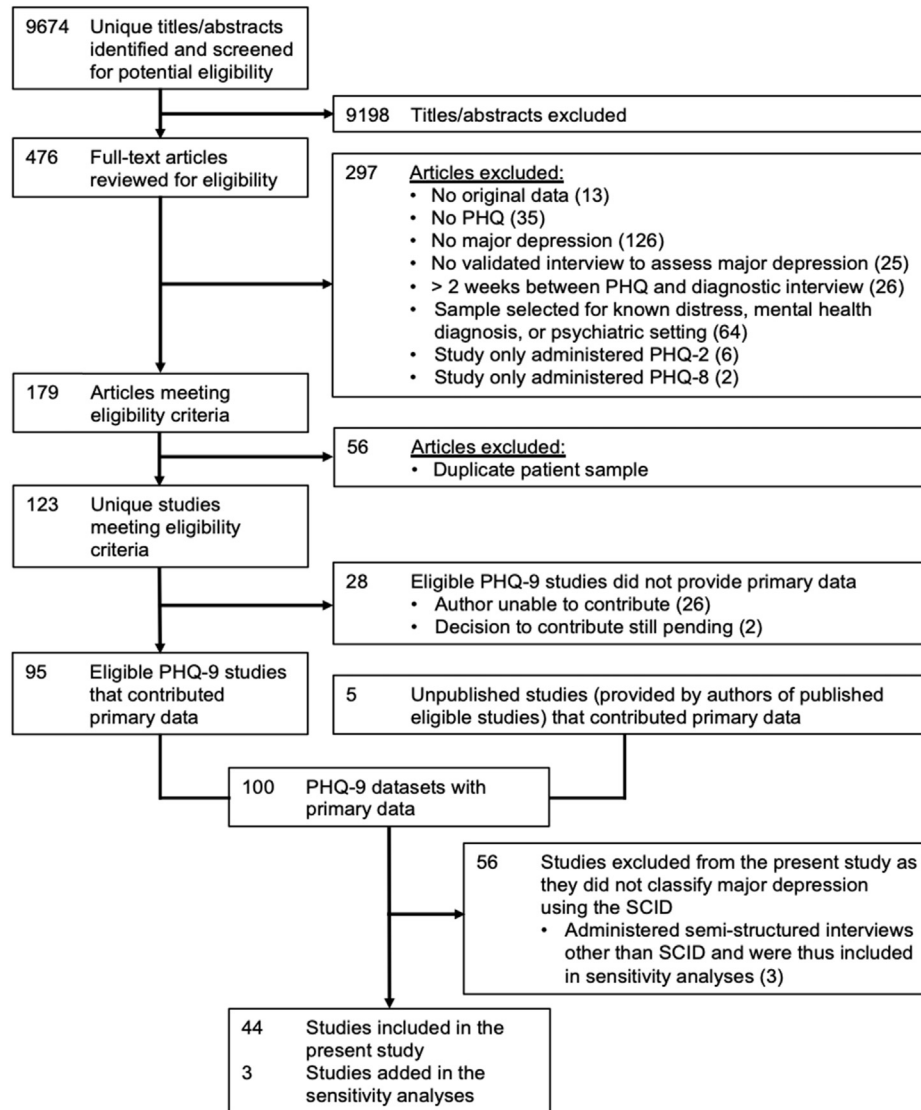


Fig. 1. Flow diagram of study selection process.

Differences in prevalence (PHQ-9 ≥ 10 minus SCID) ranged from -6.0% to 46.9% . The pooled difference was 11.9% (95% CI: 9.3% , 14.6% ; τ^2 : 0.007).

The ratio of PHQ-9 ≥ 10 prevalence to SCID-based prevalence ranged from 0.7 to 10.0 times (mean: 2.5 and median: 1.9). The mean ratio was 3.8 times for the 17 studies with SCID-based prevalence $< 10\%$ (mean difference: 13.3%), 2.0 times for the 16 studies with SCID-based prevalence between 10% and 20% (mean difference: 12.7%), and 1.3 times for the 11 studies with SCID-based prevalence $\geq 20\%$ (mean difference: 8.9%).

3.3. Prevalence matching

PHQ-9 ≥ 14 (pooled difference in prevalence: 0.5% , 95% CI: -1.7% , 2.6% , τ^2 : 0.005) and the PHQ-9 diagnostic algorithm (pooled difference in prevalence: -0.7% ,

95% CI: -3.2% , 1.8% ; τ^2 : 0.006) provided prevalence closest to SCID-based prevalence. Pooled differences in prevalence for PHQ-9 ≥ 13 and ≥ 15 compared with SCID were 2.6% and -2.0% .

In the 44 individual SCID studies, differences between the percentage of participants with PHQ-9 ≥ 14 and SCID major depression ranged from -18.7% to 29.7% (mean absolute difference: 4.8%). Of 44 prevalence estimates based on PHQ-9 ≥ 14 , 24 (54.5%) were ≤ 0.75 times or ≥ 1.25 times the SCID-based prevalence. The 95% prediction interval for the difference in prevalence was -13.6% to 14.5% . For the PHQ-9 diagnostic algorithm, study-level differences in prevalence ranged from -20.1% to 27.1% (mean absolute difference: 5.6%). Of 41 prevalence estimates based on the PHQ-9 diagnostic algorithm, 28 (68.3%) were ≤ 0.75 times or ≥ 1.25 times the SCID-based prevalence. The 95% prediction interval for the

difference in prevalence was -16.4% to 15.0% . No study or participant characteristics were significantly associated with differences in prevalence for either of the PHQ-9 prevalence match scoring approaches compared with SCID.

3.4. Sensitivity analyses

Results for all analyses were similar when data from the three studies with semistructured interviews other than the SCID were added or when the 612 participants without data to determine PHQ-9 diagnostic algorithm classification were excluded.

4. Discussion

Primary studies and meta-analyses that describe their results as reflecting the prevalence of depression or depressive disorders are frequently based on depression screening tools, which are not designed for this purpose, rather than validated diagnostic interviews [4,5]. The PHQ-9 is often used to generate what are described by researchers as depression prevalence estimates. The present study found that using PHQ-9 ≥ 10 to assess depression prevalence, which is commonly done, overestimated depression prevalence compared with prevalence based on actual diagnostic criteria by 11.9% (mean ratio: 2.5 times).

These results are consistent with what was predicted in a previous analysis that used hypothetical estimates of sensitivity and specificity to demonstrate how depression screening tools would be expected to inflate prevalence [4]. Results are also consistent with the findings of a meta-research review of prevalence estimates from 69 meta-analyses that found higher mean depression prevalence based on screening or rating tools than based on diagnostic interviews [5]. Thus, if a screening tool, such as the PHQ-9 ≥ 10 , is used to estimate prevalence, prevalence will appear to be substantial in virtually all populations, even when true prevalence is very low. This could have important ramifications in terms of policies, service planning, and health care budgets.

Identifying a PHQ-9 cutoff that could be used to match true prevalence based on a diagnostic interview would allow researchers to use inexpensive questionnaires instead of more costly interview methods for prevalence estimation. We tested a prevalence matching approach and found that PHQ-9 ≥ 14 and the PHQ-9 diagnostic algorithm provided the smallest differences in prevalence compared with SCID major depression, but heterogeneity was high and not associated with study or participant characteristics. The mean absolute difference between prevalence based on PHQ-9 vs. SCID in individual studies was 4.8% for PHQ-9 ≥ 14 and 5.6% for the PHQ-9 diagnostic algorithm, reflecting both overestimation and underestimation. For more than half of the studies examined, PHQ-9 ≥ 14 prevalence was less than 75% or more than 125% of SCID-based prevalence; for the PHQ-9 diagnostic algorithm, the fraction was over two-

thirds. The 95% prediction interval for the difference between PHQ-9 ≥ 14 and SCID-based prevalence ranged from 14% below to 15% above SCID-based prevalence; for the PHQ-9 diagnostic algorithm, it was from 16% below to 15% above SCID-based prevalence.

Researchers sometimes report prevalence estimates based on cutoffs from questionnaires, including the PHQ-9, as prevalence of “clinically significant” symptoms or “symptoms” of depression, rather than “depression” [14,76,77]. However, screening tool cutoffs do not reflect a meaningful divide between impairment and non-impairment. Patients scoring at or above virtually any cutoff would be expected to have greater impairment than patients scoring below the cutoff, but no evidence has established any single cutoff for establishing an impairment threshold or that would support clinical decision-making for individual patients without a validated clinical assessment [4].

Research on screening using the PHQ-9 would be expected to report the proportion of patients who score at or above screening cutoffs because this provides information on the number of patients who would need resources for further mental health assessment. Reporting this percentage as depression prevalence, however, would be akin, for example, to reporting the proportion of women with positive mammogram screens as the prevalence of breast cancer and, as shown in the present study, would dramatically overestimate prevalence.

This is the first study to estimate the degree to which using PHQ-9 ≥ 10 to estimate depression prevalence, a common practice, leads to overestimation of prevalence. The strengths of the study are that we incorporated data from 44 primary studies and that we directly compared PHQ-9 ≥ 10 prevalence estimates to those based on the SCID, a rigorous semistructured interview intended to facilitate the standardized application of actual diagnostic criteria by trained diagnosticians [10–12]. This study had some limitations. First, we were unable to include data from 14 of 58 published eligible datasets (24%). Second, included datasets were almost exclusively from patients in health care settings where the presence of transdiagnostic somatic symptoms and adjustment to illness or injury may have contributed to error variance [75]. Third, included datasets were from a wide range of study settings, which may account for some of the observed heterogeneity. Fourth, the overestimation of prevalence when screening tools are used is expected to be greater with lower true prevalence. This is because false positives are disproportionately high in low-prevalence populations and only minimally offset by false negative screens, which occur when true cases are missed by the screening test. However, we were unable to assess this because of the small number of heterogeneous datasets included. Fifth, not all SCID studies described interviewer qualifications; untrained interviewers may have reduced the ability to detect differences across interviews. Sixth, we only examined one depression screening tool, the PHQ-9,

Table 1. Characteristics of included studies and difference between percentage with PHQ-9 ≥ 10 and prevalence matching-based prevalence and prevalence based on diagnostic criteria for major depression

Author, year	Country	Recruited population	Total (N)	Age, mean (SD)	Female, n (%)	Major depression, n (%)	PHQ-9 ≥ 10 , n (%)
Studies from IPDMA that used the SCID and were included in main analyses							
Alamri, 2017 ^a [31]	Saudi Arabia	Hospitalized elderly in medical and surgical wards	199	70 (8)	117 (59)	24 (12.1)	44 (22.1)
Amoozegar, 2017 [34]	Canada	Migraine patients	203	43 (13)	41 (20)	49 (24.1)	72 (35.5)
Amtmann, 2015 ^b [35]	USA	Multiple sclerosis patients	164	55 (11)	127 (71)	48 (17.6)	90 (33.0)
Ayalon, 2010 [36]	Israel	Elderly primary care patients	151	76 (8)	61 (40)	6 (4.0)	14 (9.3)
Beraldi, 2014 ^c [37]	Germany	Cancer inpatients	116	52 (16)	37 (32)	7 (6.0)	21 (18.1)
Bernstein, 2018 [38]	Canada	IBD patients	240	49 (15)	151 (63)	21 (8.8)	59 (24.6)
Bhana, 2015 [39]	South Africa	Chronic care patients	679	47 (13)	509 (75)	78 (11.5)	53 (7.8)
Bombardier, 2012 [40]	USA	Inpatients with spinal cord injuries	160	42 (16)	36 (23)	14 (8.8)	43 (26.9)
Chagas, 2013 [41]	Brazil	Outpatients with Parkinson's Disease	84	59 (12)	39 (46)	19 (22.6)	30 (35.7)
Chiabanda, 2016 ^d [42]	Zimbabwe	A primary care population with high HIV prevalence	264	38 (10)	208 (79)	149 (56.4)	171 (64.8)
Eack, 2006 [43]	USA	Women seeking psychiatric services for their children at two mental health centers	48	39 (10)	48 (100)	12 (25.0)	24 (50.0)
Fann, 2005 ^{e,b,e} [32]	USA	Inpatients with traumatic brain injury	135	48 (20)	41 (28)	45 (16.2)	64 (22.5)
Fiest, 2014 ^a [44]	Canada	Epilepsy outpatients	169	39 (15)	86 (51)	23 (13.6)	37 (21.9)
Fischer, 2014 ^f [45]	Germany	Heart failure patients	194	66 (11)	40 (21)	11 (5.7)	37 (19.1)
Gjerdingen, 2009 ^g [46]	USA	Mothers registering their newborns for well-child visits at medical or pediatric clinics	419	30 (6)	419 (100)	19 (4.5)	49 (11.7)
Gräfe, 2004 ^h [47]	Germany	Medical and psychosomatic outpatients	494	42 (14)	331 (67)	67 (13.6)	166 (33.6)
Green, 2017 [48]	USA	Returning veterans	176	37 (10)	95 (54)	22 (12.5)	65 (36.9)
Green, 2018 [49]	Kenya	Pregnant women and new mothers	192	27 (6)	192 (100)	10 (5.2)	100 (52.1)
Haroz, 2017 [50]	Myanmar	Primary care patients	132	48 (14)	86 (65)	29 (22.0)	25 (18.9)
Hitchon, 2019 ^{h,i} [51]	Canada	Rheumatoid arthritis patients	148	61 (12)	124 (84)	16 (10.8)	44 (29.7)
Khamseh, 2011 ^{d,i} [52]	Iran	Type 2 diabetes patients	184	56 (9)	96 (52)	79 (42.9)	103 (56.0)
Kwan, 2012 [53]	Singapore	Post-stroke inpatients undergoing rehabilitation	113	60 (12)	37 (33)	3 (2.7)	24 (21.2)
Lambert, 2015 [54]	Australia	Cancer patients	147	58 (10)	96 (65)	21 (14.3)	38 (25.9)
Lara, 2015 [55]	Mexico	Pregnant women during the third trimester of pregnancy	280	29 (6)	280 (100)	29 (10.4)	57 (20.4)
Marrie, 2018 [56]	Canada	Multiple sclerosis patients	244	53 (13)	198 (81)	25 (10.2)	73 (29.9)
Martin-Subero, 2017 [57]	Spain	Medical inpatients	1,003	43 (14)	457 (46)	83 (8.3)	289 (28.8)
Osório, 2009 [58]	Brazil	Women in primary care	177	33 (7)	177 (100)	60 (33.9)	62 (35.0)
Osório, 2012 [59]	Brazil	Inpatients from various clinical wards	86	49 (12)	35 (41)	28 (32.6)	41 (47.7)
Patten, 2015 [60]	Canada	Multiple sclerosis patients	143	50 (12)	110 (77)	20 (14.0)	36 (25.2)
Picardi, 2005 [61]	Italy	Inpatients with skin diseases	138	37 (13)	77 (56)	12 (8.7)	38 (27.5)
Prisnie, 2016 [62]	Canada	Stroke and transient ischemic attack patients	114	60 (16)	64 (56)	11 (9.6)	16 (14.0)
Quinn, unpublished ^{h,j}	UK	Stroke patients	146	68 (13)	65 (47)	17 (11.6)	43 (29.5)
Richardson, 2010 [63]	USA	Older adults undergoing in-home aging services care management assessment	377	77 (9)	258 (68)	95 (25.2)	117 (31.0)
Rooney, 2013 [64]	UK	Adults with cerebral glioma	126	54 (12)	54 (43)	14 (11.1)	27 (21.4)
Shinn, 2017 ^k [65]	USA	Cancer patients	139	59 (11)	139 (100)	12 (8.6)	24 (17.3)
Sidebottom, 2012 ^l [66]	USA	Pregnant women	246	25 (5)	246 (100)	12 (4.9)	59 (24.0)
Simning, 2012 [67]	USA	Older adults living in public housing	190	68 (7)	110 (58)	10 (5.3)	25 (13.2)
Spangenberg, 2015 [68]	Germany	Primary care patients	160	72 (6)	97 (61)	1 (0.6)	9 (5.6)
Turner, 2012 [69]	Australia	Stroke patients	72	67 (13)	34 (47)	13 (18.1)	22 (30.6)
Turner, unpublished ^h	Australia	Cardiac rehabilitation patients	51	60 (12)	7 (14)	4 (7.8)	6 (11.8)
Vöhringer, 2013 ^h [33]	Chile	Primary care patients	190	50 (17)	143 (75)	59 (31.1)	85 (44.7)
Wagner, 2017 ^h [70]	USA	Patients starting radiotherapy for the first diagnosis of any tumor	54	59 (11)	38 (69)	6 (4.3)	13 (5.3)
Williams, 2012 [71]	USA	Parkinson's Disease patients	235	66 (10)	76 (32)	61 (26.0)	47 (20.0)
Wittkamp, 2009 ^h [72]	Netherlands	Primary care patients at risk for depression	260	51 (12)	175 (64)	45 (11.6)	90 (22.2)
Studies from IPDMA that used other semistructured interviews and were included in sensitivity analyses							
Liu, 2011 ^m [73]	Taiwan	Primary care patients	1,532	53 (19)	933 (61)	50 (3.3)	133 (8.7)
McGuire, 2013 ⁿ [74]	USA	Acute coronary syndrome inpatients	100	63 (12)	31 (31)	9 (9.0)	25 (25.0)
Twist, 2013 ^{h,c,m,o} [75]	UK	Type 2 diabetes outpatients	360	56 (11)	172 (45)	80 (7.4)	178 (14.8)

Abbreviations: DA+, positive classification based on PHQ-9 diagnostic algorithm.

^a Study did not provide item-level data necessary to determine classification based on the PHQ-9 diagnostic algorithm.

^b Sampling weights were applied. Counts are based on actual numbers, whereas percentages are weighted.

^c One participant missing data for age.

^d Ten participants missing data for age.

^e One participant missing data for both age and sex.

^f Two participants missing data for age.

^g Twenty-one participants missing data for age.

^h Unpublished at the time of electronic database search.

ⁱ Sixty-two participants missing data to determine classification based on the PHQ-9 diagnostic algorithm.

^j Two participants missing data for age, seven participants missing data for sex, 10 participants missing data to determine classification based on the PHQ-9 diagnostic algorithm, and one participant missing data for age, sex, and diagnostic algorithm.

^k Two participants missing data for age, 14 participants missing data to determine classification based on the PHQ-9 diagnostic algorithm, and one participant missing data for age and diagnostic algorithm.

^l Four participants missing data for age.

^m Diagnostic interview: Schedules for Clinical Assessment in Neuropsychiatry.

ⁿ Diagnostic interview: Depression Interview and Structured Hamilton.

^o Eight participants missing data to determine classification based on the PHQ-9 diagnostic algorithm.

PHQ-9 ≥10		Prevalence matching					
% difference: PHQ-9 ≥ 10–Major depression	Ratio: PHQ-9 ≥10/Major depression	PHQ-9 ≥14, n (%)	% difference: PHQ-9 ≥14–Major depression	Ratio: PHQ-9 ≥14/Major depression	PHQ-9 DA+, n (%)	% difference: PHQ-9 DA+ –Major depression	Ratio: PHQ-9 DA+ /Major depression
Studies from IPDMA that used the SCID and were included in main analyses							
10.1	1.8	25 (12.6)	0.5	1.0	—	—	—
11.3	1.5	40 (19.7)	–4.4	0.8	36 (17.7)	–6.4	0.7
15.4	1.9	55 (20.2)	2.6	1.1	42 (15.4)	–2.2	0.9
5.3	2.3	7 (4.6)	0.7	1.2	6 (4.0)	0.0	1.0
12.1	3.0	4 (3.4)	–2.6	0.6	2 (1.7)	–4.3	0.3
15.8	2.8	33 (13.8)	5.0	1.6	25 (10.4)	1.7	1.2
–3.7	0.7	26 (3.8)	–7.7	0.3	15 (2.2)	–9.3	0.2
18.1	3.1	23 (14.4)	5.6	1.6	17 (10.6)	1.9	1.2
13.1	1.6	16 (19.0)	–3.6	0.8	12 (14.3)	–8.3	0.6
8.3	1.1	122 (46.2)	–10.2	0.8	96 (36.4)	–20.1	0.6
25.0	2.0	17 (35.4)	10.4	1.4	17 (35.4)	10.4	1.4
6.3	1.4	33 (12.2)	–4.0	0.8	—	—	—
8.3	1.6	17 (10.1)	–3.6	0.7	17 (10.1)	–3.6	0.7
13.4	3.4	19 (9.8)	4.1	1.7	20 (10.3)	4.6	1.8
7.2	2.6	26 (6.2)	1.7	1.4	31 (7.4)	2.9	1.6
20.0	2.5	97 (19.6)	6.1	1.4	86 (17.4)	3.8	1.3
24.4	3.0	31 (17.6)	5.1	1.4	32 (18.2)	5.7	1.5
46.9	10.0	67 (34.9)	29.7	6.7	62 (32.3)	27.1	6.2
–3.0	0.9	16 (12.1)	–9.8	0.6	13 (9.8)	–12.1	0.4
18.9	2.8	22 (14.9)	4.1	1.4	26 (17.6)	6.8	1.6
13.0	1.3	81 (44.0)	1.1	1.0	55 (45.1)	6.6	1.2
18.6	8.0	9 (8.0)	5.3	3.0	7 (6.2)	3.5	2.3
11.6	1.8	21 (14.3)	0.0	1.0	18 (12.2)	–2.0	0.9
10.0	2.0	21 (7.5)	–2.9	0.7	23 (8.2)	–2.1	0.8
19.7	2.9	43 (17.6)	7.4	1.7	36 (14.8)	4.5	1.4
20.5	3.5	154 (15.4)	7.1	1.9	143 (14.3)	6.0	1.7
1.1	1.0	45 (25.4)	–8.5	0.8	43 (24.3)	–9.6	0.7
15.1	1.5	26 (30.2)	–2.3	0.9	26 (30.2)	–2.3	0.9
11.2	1.8	24 (16.8)	2.8	1.2	12 (8.4)	–5.6	0.6
18.8	3.2	21 (15.2)	6.5	1.8	18 (13.0)	4.3	1.5
4.4	1.5	11 (9.6)	0.0	1.0	9 (7.9)	–1.8	0.8
17.8	2.5	17 (11.6)	0.0	1.0	17 (12.6)	1.5	1.1
5.8	1.2	65 (17.2)	–8.0	0.7	60 (15.9)	–9.3	0.6
10.3	1.9	15 (11.9)	0.8	1.1	13 (10.3)	–0.8	0.9
8.6	2.0	11 (7.9)	–0.7	0.9	8 (6.5)	2.4	1.6
19.1	4.9	32 (13.0)	8.1	2.7	32 (13.0)	8.1	2.7
7.9	2.5	11 (5.8)	0.5	1.1	9 (4.7)	–0.5	0.9
5.0	9.0	4 (2.5)	1.9	4.0	4 (2.5)	1.9	4.0
12.5	1.7	12 (16.7)	–1.4	0.9	9 (12.5)	–5.6	0.7
3.9	1.5	2 (3.9)	–3.9	0.5	2 (3.9)	–3.9	0.5
13.7	1.4	54 (28.4)	–2.6	0.9	—	—	—
0.9	1.2	7 (2.8)	–1.5	0.7	6 (2.4)	–1.9	0.6
–6.0	0.8	17 (7.2)	–18.7	0.3	17 (7.2)	–18.7	0.3
10.6	1.9	49 (11.6)	0.0	1.0	44 (10.4)	–1.2	0.9
Studies from IPDMA that used other semistructured interviews and were included in sensitivity analyses							
5.4	2.7	46 (3.0)	–0.3	0.9	50 (3.3)	0.0	1.0
16.0	2.8	13 (13.0)	4.0	1.4	12 (12.0)	3.0	1.3
7.4	2.0	112 (9.3)	1.9	1.3	97 (8.2)	0.7	1.1

although we expect that other tools would similarly exaggerate depression prevalence [4,5].

In summary, we found that using PHQ-9 ≥ 10 to estimate depression prevalence results in estimates that are, on average, 12% greater than what would be obtained using validated semistructured diagnostic interviews. Substantial heterogeneity presents a barrier to using statistical methods to estimate major depression prevalence based on PHQ-9 ≥ 10 or based on any other PHQ-9 cutoff. Researchers should not report results from the PHQ-9 as prevalence of major depression. Users of evidence should evaluate reports of prevalence with caution and ensure that they are based on methods intended to classify major depression.

CRedit authorship contribution statement

Brooke Levis: Conceptualization, Methodology, Formal analysis, Writing - original draft, Writing - review & editing. **Andrea Benedetti:** Conceptualization, Methodology, Formal analysis, Writing - review & editing. **John P.A. Ioannidis:** Conceptualization, Methodology, Writing - review & editing. **Ying Sun:** Formal analysis, Writing - review & editing. **Zelalem Negeri:** Formal analysis, Writing - review & editing. **Chen He:** Formal analysis, Writing - review & editing. **Yin Wu:** Formal analysis, Writing - review & editing. **Ankur Krishnan:** Formal analysis, Writing - review & editing. **Parash Mani Bhandari:** Formal analysis, Writing - review & editing. **Dipika Neupane:** Formal analysis, Writing - review & editing. **Mahrukh Imran:** Formal analysis, Writing - review & editing. **Danielle B. Rice:** Formal analysis, Writing - review & editing. **Kira E. Riehm:** Formal analysis, Writing - review & editing. **Nazanin Saadat:** Formal analysis, Writing - review & editing. **Marleine Azar:** Formal analysis, Writing - review & editing. **Jill Boruff:** Conceptualization, Methodology, Investigation, Writing - review & editing. **Pim Cuijpers:** Conceptualization, Methodology, Writing - review & editing. **Simon Gilbody:** Conceptualization, Methodology, Writing - review & editing. **Lorie A. Kloda:** Conceptualization, Methodology, Investigation, Writing - review & editing. **Dean McMillan:** Conceptualization, Methodology, Writing - review & editing. **Scott B. Patten:** Conceptualization, Methodology, Data curation, Writing - review & editing. **Ian Shrier:** Conceptualization, Methodology, Writing - review & editing. **Roy C. Ziegelstein:** Conceptualization, Methodology, Writing - review & editing. **Sultan H. Alamri:** Data curation, Writing - review & editing. **Dagmar Amtmann:** Data curation, Writing - review & editing. **Liat Ayalon:** Data curation, Writing - review & editing. **Hamid R. Baradaran:** Data curation, Writing - review & editing. **Anna Beraldi:** Data curation, Writing - review & editing. **Charles N. Bernstein:** Data curation, Writing - review & editing. **Arvin Bhana:** Data curation, Writing - review & editing. **Charles H. Bombardier:** Data curation, Writing - review & editing.

Gregory Carter: Data curation, Writing - review & editing. **Marcos H. Chagas:** Data curation, Writing - review & editing. **Dixon Chibanda:** Data curation, Writing - review & editing. **Kerrie Clover:** Data curation, Writing - review & editing. **Yeates Conwell:** Data curation, Writing - review & editing. **Crisanto Diez-Quevedo:** Data curation, Writing - review & editing. **Jesse R. Fann:** Data curation, Writing - review & editing. **Felix H. Fischer:** Data curation, Writing - review & editing. **Leila Gholizadeh:** Data curation, Writing - review & editing. **Lorna J. Gibson:** Data curation, Writing - review & editing. **Eric P. Green:** Data curation, Writing - review & editing. **Catherine G. Greeno:** Data curation, Writing - review & editing. **Brian J. Hall:** Data curation, Writing - review & editing. **Emily E. Haroz:** Data curation, Writing - review & editing. **Khalida Ismail:** Data curation, Writing - review & editing. **Nathalie Jetté:** Data curation, Writing - review & editing. **Mohammad E. Khamseh:** Data curation, Writing - review & editing. **Yunxin Kwan:** Data curation, Writing - review & editing. **Maria Asunción Lara:** Data curation, Writing - review & editing. **Shen-Ing Liu:** Data curation, Writing - review & editing. **Sonia R. Loureiro:** Data curation, Writing - review & editing. **Bernd Löwe:** Data curation, Writing - review & editing. **Ruth Ann Marrie:** Data curation, Writing - review & editing. **Laura Marsh:** Data curation, Writing - review & editing. **Anthony McGuire:** Data curation, Writing - review & editing. **Kumiko Muramatsu:** Data curation, Writing - review & editing. **Laura Navarrete:** Data curation, Writing - review & editing. **Flávia L. Osório:** Data curation, Writing - review & editing. **Inge Petersen:** Data curation, Writing - review & editing. **Angelo Picardi:** Data curation, Writing - review & editing. **Stephanie L. Pugh:** Data curation, Writing - review & editing. **Terence J. Quinn:** Data curation, Writing - review & editing. **Alasdair G. Rooney:** Data curation, Writing - review & editing. **Eileen H. Shinn:** Data curation, Writing - review & editing. **Abbey Sidebottom:** Data curation, Writing - review & editing. **Lena Spangenberg:** Data curation, Writing - review & editing. **Pei Lin Lynnette Tan:** Data curation, Writing - review & editing. **Martin Taylor-Rowan:** Data curation, Writing - review & editing. **Alyna Turner:** Data curation, Writing - review & editing. **Henk C. van Weert:** Data curation, Writing - review & editing. **Paul A. Vöhringer:** Data curation, Writing - review & editing. **Lynne I. Wagner:** Data curation, Writing - review & editing. **Jennifer White:** Data curation, Writing - review & editing. **Kirsty Winkley:** Data curation, Writing - review & editing. **Brett D. Thombs:** Conceptualization, Methodology, Formal analysis, Writing - original draft, Writing - review & editing.

Acknowledgments

Authors' contributions: B.Le., A.Ben., J.P.A.I., J.B., P.C., S.G., L.A.K., D.M., S.B.P., I.S., R.C.Z., and B.D.T.

were responsible for the study conceptualization and methodology. J.B. and L.A.K. were responsible for study investigation, and they designed and conducted database searches to identify eligible studies. S.B.P., S.H.A., D.A., L.A., H.R.B., A.Ber., C.N.B., A.B., C.H.B., G.C., M.H.C., D.C., K.C., Y.C., C.D.Q., J.R.F., F.H.F., L.G., L.J.G., E.P.G., C.G.G., B.J.H., E.E.H., K.I., N.J., M.E.K., Y.K., M.A.L., S.I.L., S.R.L., B.Lö., R.A.M., L.M., A.M., K.M., L.N., F.L.O., I.P., A.P., S.L.P., T.J.Q., A.G.R., E.H.S., A.S., L.S., P.L.L.T., M.T.R., A.T., H.C.v.W., P.A.V., L.I.W., J.W., and K.W. contributed to data curation of primary datasets that were included in this study. B.Le., Y.S., Z.N., C.H., Y.W., A.K., P.M.B., D.N., M.I., D.B.R., K.E.R., N.S., M.A., and B.D.T. contributed to formal analysis including data extraction and coding for the individual participant data meta-analysis. B.Le., A.Ben., and B.D.T. conducted formal analyses and interpreted results. B.Le. and B.D.T. drafted the original manuscript. All authors provided a critical review and approved the final manuscript. B.D.T. is the guarantor.

Funding: This study was funded by the Canadian Institutes of Health Research (CIHR; KRS-134297, PCG-155468, PJT-162206). Dr. Levis and Dr. Wu were supported by Fonds de recherche du Québec - Santé (FRQS) Postdoctoral Training Fellowships. Drs. Benedetti and Thombs were supported by an FRQS researcher salary award. Mr. Bhandari was supported by a studentship from the Research Institute of the McGill University Health Centre. Ms. Neupane was supported by G.R. Caverhill Fellowship from the Faculty of Medicine, McGill University. Ms. Rice was supported by a Vanier Canada Graduate Scholarship. Ms. Riehm and Ms. Saadat were supported by CIHR Frederick Banting and Charles Best Canada Graduate Scholarship master's awards. The primary studies by Fiest et al., Patten et al., Amoozegar et al. and Prinsie et al. were supported by the Cumming School of Medicine, University of Calgary, and Alberta Health Services through the Calgary Health Trust, as well as the Hotchkiss Brain Institute. Dr. Patten was supported by a Senior Health Scholar award from Alberta Innovates Health Solutions. Dr. Jetté was supported by a Canada Research Chair in Neurological Health Services Research and an AIHS Population Health Investigator Award. The primary study by Amtmann et al. was supported by a grant from the Department of Education (NIDRR grant number H133B080025) and by the National Multiple Sclerosis Society (MB 0008). 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 studies by Marrie et al. and Bernstein et al. were supported by CIHR (THC-135234) and Crohn's and Colitis Canada. Dr. Bernstein was supported in part by the Bingham Chair in Gastroenterology. Dr. Marrie was supported by the Waugh Family Chair in Multiple Sclerosis. The primary study by Bhana et al. was output of the Program for

Improving Mental health care (PRIME) and was supported by the UK Department for International Development (201446). The views expressed do not necessarily reflect the UK Government's official policies. 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). The primary study by Chibanda et al. was supported by a grant from Grand Challenges Canada (0087-04). Dr. Conwell received support from NIMH (R24MH071604) and the Centers for Disease Control and Prevention (R49 CE002093). The primary study by Martin-Subero et al. was supported in part by a grant from the Spanish Ministry of Health's Health Research Fund (Fondo de Investigaciones Sanitarias, project 97/1184). Collection of data for the primary study by Fann et al. was supported by grant RO1 HD39415 from the US National Center for Medical Rehabilitation Research. The primary study by Fischer et al. was funded by the German Federal Ministry of Education and Research (01GY1150). 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 Gholizadeh et al. was supported by University of Technology Sydney under UTS Research Reestablishment Grants. The primary study by Green et al. (2018) was supported by a grant from the Duke Global Health Institute (453-0751). The primary study by Eack et al. was funded by the NIMH (R24 MH56858). The primary study by Garabiles et al. was supported by the Macao (SAR) Government, through the University of Macau RSKTO grants: MYRG-2014-111. The primary study by Haroz et al. was supported by the United States Agency for International Development Victims of Torture Fund: AID-DFD A-00-08-00308. Dr. Haroz was supported by a NIMH T32 predoctoral training grant (MH014592-38) and postdoctoral training grant (MH103210) during the conduct of primary study. 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 primary study by Lara et al. was supported by the Consejo Nacional de Ciencia y Tecnología/National Council for Science and Technology (CB-2009-133923-H). The primary study by Liu et al. (2011) was funded by a grant from the National Health Research Institute, Republic of China (NHRI-EX97-9706PI). The primary studies by Osório et al. (2012) 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. Bernd Löwe received research grants from Pfizer, Germany, and from the medical faculty of the University of Heidelberg, Germany (project 121/2000) for the study by Gräfe et al. Dr. Marrie was supported

by the Waugh Family Chair in Multiple Sclerosis and the Research Manitoba Chair, and CIHR grants, during the conduct of the study. 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 Muramatsu et al. (2018) was supported by grants from Niigata Seiryō University. Dr. Osório was supported by Productivity Grants (PQ-CNPq-2 -number 301321/2016-7). The primary study by Picardi et al. was supported by funds for current research from the Italian Ministry of Health. The primary study by Wagner et al. was supported by grants U10CA21661, U10CA180868, U10CA180822, and U10CA37422 from the National Cancer Institute. The study was also funded in part by a grant from the Pennsylvania Department of Health. The department specifically disclaims responsibility for any analyses, interpretations, or conclusions of the primary study. 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 Shinn et al. was supported by grant NCI K07 CA 093512 and the Lance Armstrong Foundation. 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). The primary study by Spangenberg et al. was supported by a junior research grant from the medical faculty, University of Leipzig. Collection of data for the studies by Turner et al. (2012) were funded by a bequest from Jennie Thomas through the Hunter Medical Research Institute. 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. Dr Vöhringer was supported by the Fund for Innovation and Competitiveness of the Chilean Ministry of Economy, Development and Tourism, through the Millennium Scientific Initiative (grant number IS130005). No other authors reported funding for primary studies or for their work on the present study. No sponsor or funder was involved in the study design; in the collection, analysis and interpretation of the data; in the writing of the report; or in the decision to submit the paper for publication.

Conflicts of interests: All authors have completed the ICJME uniform disclosure form and declare: no support from any organization for the submitted work; no financial relationships with any organizations that might have an interest in the submitted work in the previous 3 years with the following exceptions: C.N.B. declares that he receives grants and personal fees from Abbvie, Janssen, Pfizer, and Takeda; grants from Shire Canada, Celgene,

Boehringer Ingelheim, and Roche; and personal fees from Mylan Pharmaceuticals; outside the submitted work. K.I. declares that she has received an honorarium for speaker fees for educational lectures for Sanofi, Sunovion, Janssen, and Novo Nordisk. S.L.P. declares that she received salary support from Pfizer-Astellia and Millennium, outside the submitted work. L.I.W. declares that she receives personal fees from Celgene, outside 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.

Data sharing: Statistical codes and dataset used in the individual patient data meta-analysis can be requested from the corresponding author, Dr. Brett D. Thombs.

References

- [1] Rogan WJ, Gladen B. Estimating prevalence from the results of a screening test. *Am J Epidemiol* 1978;107:71–6.
- [2] Wittchen H-U. Reliability and validity studies of the WHO-Composite International Diagnostic Interview (CIDI): a critical review. *J Psychiatr Res* 1994;28(1):57–84.
- [3] Spitzer RL, Williams JBW, Gibbon M, First MB. The structured clinical interview for DSM-III-R (SCID) – I: history, rationale, and description. *Arch Gen Psychiatry* 1992;49:624–9.
- [4] Thombs BD, Kwakkenbos L, Levis AW, Benedetti A. Addressing overestimation of the prevalence of depression prevalence based on self-report screening questionnaires. *CMAJ* 2018;190:E44–9.
- [5] Levis B, Yan XW, He C, Sun Y, Benedetti A, Thombs BD. A comparison of depression prevalence estimates in meta-analyses based on screening tools and rating scales versus diagnostic interviews: a meta-research review. *BMC Med* 2019;17:65.
- [6] Kroenke K, Spitzer RL, Williams JB. The PHQ-9: validity of a brief depression severity measure. *J Gen Intern Med* 2001;16:606–13.
- [7] Kroenke K, Spitzer RL. The PHQ-9: a new depression diagnostic and severity measure. *Psychiatr Ann* 2002;32(9):1–7.
- [8] Spitzer RL, Kroenke K, Williams JB. Validation and utility of a self-report version of PRIME-MD: the PHQ primary care study. Primary care evaluation of mental disorders. Patient health questionnaire. *JAMA* 1999;282:1737–44.
- [9] Maurer DM, Raymond TJ, Davis BN. Depression: screening and diagnosis. *Am Fam Physician* 2018;98(8):508–15.
- [10] Diagnostic and statistical manual of mental disorders: DSM-III Washington, DC: American Psychiatric Association; 1987.
- [11] Diagnostic and statistical manual of mental disorders: DSM-IV Washington, DC: American Psychiatric Association; 1994.
- [12] Diagnostic and statistical manual of mental disorders: DSM-IV Washington, DC: American Psychiatric Association; 2000.
- [13] Levis B, Benedetti A, Thombs BD. DEPRESSION Screening Data (DEPRESSD) Collaboration. The diagnostic accuracy of the Patient Health Questionnaire-9 (PHQ-9) for screening to detect major depression: an individual participant data meta-analysis. *BMJ* 2019;365:11476.
- [14] Mata DA, Ramos MA, Bansal N, Khan R, Guille C, Di Angelantonio E, et al. Prevalence of depression and depressive symptoms among resident physicians: a systematic review and meta-analysis. *JAMA* 2015;314:2373–83.
- [15] Rotenstein LS, Ramos MA, Torre M, Segal JB, Peluso MJ, Guille C, et al. Prevalence of depression, depressive symptoms, and suicidal ideation among medical students: a systematic review and meta-analysis. *JAMA* 2016;316:2214–36.

- [16] Qato DM, Ozenberger K, Olfson M. Prevalence of prescription medications with depression as a potential adverse effect among adults in the United States. *JAMA* 2018;319:2289–98.
- [17] Dejesus RS, Vickers KS, Melin GJ, Williams MD. A system-based approach to depression management in primary care using the Patient Health Questionnaire-9. *Mayo Clin Proc* 2007;82(11):1395–402.
- [18] Kroenke K, Spitzer RL, Williams JB. The patient health questionnaire-2. *Med Care* 2003;41:1284–92.
- [19] Whooley MA. Depression and cardiovascular disease: healing the broken-hearted. *JAMA* 2006;295:2874–81.
- [20] First MB, Gibbon M. The Structured Clinical Interview for DSM-IV Axis I Disorders (SCID-I) and the Structured Clinical Interview for DSM-IV Axis II Disorders (SCID-II). *Comprehensive handbook of psychological assessment: Vol. 2. Personality assessment*. Hoboken, NJ: John Wiley & Sons, Inc.; 2004:134–43.
- [21] Kelly MJ, Dunstan FD, Lloyd K, Fone DL. Evaluating cutpoints for the MHI-5 and MCS using the GHQ-12: a comparison of five different methods. *BMC Psychiatry* 2008;8:10.
- [22] Thombs BD, Benedetti A, Kloda LA, Levis B, Nicolau I, Cuijpers P, et al. 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-analysis. *Syst Rev* 2014;3:124.
- [23] The ICD-10 Classifications of Mental and Behavioural Disorder: Clinical Descriptions and Diagnostic Guidelines. Geneva: World Health Organization; 1992.
- [24] Levis B, Benedetti A, Riehm KE, Saadat N, Levis AW, Azar M, et al. Probability of major depression diagnostic classification using semi-structured versus fully structured diagnostic interviews. *Br J Psychiatry* 2018;212(6):377–85.
- [25] Levis B, McMillan D, Sun Y, He C, Rice DB, Krishnan A, et al. Comparison of major depression diagnostic classification probability using the SCID, CIDI, and MINI diagnostic interviews among women in pregnancy or postpartum: an individual participant data meta-analysis. *Int J Methods Psychiatr Res* 2019;28(4): e1803.
- [26] Wu Y, Levis B, Sun Y, Krishnan A, He C, Riehm KE, et al. Probability of major depression diagnostic classification based on the SCID, CIDI and MINI diagnostic interviews controlling for Hospital Anxiety and Depression Scale – depression subscale scores: an individual participant data meta-analysis of 73 primary studies. *J Psychosom Res* 2020;129:109892.
- [27] Phelan E, Williams B, Meeker K, Bonn K, Frederick J, LoGerfo J, et al. A study of the diagnostic accuracy of the PHQ-9 in primary care elderly. *BMC Fam Pract* 2010;11:63.
- [28] Watnick S, Wang PL, Demadura T, Ganzini L. Validation of 2 depression screening tools in dialysis patients. *Am J Kidney Dis* 2005;46:919–24.
- [29] Liu ZW, Yu Y, Hu M, Liu HM, Zhou L, Xiao SY. PHQ-9 and PHQ-2 for screening depression in Chinese rural elderly. *PLoS One* 2016;11: e0151042.
- [30] PRESS – Peer Review of Electronic Search Strategies: 2015 Guideline Explanation and Elaboration (PRESS E&E). Ottawa: CADTH; 2016.
- [31] Alamri SH, Bari AI, Ali AT. Depression and associated factors in hospitalized elderly: a cross-sectional study in a Saudi teaching hospital. *Ann Saudi Med* 2017;37:122–9.
- [32] Fann JR, Bombardier CH, Dikmen S, Esselman P, Warms CA, Pelzer E, et al. Validity of the Patient Health Questionnaire-9 in assessing depression following traumatic brain injury. *J Head Trauma Rehabil* 2005;20:501–11.
- [33] Vöhringer PA, Jimenez MI, Igor MA, Fores GA, Correa MO, Sullivan MC, et al. Detecting mood disorder in resource-limited primary care settings: comparison of a self-administered screening tool to general practitioner assessment. *J Med Screen* 2013;20:118–24.
- [34] Amoozegar F, Patten SB, Becker WJ, Bulloch AG, Fiest KM, Davenport WJ, et al. The prevalence of depression and the accuracy of depression screening tools in migraine patients. *Gen Hosp Psychiatry* 2017;48:25–31.
- [35] Amtmann D, Bamer AM, Johnson KL, Ehde DM, Beier ML, Elzea JL, et al. A comparison of multiple patient reported outcome measures in identifying major depressive disorder in people with multiple sclerosis. *J Psychosom Res* 2015;79:550–7.
- [36] Ayalon L, Goldfracht M, Bech P. ‘Do you think you suffer from depression?’ Re-evaluating the use of a single item question for the screening of depression in older primary care patients. *Int J Geriatr Psychiatry* 2010;25:497–502.
- [37] Beraldi A, Baklayan A, Hoster E, Hiddemann W, Heussner P. Which questionnaire is most suitable for the detection of depressive disorders in haemato-oncological patients? Comparison between HADS, CES-D and PHQ-9. *Oncol Res Treat* 2014;37:108–9.
- [38] Bernstein CN, Zhang L, Lix LM, Graff LA, Walker JR, Fisk JD, et al. The validity and reliability of screening measures for depression and anxiety disorders in inflammatory bowel disease. *Inflamm Bowel Dis* 2018;24:1867–75.
- [39] Bhana A, Rathod SD, Selohilwe O, Kathree T, Petersen I. The validity of the Patient Health Questionnaire for screening depression in chronic care patients in primary health care in South Africa. *BMC Psychiatry* 2015;15:118.
- [40] Bombardier CH, Kalpakjian CZ, Graves DE, Dyer JR, Tate DG, Fann JR. Validity of the Patient Health Questionnaire-9 in assessing major depressive disorder during inpatient spinal cord injury rehabilitation. *Arch Phys Med Rehabil* 2012;93:1838–45.
- [41] Chagas MH, Tumas V, Rodrigues GR, Machado-de-Sousa JP, Filho AS, Hallak JE, et al. Validation and internal consistency of Patient Health Questionnaire-9 for major depression in Parkinson’s disease. *Age Ageing* 2013;42:645–9.
- [42] Chibanda D, Verhey R, Gibson LJ, Munetsi E, Machando D, Rusakaniko S, et al. Validation of screening tools for depression and anxiety disorders in a primary care population with high HIV prevalence in Zimbabwe. *J Affect Disord* 2016;198:50–5.
- [43] Eack SM, Greeno CG, Lee BJ. Limitations of the Patient Health Questionnaire in identifying anxiety and depression in community mental health: many cases are undetected. *Res Soc Work Pract* 2006;16:625–31.
- [44] Fiest KM, Patten SB, Wiebe S, Bulloch AG, Maxwell CJ, Jette N. Validating screening tools for depression in epilepsy. *Epilepsia* 2014;55:1642–50.
- [45] Fischer HF, Klug C, Roeper K, Blozik E, Edelmann F, Eisele M, et al. Screening for mental disorders in heart failure patients using computer-adaptive tests. *Qual Life Res* 2014;23:1609–18.
- [46] Gjerdingen D, Crow S, McGovern P, Miner M, Center B. Postpartum depression screening at well-child visits: validity of a 2-question screen and the PHQ-9. *Ann Fam Med* 2009;7:63–70.
- [47] Gräfe K, Zipfel S, Herzog W, Löwe B. Screening for psychiatric disorders with the Patient Health Questionnaire (PHQ). Results from the German validation study. *Diagnostica* 2004;50:171–81.
- [48] Green JD, Annunziata A, Kleiman SE, Bovin MJ, Harwell AM, Fox AM, et al. Examining the diagnostic utility of the DSM-5 PTSD symptoms among male and female returning veterans. *Depress Anxiety* 2017;34:752–60.
- [49] Green EP, Tuli H, Kwobah E, Menya D, Chesire I, Schmidt C. Developing and validating a perinatal depression screening tool in Kenya blending Western criteria with local idioms: a mixed methods study. *J Affect Disord* 2018;228:49–59.
- [50] Haroz EE, Bass J, Lee C, Oo SS, Lin K, Kohrt B, et al. Development and cross-cultural testing of the International Depression Symptom Scale (IDSS): a measurement instrument designed to represent global presentations of depression. *Glob Ment Health* 2017;4:e17.
- [51] Hitchon CA, Zhang L, Peschken CA, Lix LM, Graff LA, Fisk JD, et al. The validity and reliability of screening measures for depression and anxiety disorders in rheumatoid arthritis. *Arthritis Care Res* 2019.
- [52] Khamseh ME, Baradaran HR, Javanbakht A, Mirghorbani M, Yadollahi Z, Malek M. Comparison of the CES-D and PHQ-9

- depression scales in people with type 2 diabetes in Tehran, Iran. *BMC Psychiatry* 2011;11:61.
- [53] Kwan Y, Tham WY, Ang A. Validity of the Patient Health Questionnaire-9 (PHQ-9) in the screening of post-stroke depression in a multi-ethnic population. *Biol Psychiatry* 2012;71:141S.
- [54] Lambert SD, Clover K, Pallant JF, Britton B, King MT, Mitchell AJ, et al. Making sense of variations in prevalence estimates of depression in cancer: a co-calibration of commonly used depression scales using Rasch analysis. *J Natl Compr Canc Netw* 2015;13:1203–11.
- [55] Lara MA, Navarrete L, Nieto L, Martín JP, Navarro JL, Lara-Tapia H. Prevalence and incidence of perinatal depression and depressive symptoms among Mexican women. *J Affect Disord* 2015;175:18–24.
- [56] Marrie RA, Zhang L, Lix LM, Graff LA, Walker JR, Fisk JD, et al. The validity and reliability of screening measures for depression and anxiety disorders in multiple sclerosis. *Mult Scler Relat Disord* 2018;20:9–15.
- [57] Martin-Subero M, Kroenke K, Diez-Quevedo C, Rangil T, de Antonio M, Morillas RM, et al. Depression as measured by PHQ-9 versus clinical diagnosis as an independent predictor of long-term mortality in a prospective cohort of medical inpatients. *Psychosom Med* 2017;79:273–82.
- [58] Osório FL, Vilela Mendes A, Crippa JA, Loureiro SR. Study of the discriminative validity of the PHQ-9 and PHQ-2 in a sample of Brazilian women in the context of primary health care. *Perspect Psychiatr Care* 2009;45:216–27.
- [59] Osório FL, Carvalho AC, Fracalossi TA, Crippa JA, Loureiro ES. Are two items sufficient to screen for depression within the hospital context? *Int J Psychiatry Med* 2012;44:141–8.
- [60] Patten SB, Burton JM, Fiest KM, Wiebe S, Bulloch AG, Koch M, et al. Validity of four screening scales for major depression in MS. *Mult Scler* 2015;21:1064–71.
- [61] Picardi A, Adler DA, Abeni D, Chang H, Pasquini P, Rogers WH, et al. Screening for depressive disorders in patients with skin diseases: a comparison of three screeners. *Acta Derm Venereol* 2005;85:414–9.
- [62] Prisdie JC, Fiest KM, Coutts SB, Patten SB, Atta CA, Blaikie L, et al. Validating screening tools for depression in stroke and transient ischemic attack patients. *Int J Psychiatry Med* 2016;51:262–77.
- [63] Richardson TM, He H, Podgorski C, Tu X, Conwell Y. Screening depression aging services clients. *Am J Geriatr Psychiatry* 2010;18:1116–23.
- [64] Rooney AG, McNamara S, Mackinnon M, Fraser M, Rampling R, Carson A, et al. Screening for major depressive disorder in adults with cerebral glioma: an initial validation of 3 self-report instruments. *Neuro Oncol* 2013;15:122–9.
- [65] Shinn EH, Valentine A, Baum G, Carmack C, Kilgore K, Bodurka D, et al. Comparison of four brief depression screening instruments in ovarian cancer patients: diagnostic accuracy using traditional versus alternative cutpoints. *Gynecol Oncol* 2017;145:562–8.
- [66] Sidebottom AC, Harrison PA, Godecker A, Kim H. Validation of the patient health questionnaire (PHQ)-9 for prenatal depression screening. *Arch Womens Ment Health* 2012;15:367–74.
- [67] Simning A, van Wijngaarden E, Fisher SG, Richardson TM, Conwell Y. Mental healthcare need and service utilization in older adults living in public housing. *Am J Geriatr Psychiatry* 2012;20:441–51.
- [68] Spangenberg L, Glaesmer H, Boecker M, Forkmann T. Differences in patient health questionnaire and Aachen depression item bank scores between tablet versus paper-and-pencil administration. *Qual Life Res* 2015;24:3023–32.
- [69] Turner A, Hambridge J, White J, Carter G, Clover K, Nelson L, et al. Depression screening in stroke: a comparison of alternative measures with the structured diagnostic interview for the Diagnostic and Statistical Manual of Mental Disorders, Fourth Edition (major depressive episode) as criterion standard. *Stroke* 2012;43:1000–5.
- [70] Wagner LI, Pugh SL, Small W Jr, Kirshner J, Sidhu K, Bury MJ, et al. Screening for depression in cancer patients receiving radiotherapy: feasibility and identification of effective tools in the NRG Oncology RTOG 0841 trial. *Cancer* 2017;123:485–93.
- [71] Williams JR, Hirsch ES, Anderson K, Bush AL, Goldstein SR, Grill S, et al. A comparison of nine scales to detect depression in Parkinson disease: which scale to use? *Neurology* 2012;78:998–1006.
- [72] Wittkamp K, van Ravesteijn H, Baas K, van de Hoogen H, Schene A, Bindels P, et al. The accuracy of Patient Health Questionnaire-9 in detecting depression and measuring depression severity in high-risk groups in primary care. *Gen Hosp Psychiatry* 2009;31:451–9.
- [73] Liu SI, Yeh ZT, Huang HC, Sun FJ, Tjung JJ, Hwang LC, et al. Validation of Patient Health Questionnaire for depression screening among primary care patients in Taiwan. *Compr Psychiatry* 2011;52:96–101.
- [74] McGuire AW, Eastwood JA, Macabasco-O'Connell A, Hays RD, Doering LV. Depression screening: utility of the Patient Health Questionnaire in patients with acute coronary syndrome. *Am J Crit Care* 2013;22:12–9.
- [75] Twist K, Stahl D, Amiel SA, Thomas S, Winkley K, Ismail K. Comparison of depressive symptoms in type 2 diabetes using a two-stage survey design. *Psychosom Med* 2013;75:791–7.
- [76] Scott JE, Mathias JL, Kneebone AC. Depression and anxiety after total hip replacement among older adults; a meta-analysis. *Aging Ment Health* 2016;20(12):1243–54.
- [77] Buchberger B, Huppertz H, Krabbe L, Lux B, Mattivi JT, Sifarikas A. Symptoms of depression and anxiety in youth with type 1 diabetes: a systematic review and meta-analysis. *Psychoneuroendocrinology* 2016;70:70–84.

Appendix Methods. Search Strategies**MEDLINE (OvidSP)**

1. PHQ*.af.
2. patient health questionnaire*.af.
3. 1 or 2.
4. Mass Screening/.
5. Psychiatric Status Rating Scales/.
6. "Predictive Value of Tests"/.
7. "Reproducibility of Results"/.
8. exp "Sensitivity and Specificity"/.
9. Psychometrics/.
10. Prevalence/.
11. Reference Values/.
12. Reference Standards/.
13. exp Diagnostic Errors/.
14. Mental Disorders/di, pc [Diagnosis, Prevention & Control].
15. Mood Disorders/di, pc [Diagnosis, Prevention & Control].
16. Depressive Disorder/di, pc [Diagnosis, Prevention & Control].
17. Depressive Disorder, Major/di, pc [Diagnosis, Prevention & Control].
18. Depression, Postpartum/di, pc [Diagnosis, Prevention & Control].
19. Depression/di, pc [Diagnosis, Prevention & Control].
20. validation studies.pt.
21. comparative study.pt.
22. screen*.af.
23. prevalence.af.
24. predictive value*.af.
25. detect*.ti.
26. sensitiv*.ti.
27. valid*.ti.
28. revalid*.ti.
29. predict*.ti.
30. accur*.ti.
31. psychometric*.ti.
32. identif*.ti.
33. specificit*.ab.
34. cut?off*.ab.
35. cut* score*.ab.
36. cut?point*.ab.
37. threshold score*.ab.
38. reference standard*.ab.
39. reference test*.ab.
40. index test*.ab.
41. gold standard.ab.
42. or/4-41.
43. 3 and 42.
44. limit 43 to yr = "2000-Current".

PsycINFO (OvidSP)

1. PHQ*.af.
2. patient health questionnaire*.af.
3. 1 or 2.
4. Diagnosis/.
5. Medical Diagnosis/.
6. Psychodiagnosis/.
7. Misdiagnosis/.
8. Screening/.
9. Health Screening/.
10. Screening Tests/.
11. Prediction/.
12. Cutting Scores/.
13. Psychometrics/.
14. Test Validity/.
15. screen*.af.
16. predictive value*.af.
17. detect*.ti.
18. sensitiv*.ti.
19. valid*.ti.
20. revalid*.ti.
21. accur*.ti.
22. psychometric*.ti.
23. specificit*.ab.
24. cut?off*.ab.
25. cut* score*.ab.
26. cut?point*.ab.
27. threshold score*.ab.
28. reference standard*.ab.
29. reference test*.ab.
30. index test*.ab.
31. gold standard.ab.
32. or/4-31.
33. 3 and 32.
38. Limit 33 to "2000 to current".

Web of Science (Web of Knowledge)

#1: TS= (PHQ* OR "Patient Health Questionnaire*").
 #2: TS= (screen* OR prevalence OR "predictive value*" OR detect* OR sensitiv* OR valid* OR revalid* OR predict* OR accur* OR psychometric* OR identif* OR specificit* OR cutoff* OR "cut off*" OR "cut* score*" OR cutpoint* OR "cut point*" OR "threshold score*" OR "reference standard*" OR "reference test*" OR "index test*" OR "gold standard").

#1 AND #2.

Indexes = SCI-EXPANDED, SSCI, A&HCI, CPCI-S, CPCI-SSH.