

## Research paper

## Identifying depression with the PHQ-2: A diagnostic meta-analysis



Laura Manea<sup>a</sup>, Simon Gilbody<sup>a</sup>, Catherine Hewitt<sup>b</sup>, Alice North<sup>b</sup>, Faye Plummer<sup>b</sup>,  
Rachel Richardson<sup>b</sup>, Brett D. Thombs<sup>a,b</sup>, Bethany Williams<sup>b</sup>, Dean McMillan<sup>a,\*</sup>

<sup>a</sup> Hull York Medical School and Department of Health Sciences, University of York, United Kingdom

<sup>b</sup> Department of Health Sciences, University of York, United Kingdom

## ARTICLE INFO

## Article history:

Received 22 March 2016

Received in revised form

28 May 2016

Accepted 3 June 2016

Available online 6 June 2016

## Keywords:

Major depression

Screening

Diagnostic accuracy

PHQ-2

Ultra-brief screening instruments

Diagnostic meta-analysis

## ABSTRACT

**Background:** There is interest in the use of very brief instruments to identify depression because of the advantages they offer in busy clinical settings. The PHQ-2, consisting of two questions relating to core symptoms of depression (low mood and loss of interest or pleasure), is one such instrument.

**Method:** A systematic review was conducted to identify studies that had assessed the diagnostic performance of the PHQ-2 to detect major depression. Embase, MEDLINE, PsychINFO and grey literature databases were searched. Reference lists of included studies and previous relevant reviews were also examined. Studies were included that used the standard scoring system of the PHQ-2, assessed its performance against a gold-standard diagnostic interview and reported data on its performance at the recommended ( $\geq 3$ ) or an alternative cut-off point ( $\geq 2$ ). After assessing heterogeneity, where appropriate, data from studies were combined using bivariate diagnostic meta-analysis to derive sensitivity, specificity, likelihood ratios and diagnostic odds ratios.

**Results:** 21 studies met inclusion criteria totalling  $N=11,175$  people out of which 1529 had major depressive disorder according to a gold standard. 19 of the 21 included studies reported data for a cut-off point of  $\geq 3$ . Pooled sensitivity was 0.76 (95% CI = 0.68–0.82), pooled specificity was 0.87 (95% CI = 0.82–0.90). However there was substantial heterogeneity at this cut-off ( $I^2=81.8\%$ ). 17 studies reported data on the performance of the measure at cut-off point  $\geq 2$ . Heterogeneity was  $I^2=43.2\%$  pooled sensitivity at this cut-off point was 0.91 (95% CI = 0.85–0.94), and pooled specificity was 0.70 (95% CI = 0.64–0.76).

**Conclusion:** The generally lower sensitivity of the PHQ-2 at cut-off  $\geq 3$  than the original validation study (0.83) suggests that  $\geq 2$  may be preferable if clinicians want to ensure that few cases of depression are missed. However, in situations in which the prevalence of depression is low, this may result in an unacceptably high false-positive rate because of the associated modest specificity. These results, however, need to be interpreted with caution given the possibility of selectively reported cut-offs.

© 2016 Published by Elsevier B.V.

## 1. Introduction

Depression is common and disabling, but its management is suboptimal in primary and secondary care (Gilbody et al., 2008). Screening has been proposed as a solution to improving depression care, but the value of routine screening and case finding procedures to detect depression has not been proven (Gilbody et al., 2008; Thombs et al., 2012). Some national guidelines recommend it in primary care (U.S. Preventive Services Task Force, 2009), whereas others do not (Joffres et al., 2013; Allaby 2010).

Recently there has been an increased interest in the potential of using very brief instruments to identify patients with major

depression, because of the advantages they may offer in busy clinical settings in which time is limited (Mitchell and Coyne, 2007). One such very brief screening measure for depression is the two-item Patient Health Questionnaire (PHQ-2) (Kroenke et al., 2003), an abbreviated version of the widely used PHQ-9 (Kroenke et al., 2001). It is comprised of the first two questions of the PHQ-9, which reflect the core symptoms of depression (low mood, loss of interest/pleasure). The original validation study of the PHQ-2 provided preliminary evidence that it may be an effective screen for depression (Kroenke et al., 2003). In that study, a cut-off point of  $\geq 3$  (out of a possible score of 6) had a sensitivity of 0.83 and a specificity of 0.90 to identify major depression in a sample of 580 primary and secondary care patients, although this included only 41 patients with major depression, a small number for estimating diagnostic accuracy. This contrasts favourably with sensitivity of 0.88 and specificity of 0.88 in the nine-item PHQ-9 among the same patients (Kroenke et al., 2001).

\* Correspondence to: Hull York Medical School and Department of Health Sciences, ARRC Building, University of York, YO10 5DD, United Kingdom.

E-mail address: [dean.mcmillan@york.ac.uk](mailto:dean.mcmillan@york.ac.uk) (D. McMillan).

A previous systematic review of the diagnostic properties of the PHQ-2 identified only a small number of studies ( $N = 3$ ) that had examined the diagnostic performance of the PHQ-2 (Gilbody et al., 2007). The review concluded that no recommendations could be made about the PHQ-2 without further validation studies across a range of clinical settings and populations. The authors of the review, however, did suggest that preliminary evidence suggested that the PHQ-2 could be a brief, yet accurate tool. Since that initial review the PHQ-2 has been much more widely evaluated in primary studies, but there is not an updated systematic review. The current systematic review aims to evaluate the current evidence base for the PHQ-2 to identify patients with major depression.

## 2. Methods

### 2.1. Literature search

We searched Embase, MEDLINE, PsycINFO and grey literature databases (OIASTER, OpenGrey, ZETOC) from inception to August 2014. The search terms used for Embase, Medline and PsycINFO are given in Appendix A. The terms were adapted as necessary for the grey databases. In addition, we examined the reference lists of all included studies and previous relevant reviews, including reviews of the PHQ-9 (Gilbody et al., 2007; Wittkamp et al., 2007; Kroenke et al., 2010; Manea et al., 2012) and a review of ultra-brief screening instruments for depression (Mitchell and Coyne, 2007).

### 2.2. Study selection

A pre-piloted coding manual outlining a priori inclusion-exclusion criteria along with operational definitions of each was developed. *Population*: Any population or setting was included. *Instrument*: We included studies that used the PHQ-2 scored in the standard way (each item scored 0–3 and summed to give a total score between 0 and 6). Studies that used atypical methods of scoring the PHQ-2 (e.g., scored as positive if either item was scored as two or above) were excluded. *Comparison (reference standard)*: The accuracy of the PHQ-2 had to be assessed against a recognised gold-standard instrument for the diagnosis of either Diagnostic and Statistical Manual (DSM) or International Classification of Disease (ICD) criteria for major depression. Studies that used other reference standards, such as unaided clinician diagnosis or scores above a cut-off point on another self-report instrument, were excluded. Studies were also excluded if the target diagnosis was not major depression (e.g., any depressive disorder). *Outcome*: Studies had to report sufficient information to calculate a 2\*2 contingency table for the cut-off point  $\geq 3$  recommended by the original validation study or the lower, alternative cut-off recommended by some studies ( $\geq 2$ ). *Study design*: Any design. *Additional criterion*: Studies were excluded if the sample overlapped with that used in another included study. Citations with overlapping samples were examined to establish whether they contained information relevant to the research question that was not contained in the included report. We included in the review the study that had the larger sample or, if the samples were the same size, the study that provided all the details required for its review. No restrictions were made in terms of publication status, publication year or language.

All identified citations were first assessed on the basis of title and abstract. At this stage, the inclusion-exclusion criteria were interpreted liberally; if there was doubt about whether a citation met the criteria it was included. Full paper copies of those that passed this first sift were obtained and examined in detail against the inclusion-exclusion criteria. Studies that met this second sift were included in the systematic review. Where necessary authors

were contacted to provide further clarification or to obtain additional information.

### 2.3. Data extraction

We extracted the following data to a pre-piloted, standardised form: sample characteristics (country, setting, age, gender), sample size and percentage with major depression according to the gold standard, information on the PHQ-2 (method of administration, cut-offs reported, language), and details of the reference standard. In addition, we calculated cell Ns of the 2\*2 tables at cut-offs  $\geq 2$  and  $\geq 3$ . Again, where necessary authors were contacted to provide clarification.

### 2.4. Quality assessment

Quality assessment was conducted at the study level and used criteria based on the QUADAS-2 (the revised tool for the Quality Assessment of Diagnostic Accuracy Studies) (Whiting et al., 2011). QUADAS-2 incorporates assessments of risk of bias across four core domains: patient selection, the index test, the reference standard, and the flow and timing of assessments. The QUADAS-2 guidelines require that it is adapted for each specific review; this can involve adding or omitting questions and providing clarification about how specific questions are to be rated. We retained all of the risk of bias signaling questions and applicability questions, for which we developed specific guidance on coding in the form of a brief field guide. For the signaling question 'Is the reference standard likely to correctly classify the target condition?' we operationalised this as whether the researchers who conducted the gold standard interview had received appropriate training. For the signaling question 'Was there an appropriate interval between the index test and reference standard?' we defined an appropriate interval as less than two weeks in keeping with how this item has been applied in previous diagnostic test accuracy studies of depression (Mann et al., 2009).

We added four additional questions that were applied to studies using translated versions of the PHQ-2 and reference test. For translations of the PHQ-2, we asked whether appropriate translation methods were used and whether psychometric properties of the translated version were reported. The same two questions (appropriate translation, psychometric properties) were also applied to any translated version of the reference test.

### 2.5. Data analysis and synthesis

Sensitivity, specificity, positive and negative likelihood ratios and diagnostic odds ratios along with their associated 95% confidence intervals were calculated for cut-off points  $\geq 2$  and  $\geq 3$ . Heterogeneity was assessed using  $I^2$  for the diagnostic odds ratio, an estimate of the proportion of study variability that is due to between-study variability rather than sampling error. We considered values of  $\geq 50\%$  to indicate substantial heterogeneity (Centre for Reviews and Dissemination, 2009). Where heterogeneity was not substantial we used bivariate diagnostic meta-analyses to generate pooled estimates of sensitivity and specificity. Summary Receiver Operating Characteristics (sROC) were calculated to produce 95% confidence interval ellipses within ROC space.

Where substantial heterogeneity was identified, we conducted pre-planned subgroup analyses based on clinical setting. We further explored possible reasons for heterogeneity by conducting pre-planned meta-regressions of key descriptive variables and the quality assessment criteria (Centre for Reviews and Dissemination, 2009).

We attempted to limit publication bias by searching a range of grey literature databases. The potential for selective outcome

reporting bias related to the reporting of results for some but not other cut-off points is explored in the discussion section.

Bayesian nomograms were generated to examine the performance of the PHQ-2 at different prevalence estimates.

### 3. Results

The initial search identified 1054 unique citations (2882 citations before de-duplication). 59 of these citations met initial inclusion criteria and were selected for further screening of the full article. 21 of the 59 met final stage inclusion criteria (Kroenke et al., 2003; Arroll et al., 2010; Chagas et al., 2011; de Lima Osorio et al., 2009, Osorio et al., 2012; de Man-van Ginkel et al., 2012; Delgadillo et al., 2011; Fiest et al., 2014; Inagaki et al., 2013; Lowe et al., 2005; Margrove et al., 2011a; Phelan et al., 2010a; Richardson et al., 2010a, 2010b; Smith et al., 2010; Thombs et al., 2008a; Tsai et al., 2014; Williams et al., 2005; Zhang et al., 2013; Zuithoff et al., 2010a).

The remaining 38 were excluded for the following reasons: screening instrument was not the PHQ-2 ( $N = 9$ ), PHQ-2 was

scored in a non-standard way ( $N = 7$ ), reference standard was not a recognised gold-standard instrument ( $N = 7$ ), reference standard diagnosis was not solely major depression ( $N = 3$ ), study reported insufficient information to calculate a 2\*2 table for at least one of the cut-off points ( $N = 2$ ), and overlap in samples with included studies ( $N = 7$ ). Two additional citations were excluded because we were unable to obtain further information from the authors to establish whether they met inclusion criteria. Finally, one study was excluded, as all included patients were known to have depression and would, thus, not be screened in practice. The selection of studies is summarised in the PRISMA flowchart (Moher et al., 2009) in Fig. 1 and further details about the reasons for exclusion are given in Appendix B.

#### 3.1. Overview of included studies

Table 1 summarises the characteristics of the included studies. Three studies used general primary care samples (Kroenke et al., 2003; Arroll et al., 2010; Zuithoff et al., 2010b), with a further one focused on older adults in primary care (Phelan et al., 2010b). One study focused on patients with epilepsy, but recruited these from

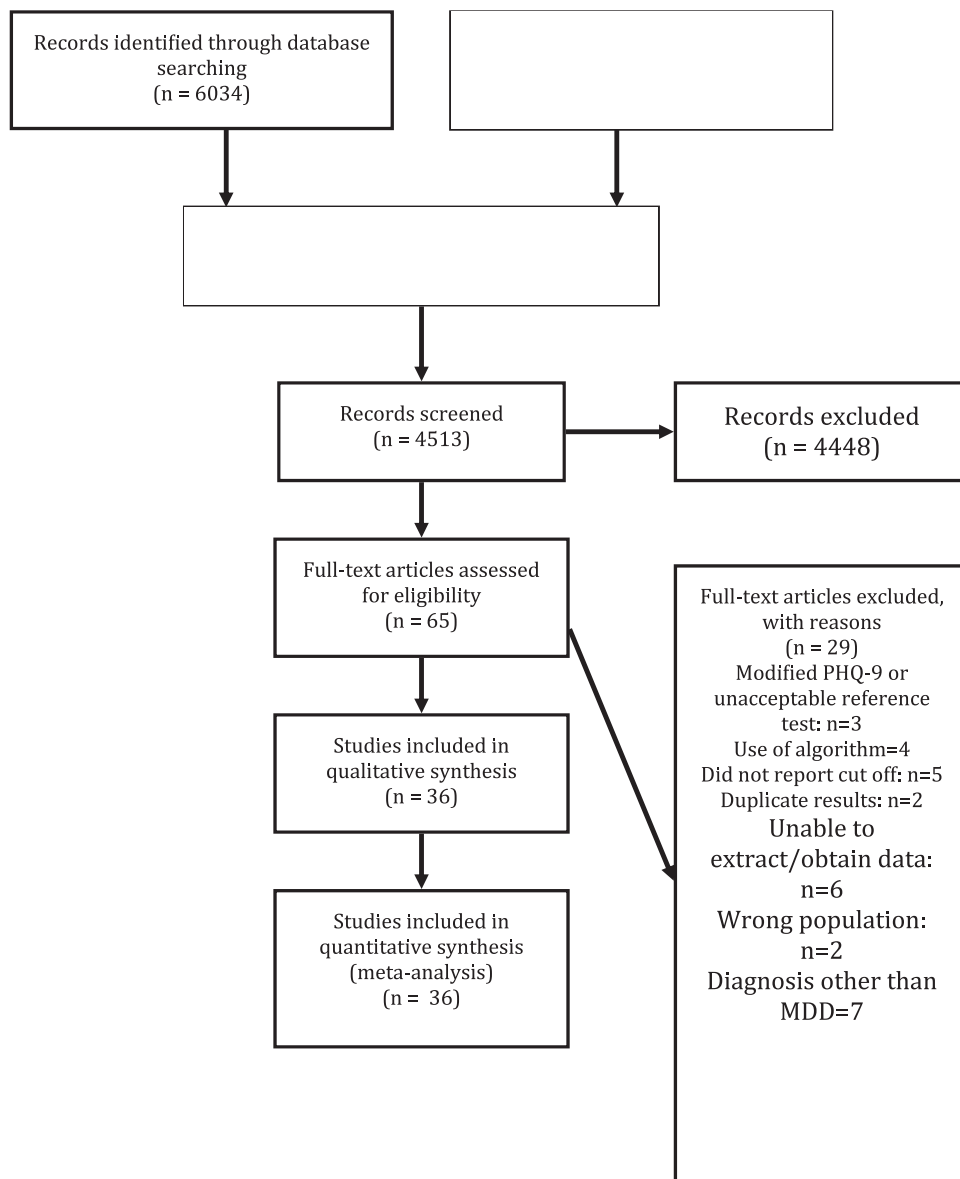


Fig. 1. PRISMA Flow diagram outlining study selection.

**Table 1**  
Descriptive characteristics of the included studies.

Study	Sample characteristics (Country, setting, age, sex)	Sample size and % depressed	PHQ-2 characteristics	Diagnostic standard
Arroll et al. (2010)	Country: New Zealand Setting: Primary care Age (yrs): Av. = 49 (range = 17–99) Female: 61%	N = 2642 Depressed: 6.2%	Administration: Not stated Language: English	DSM-IV CIDI
Chagas et al. (2011)	Country: Brazil  Setting: Movement disorders outpatient clinic Age (yrs): M = 71.09 (sd = 12.62) Female: 53%	N = 110 Depressed: 25.5%	Administration: Neurologist administered Language: Brazilian	DSM-IV SCID
De Lima Osorio et al. (2009)	Country: Brazil  Setting: Gynaecology and General Practice Age (yrs): 48% < 30 Female: 100%	N = 177 Depressed: 34%	Administration: Not stated Language: Brazilian Portuguese	DSM-IV SCID
De Lima Osorio et al. (2012)	Country: Brazil Setting: General hospital Age (yrs): M = 49 (SD = 12.4) Female: 39%	N = 100 Depressed: 2%	Administration: Not stated Language: Brazilian Portuguese	DSM-IV SCID CIDI
De Man-van Ginkel et al. (2012)	Country: Netherlands  Setting: Stroke patients  Age (yrs): M = not specified Female: % not specified	N = 164 Depressed: 12.2%	Administration: Face to face Language: Unclear (?Dutch and English)	
Delgadillo et al. (2011)	Country: UK  Setting: Community drug treatment service Age (yrs): M = 35 (range: 23–54) Female: 23%	N = 103 Depressed: 61.2%	Administration: Self-report (assistance if required) Language: English	ICD-10 CIS-R
Fiest et al. (2014)	Country: Canada Setting: Secondary care (epilepsy clinic) Age (yrs): M = 40.3 (range: 18.2–78.1) Female: 51.4%	N = 185 Depressed: 14.6%	Administration: Self-report Language: English	DSM IV/V SCID
Inagaki et al. (2013)	Country: Japan Setting: Secondary care (general medical clinic) Age (yrs): M = 73.5 (SD 12.3) Female: 59.3%	N = 104 Depressed: 7.4%	Administration: Face to face Language: Japanese	MINI
Kroenke et al. (2003)	Country: US Setting: Primary care Age (yrs): Primary: M = 46 Female: Primary = 66%	N = 580 Depressed: 7.1%	Administration: Self-report Language: English	DSM-III-R PRIME-MD
Liu et al. (2011)	Country: Taiwan Setting: Community-based primary care and hospital-based family physician clinics Age (yrs): Not reported Female: % not reported	N = 1532 Depressed: 3.3%	Administration: Not stated Language: Chinese	DSM-IV SCAN
Lowe et al. (2005)	Country: Germany Setting: Outpatient clinics and family practices Age (yrs): M = 42.0 (sd = 13.8) Female: 67.5%	N = 520 Depressed: 13.7%	Administration: Self-report Language: German	DSM-IV SCID
Margrove et al. (2011)	Country: UK Setting: Diagnosis of epilepsy in primary care Age (yrs): M = 49 (sd = 16) Female: 49.8%	N = 52 Depressed: 48.1%	Administration: Self-report Language: English	DSM-IV SCID
Phelan et al. (2010)	Country: US  Setting: Older adults in primary care clinics Age (yrs): M = 78 (sd = 7) Female: 62%	N = 69 Depressed: 12%	Administration: Self-report (assistance if required) Language: English	DSM-IV SCID
Richardson et al. (2010)	Country: US Setting: Group Health Research Institute Age (yrs): M = 15.3 (sd = 1.1) Female: 60%	N = 444 Depressed: 54.5%	Administration: Telephone administered Language: English	DSM-IV DISC
Richardson et al. (2010)	Country: US Setting: Community-based aging services agency Age (yrs): M = 76.5 (sd = 9.2) Female: 68.5%	N = 378 Depressed: 26.7%	Administration: Unclear Cut-offs: ≥ 1 to 6 Language: English	DSM-IV SCID
Smith et al. (2010)	Country: US Setting: Obstetrical settings Age (yrs): Depressed: 29.31 (sd = 5.98) Non depressed: 28.87 (sd = 6.72) Female: 100%	N = 213 Depressed: 6.1%	Administration: Not stated Language: English	DSM-IV CIDI
Thombs et al. (2008)	Country: US	N = 1024	Administration: Not stated	DSM C-DIS

Table 1 (continued)

Study	Sample characteristics (Country, setting, age, sex)	Sample size and % depressed	PHQ-2 characteristics	Diagnostic standard
	Setting: Outpatients with coronary heart disease Age (yrs): M = 67 (sd = 11) Female: 18%	Depressed: 22%	Language: English	
Tsai et al. (2014)	Country: Taiwan Setting: Community (high-schools) Age (yrs): M = 16.9 (sd = 0.6) Female: 59.6%	N = 165 Depressed 10%	Administration: Self-report Language: Chinese	DSM K-SADS-E
Williams et al. (2005)	Country: US Setting: Inpatient stroke Age (yrs): 42% < 60 Female: 51%	N = 316 Depressed: 34%	Administration: Not stated Language: English	DSM-IV SCID
Zhang et al. (2013)	Country: China Setting: Community (university students) Age (yrs): M = 21.45 (sd = 1.04) Female: 54.3%	N = 959 Depressed: 8.8%	Administration: Face to face Language: Chinese	DSM-IV SCID
Zuithoff et al. (2010)	Country: Netherlands Setting: Primary care Age (yrs): M = 51 (sd = 16.7) Female: 63%	N = 1338 Depressed: 13%	Administration: Self-report Language: Dutch	DSM-IV CIDI

Abbreviations: C-DIS = Computerised Diagnostic Interview Schedule; CIDI = Composite International Diagnostic Interview; CIS-R = Clinical Interview Schedule (Revised); DISC = Diagnostic Interview Schedule for Children; DSM-III-R = Diagnostic and Statistical Manual (Version III Revised); DSM-IV = Diagnostic and Statistical Manual (Version IV); International Classification of Diseases (Version 10); PHQ-2 = Patient Health Questionnaire two-item version; PRIME-MD = Primary Care Evaluation of Mental Disorders; SCAN = Schedule for Clinical Assessments in Neuropsychiatry; SCID = Structured Clinical Interview for DSM

primary care (Margrove et al., 2011b). A further three studies used a combination of a primary care setting and another setting, such as outpatient clinics (Lowe et al., 2005; De Lima Osorio et al., 2009; Liu et al., 2011). Eight studies recruited from hospital- or out-patient-based medical specialties (Osorio et al., 2012; de Man-van Ginkel et al., 2012; Fiest et al., 2014; Inagaki et al., 2013; Smith et al., 2010; Williams et al., 2005; Chagas et al., 2011; Thombs et al., 2008b). Of the remainder, one recruited from a community-drug treatment service (Delgadillo et al., 2011), one from a community-based aging service (Richardson et al., 2010b), one from a research institute focusing on adolescents (Richardson et al., 2010a) and two from community settings (students) (Tsai et al., 2014; Zhang et al., 2013).

All of the studies apart from two (Richardson et al., 2010; Tsai et al., 2014) had working age or older adult samples. In the majority of studies, there were markedly more females than males or the samples were entirely female. The proportion of the sample that met reference standard criteria for major depression ranged from 2% (Osorio et al., 2012) to 61.2% (Delgadillo et al., 2011). Some of the studies had a high prevalence of depression because the study design over-sampled people with positive PHQ-2 scores for administration of the reference standard (Richardson et al., 2010a; Williams et al., 2005; Margrove et al., 2011b).

Six studies stated that a self-report version of the PHQ-2 was used (Kroenke et al., 2003; Delgadillo et al., 2011; Fiest et al., 2014; Lowe et al., 2005; Tsai et al., 2014; Zuithoff et al., 2010b; Phelan et al., 2010b). In one study it was administered over the telephone (Richardson et al., 2010a) and in four studies it was administered face to face (Chagas et al., 2011; de Man-van Ginkel et al., 2012; Inagaki et al., 2013; Zhang et al., 2013); the remaining studies did not clearly state the method of administration. Translated versions of the PHQ-2 were used in ten studies (Chagas et al., 2011; Osorio et al., 2012; de Man-van Ginkel et al., 2012; Delgadillo et al., 2011; Inagaki et al., 2013; Lowe et al., 2005; Tsai et al., 2014; Zhang et al., 2013; Zuithoff et al., 2010b; Liu et al., 2011), including Brazilian, Chinese, Dutch, Japanese and German versions.

### 3.2. Quality assessment

Table 2 summarises the results of the quality assessment using QUADAS-2. The studies varied in quality. Only two of the studies

were judged to be at a low risk of bias across all of the domains (Arroll et al., 2010; Zuithoff et al., 2010b). One of these studies (Zuithoff et al., 2010b), however, was the only one not to meet all of the applicability criteria. The reference standard in Zuithoff et al. (2010b) assessed major depression over a one-year time-frame, so, unlike the PHQ-2, is not assessing current depression. This may have lowered the observed accuracy of the PHQ-2 in that study. A number of studies had high prevalence rates of depression because the studies use a design in which participants who are at an increased risk of depression (e.g. those scoring above a threshold on the PHQ-2) were more likely to be given the reference standard (Richardson et al., 2010a; Williams et al., 2005; Margrove et al., 2011b).

### 3.3. Narrative overview of diagnostic performance

Table 3 summarises the test accuracy characteristics of the PHQ-2 at the standard cut-off point of  $\geq 3$ ; Table 4 gives the same data for the alternative cut-off point of  $\geq 2$ .

Nineteen studies reported the performance of the PHQ-2 at cut-off point  $\geq 3$ . At this cut-off, sensitivity ranged from 0.39 (Thombs et al., 2008a) to 1 (Osorio et al., 2012) and specificity from 0.59 (Smith et al., 2010) to 1 (Margrove et al., 2011b). Five studies, one of which was the original validation study, were conducted in primary care. Of these, one study focused solely on people with epilepsy (Margrove et al., 2011b) so was not considered a general primary care sample.

Seventeen studies reported details of the performance of the PHQ-2 at cut-off point  $\geq 2$  (see Table 4). The distinction between the performance of the PHQ-2 in the original validation study and the other studies was less marked than at cut-off point  $\geq 3$ , though for those studies in which a diagnostic odds ratio could be calculated, the value was higher in the original validation studies than the subsequent studies.

### 3.4. Diagnostic meta-analyses

An initial diagnostic meta-analysis was run including all 19 studies reporting the performance of the PHQ-2 at cut-off point  $\geq 3$ . Pooled sensitivity was 0.76 (95% CI 0.68–0.82), pooled specificity 0.87 (95% CI 0.82–0.90), pooled positive likelihood ratio 6.02



**Table 2**  
Quality assessment of included studies.

Study	Patient selection: Consecutive or random sample	Patient selection: Avoid case-control / avoid artificially inflated base rate	Patient selection: Avoided inappropriate exclusions	Patient selection: Overall risk of bias	Index test: PHQ-2 interpreted blind to reference test	Index test: Threshold pre-specified or multiple cut-offs reported	Index test: If translated, appropriate translation	Index test: If translated, psychometric properties reported	Index test: Overall risk of bias
Arroll et al. (2010)	✓	✓	✓	Low	✓	✓	n/a	n/a	Low
Chagas et al. (2011)	✓	✓	✓	Low	✓	✓	✗	✓	Low
De Lima Osorio et al. (2009)	✓	✓	✗	Low	?	✓	?	?	Unclear
De Lima Osorio et al. (2012)	?	?	✗	High	?	✓	✓	?	Unclear
De Man-van Ginkel et al. (2012)	✓	✓	✓	Low	✓	✓	?	?	Unclear
Delgadillo et al. (2011)	✗	✓	✓	Low	✓	✓	n/a	n/a	Low
Fiest et al. (2014)	✓	✓	✓	Low	✓	✗	n/a	n/a	High
Inagaki et al. (2013)	✗	✗	✓	High	?	✓	?	?	Unclear
Kroenke et al. (2003)	✗	✓	✗	High	✓	✓	n/a	n/a	Low
Liu et al. (2011)	?	✓	?	Unclear	✓	✓	✓	✓	Low
Lowe et al. (2005)	✗	✓	✓	Low	✓	✓	✓	✓	Low
Margrove et al. (2011)	✗	✗	✓	High	✓	✓	n/a	n/a	Low
Phelan et al. (2010)	✗	✓	✓	Low	?	✓	n/a	n/a	Unclear
Richardson et al. (2010)	✗	✗	✓	High	✓	✓	n/a	n/a	Low
Richardson et al. (2010)	✗	✓	✓	Low	✓	✓	n/a	n/a	Low
Smith et al. (2010)	?	✓	?	Unclear	✓	✓	n/a	n/a	Low
Thombs et al. (2008)	✗	✓	?	Unclear	?	✓	n/a	n/a	Unclear
Tsai et al. (2014)	?	✗	✓	High	✓	✓	?	?	Unclear
Williams et al. (2005)	✗	?	✓	Unclear	✓	✓	n/a	n/a	Low
Zhang et al. (2013)	?	✓	✓	Unclear	✓	✓	✓	?	Unclear
Zuithoff et al. (2010)	✗	✓	✓	Low	✓	✓	✓	?	Low
Study	Reference test: Reference test correctly classifies target condition	Reference test: Reference test interpreted blind to PHQ-2	Reference test: If translated, appropriate translation	Reference test: If translated, psychometric properties reported	Reference test: Overall risk of bias	Flow / timing: Interval of two weeks or less	Flow / timing: All participants receive same reference test	Flow / timing: All participants included in analysis?	Flow / timing: Overall risk of bias
Arroll et al. (2010)	✓	✓	n/a	n/a	Low	✓	✓	✓	Low
Chagas et al. (2011)	✓	?	✗	✓	Unclear	✓	✓	✗	Low
De Lima Osorio et al. (2009)	✓	?	?	?	Unclear	?	✓	✓	Unclear
De Lima Osorio et al. (2012)	✓	?	?	?	Unclear	✓	✓	✗	High
De Man-van Ginkel	✓	✓	?	?	Unclear	✓	✓	✗	High

Table 2 (continued)

Study	Reference test: Re- ference test correctly classifies target condition	Reference test: Re- ference test inter- preted blind to PHQ-2	Reference test: If translated, appro- priate translation	Reference test: If trans- lated, psychometric properties reported	Reference test: Overall risk of bias	Flow / timing: Interval of two weeks or less	Flow / timing: All participants receive same reference test	Flow / timing: All participants included in analysis?	Flow / timing: Overall risk of bias
et al. (2012)									
Delgadillo et al. (2011)	✓	?	n/a	n/a	Unclear	✓	✓	✓	Low
Fiest et al. (2014)	✓	✓	n/a	n/a	Low	✓	✓	✗	High
Inagaki et al. (2013)	✓	?	✓	?	Unclear	✓	✓	✗	High
Kroenke et al. (2003)	✓	✓	n/a	n/a	Low	✓	✓	✓	Low
Liu et al. (2011)	✓	✓	?	✓	Low	✓	✓	✗	Low
Lowe et al. (2005)	✓	✓	?	?	Unclear	✓	✓	✓	Low
Margrove et al. (2011)	✓	?	n/a	n/a	Unclear	?	✓	✗	Unclear
Phelan et al. (2010)	✓	✓	n/a	n/a	Low	✓	✓	✓	Low
Richardson et al. (2010)	✓	✗	n/a	n/a	High	✓	✓	✓	Low
(Richardson et al., 2010)									
Richardson et al. (2010)	✓	?	n/a	n/a	Unclear	✓	✓	✓	Low
(Richardson et al., 2010)									
Smith et al. (2010)	✓	?	n/a	n/a	Unclear	✗	✓	✓	Low
Thombs et al. (2008)	?	✓	n/a	n/a	Unclear	✓	✓	✓	Low
Tsai et al. (2014)	✓	✓	✓	✓	Low	?	✓	✗	High
Williams et al. (2005)	✓	✗	n/a	n/a	High	✓	✓	✓	Low
Zhang et al. (2013)	✓	✓	?	?	Unclear	✓	✓	✗	High
Zuithoff et al. (2010)	✓	✓	✓	✓	Low	?	✓	✓	Low
Study	Patient selection: Applicability			Index test: Applicability			Reference test: Applicability		
Arroll et al. (2010)		✓			✓			✓	
Chagas et al. (2011)		✓			✓			✓	
De Lima Osorio et al. (2009)		✓			✓			✓	
De Lima Osorio et al. (2012)		✓			✓			✓	
De Man-van Ginkel et al. (2012)		✓			✓			✓	
Delgadillo et al. (2011)		✓			✓			✓	
Inagaki et al. (2013)		✓			✓			✓	
Fiest et al. (2014)		✓			✓			✓	
Kroenke et al. (2003)		✓			✓			✓	
Liu et al. (2011)		✓			✓			✓	
Lowe et al. (2005)		✓			✓			✓	
Margrove et al. (2011)		✓			✓			✓	
Phelan et al. (2010)		✓			✓			✓	
Richardson et al. (2010)		✓			✓			✓	
Richardson et al. (2010)		✓			✓			✓	
Smith et al. (2010)		✓			✓			✓	
Thombs et al. (2008)		✓			✓			✓	
Tsai et al. (2014)		✓			✓			✓	
Williams et al. (2005)		✓			✓			✓	
Zhang et al. (2013)		✓			✓			✓	
Zuithoff et al. (2010)		✓			✓			✗	

✓ = criterion met; ✗ = criterion not met; ? = insufficient information to code whether criterion met; n/a = not applicable

<sup>1</sup>If studies reported multiple cut-off points, 'threshold pre-specified' is coded as not applicable.

**Table 3**Diagnostic test accuracy of the PHQ-2 at cut off point  $\geq 3$ .

	Sensitivity (95% CI)	Specificity (95% CI)	+ve LR (95% CI)	-ve LR (95% CI)	DOR (95% CI)
Arroll et al. (2010)	0.61 (0.53–0.69)	0.92 (0.91–0.93)	7.68 (6.41–9.2)	0.42 (0.35–0.51)	18.3 (12.9–25.8)
Chagas et al. (2011)	0.75 (0.55–0.89)	0.89 (0.80–0.95)	6.83 (3.56–13.1)	0.28 (0.15–0.54)	24.3 (8.22–72)
De Lima Osorio et al. (2009)	0.97 (0.89–1)	0.88 (0.81–0.93)	8.08 (4.93–13.2)	0.04 (0.01–0.14)	213 (50.9–*)
De Lima Osorio et al. (2012)	1 (0.15–1)	0.75 (0.65–0.83)	4.08 (2.88–5.78)	0 (*–*)	* (1.53–*)
Delgadillo et al. (2011)	0.68 (0.55–0.79)	0.68 (0.51–0.81)	2.1 (1.3–3.4)	0.47 (0.31–0.72)	4.47 (1.93–10.3)
Inagaki et al. (2013)	0.78 (0.61–0.90)	0.85 (0.87–0.99)	17.50 (5.72–53.6)	0.22 (0.12–0.41)	77.3 (19.9–294)
Kroenke et al. (2003)	0.83 (0.68–0.93)	0.90 (0.87–0.92)	8.28 (6.2–11)	0.19 (0.1–0.37)	43.6 (18.8–101)
Liu et al. (2011)	0.64 (0.49–0.77)	0.94 (0.92–0.95)	9.98 (7.51–13.3)	0.39 (0.27–0.56)	26 (14.1–47.6)
Lowe et al. (2005)	0.87 (0.77–0.94)	0.78 (0.74–0.82)	3.96 (3.26–4.81)	0.16 (0.09–0.3)	24.4 (11.8–50)
Margrove et al. (2011)	0.8 (0.59–0.93)	1 (0.87–1)	* (0.87–1)	0.2 (0.91–0.44)	* (23.6–*)
Phelan et al. (2010)	0.63 (0.24–0.92)	0.85 (0.74–0.93)	4.24 (1.89–9.5)	0.44 (0.18–1.08)	9.63 (2.12–43.5)
Richardson et al. (2010a)	0.74 (0.67–0.79)	0.75 (0.69–0.81)	2.97 (2.31–3.82)	0.35 (0.28–0.44)	8.46 (5.51–13)
Richardson et al. (2010b)	0.80 (0.71–0.88)	0.78 (0.73–0.83)	3.63 (2.85–4.62)	0.25 (0.17–0.38)	14.3 (8.13–25)
Smith et al. (2010)	0.77 (0.46–0.95)	0.59 (0.52–0.66)	1.88 (1.33–2.64)	0.39 (0.14–1.06)	4.8 (1.37–16.6)
Thombs et al. (2008)	0.39 (0.32–0.46)	0.93 (0.91–0.95)	5.55 (4.1–7.5)	0.66 (0.59–0.73)	8.4 (0.58–12.3)
Tsai et al. (2014)	0.94 (0.72–0.99)	0.82 (0.75–0.88)	5.34 (3.7–7.7)	0.06 (0.01–0.45)	79.1 (12.7–*)
Williams et al. (2005)	0.83 (0.75–0.90)	0.84 (0.78–0.89)	5.13 (3.73–7.06)	0.20 (0.13–0.31)	25.3 (13.6–47.1)
Zhang et al. (2013)	0.79 (0.69–0.87)	0.96 (0.94–0.97)	19.9 (14.2–28.1)	0.21 (0.13–0.32)	94.6 (50.5–177)
Zuithoff et al. (2010)	0.42 (0.34–0.50)	0.94 (0.92–0.95)	6.98 (5.24–9.29)	0.62 (0.54–0.7)	11.3 (7.71–16.6)

Abbreviations: -ve LR: Negative likelihood ratio; +ve LR: Positive likelihood ratio; DOR: Diagnostic odds ratio.

\* Value could not be estimated.

**Table 4**Diagnostic test accuracy of the PHQ-2 at cut off point  $\geq 2$ .

	Sensitivity (95% CI)	Specificity (95% CI)	+ve LR (95% CI)	-ve LR (95% CI)	DOR (95% CI)
Arroll et al. (2010)	0.86 (0.80–0.91)	0.78 (0.77–0.80)	3.95 (3.58–4.35)	0.18 (0.12–0.26)	21.9 (14.0–34.3)
Chagas et al. (2011)	0.93 (0.77–0.99)	0.70 (0.58–0.79)	3.05 (2.16–4.29)	0.10 (0.03–0.39)	29.6 (7.15–*)
De Lima Osorio et al. (2009)	1 (0.94–1)	0.78 (0.70–0.86)	4.64 (3.28–6.57)	0 (*–*)	* (55.6–*)
De Lima Osorio et al. (2012)	1 (0.15–1)	0.50 (0.39–0.60)	2 (1.64–2.44)	0 (*–*)	* (50.3–*)
De Man-van Ginkel et al. (2012)	0.75 (0.50–0.91)	0.76 (0.67–0.82)	3.09 (2.1–4.53)	0.33 (0.15–0.71)	9.34 (3.27–26.50)
Fiest et al. (2014)	0.40 (0.22–0.61)	0.88 (0.82–0.92)	3.47 (1.89–6.37)	0.67 (0.48–0.92)	5.17 (2.15–12.50)
Inagaki et al. (2013)	0.78 (0.61–0.90)	0.89 (0.79–0.95)	7.50 (3.65–15.4)	0.24 (0.13–0.44)	31.1 (10.4–92.7)
Kroenke et al. (2003)	0.93 (0.80–0.99)	0.74 (0.70–0.77)	3.52 (2.98–4.15)	0.10 (0.03–0.30)	35.4 (11.4–110)
Liu et al. (2011)	0.88 (0.76–0.96)	0.82 (0.80–0.84)	4.87 (4.19–5.65)	0.15 (0.07–0.31)	33.3 (14.3–76.8)
Lowe et al. (2005)	1 (0.95–1)	0.51 (0.46–0.56)	2.04 (1.86–2.24)	0 (*–*)	* (19.2–*)
Phelan et al. (2010)	0.75 (0.35–0.97)	0.67 (0.54–0.79)	2.29 (1.34–3.92)	0.37 (0.11–1.25)	6.15 (1.28–*)
Richardson et al. (2010)	0.90 (0.85–0.93)	0.57 (0.50–0.64)	2.08 (1.77–2.45)	0.18 (0.12–0.29)	11.5 (6.98–18.8)
Richardson et al. (2010)	0.95 (88.8–0.98)	0.58 (0.52–0.64)	2.26 (1.96–2.62)	0.9 (0.04–0.20)	26.5 (10.7–65.2)
Thombs et al. (2008)	0.82 (0.77–0.87)	0.79 (0.76–0.82)	3.91 (3.37–4.53)	0.23 (0.17–0.3)	17.3 (11.8–25.3)
Tsai et al. (2014)	1 (0.81–1)	0.49 (0.41–0.58)	1.99 (1.69–2.33)	0 (*–*)	* (4.55–*)
Zhang et al. (2013)	0.96 (0.89–0.99)	0.57 (0.53–0.60)	2.24 (2.06–2.44)	0.06 (0.02–0.19)	35.8 (11.9–108)
Zuithoff et al. (2010)	0.81 (0.75–0.87)	0.76 (0.73–0.78)	3.38 (2.99–3.83)	0.25 (0.18–0.34)	13.7 (9.2–20.5)

Abbreviations: -ve LR: Negative likelihood ratio; +ve LR: Positive likelihood ratio; DOR: Diagnostic odds ratio.

\* Value could not be estimated.

(95% CI 4.44–8.18), pooled negative likelihood ratio 0.27 (95% CI 0.20–0.36) and pooled diagnostic odds ratio 22.20 (95% CI 14.00–35.19).

One of the possible reasons for heterogeneity is the various clinical settings in which the PHQ-2 has been validated. On a priori grounds we conducted subgroup analyses to examine the diagnostic performance of the PHQ-2 in similar clinical settings. As described above, of the five primary care studies one focused solely on people with epilepsy so could not be considered a general primary care sample and was excluded (Margrove et al., 2011b). A diagnostic meta-analysis was conducted for the remaining four primary care studies (Kroenke et al., 2003; Arroll et al., 2010; Zuithoff et al., 2010b; Phelan et al., 2010b); however, heterogeneity remained substantial ( $I^2=67.7\%$ ). Pooled sensitivity was 0.64 (95% CI =0.46–0.78) and pooled specificity was 0.91 (95% CI =0.89–0.93). Six studies that reported cut-off point 3 were conducted in secondary care (Osorio et al., 2012; Inagaki et al., 2013; Smith et al., 2010; Williams et al., 2005; Chagas et al., 2011; Thombs et al., 2008b). Pooled sensitivity was 0.74 (95% CI =0.57–0.86) and pooled specificity was 0.85 (95% CI =0.74–0.91). Heterogeneity

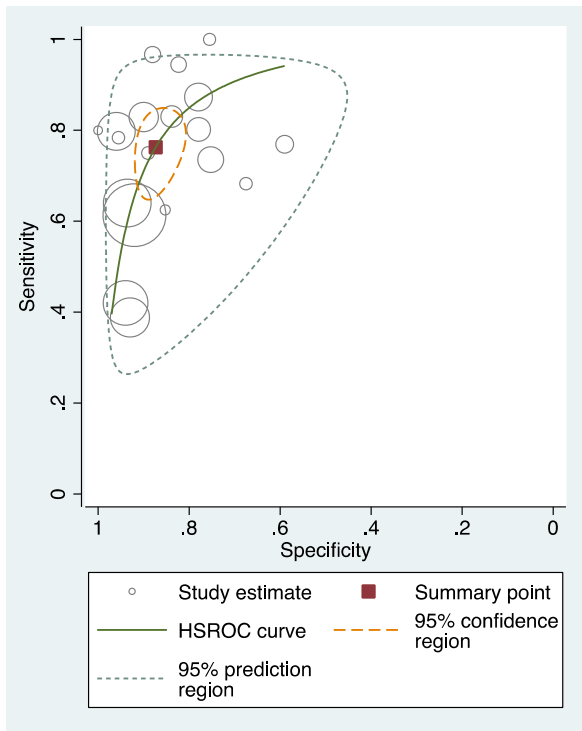
was high for this group as well ( $I^2=73.3\%$ ). We did not identify a sufficient number of studies (minimum of four studies for a diagnostic meta-analysis to be carried out in STATA) using a comparable clinical setting to conduct further subgroup analyses for other settings.

We conducted a meta-regression to further explore other possible sources of heterogeneity. Descriptive variables (setting, age, proportion female, language) were examined as predictors as were the individual quality criteria. P values were calculated using STATA metareg hand written command. None was significant at  $p < 0.05$ .

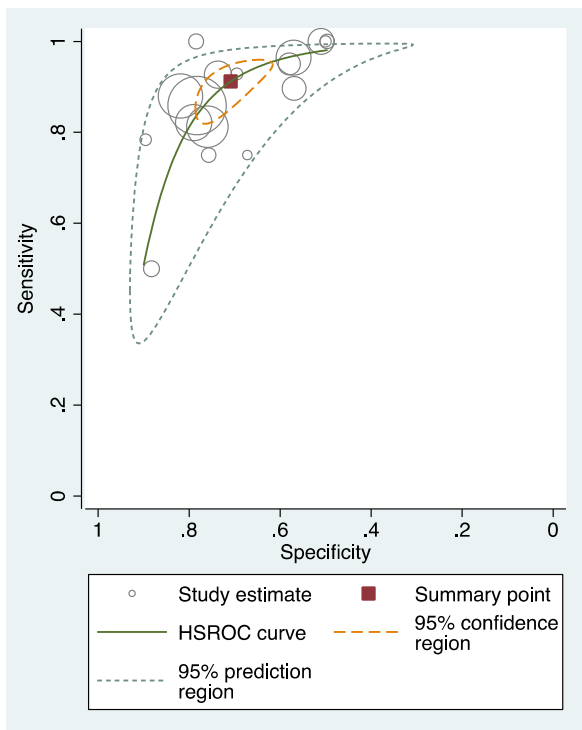
As previously mentioned, in one study (Zuithoff et al., 2010b) the reference standard assessed major depression over a one-year time-frame. Excluding this study from the meta-analyses did not significantly alter the pooled results.

An initial diagnostic meta-analysis was run for the 17 studies reporting the performance of the PHQ-2 at cut-off point  $\geq 2$ . Pooled sensitivity was 0.91 (95% CI =0.85–0.94) and pooled specificity was 0.70 (95% CI =0.64–0.76) (see Fig. 2 for sROC). Heterogeneity was moderate ( $I^2=43.5\%$ ). When the analysis was





**Fig. 2.** PHQ-2 at  $\geq 3$  summary ROC plot of diagnosis of major depressive disorder. Pooled sensitivity and specificity using a bi-variate meta-analysis.



**Fig. 3.** PHQ-2 at  $\geq 2$  summary ROC plot of diagnosis of major depressive disorder. Pooled sensitivity and specificity using a bi-variate meta-analysis.

rerun for the four primary care studies (Kroenke et al., 2003; Arroll et al., 2010; Zuithoff et al., 2010b; Phelan et al., 2010b), this gave a pooled sensitivity of 0.84 (95% CI = 0.80–0.88) and pooled specificity of 0.76 (95% CI = 0.74–0.79) (see Fig. 3 for sROC). Heterogeneity was still moderate ( $I^2=42.3\%$ ). Five studies that reported cut-off point of 2 were conducted in secondary care settings

(Osorio et al., 2012; de Man-van Ginkel et al., 2012; Fiest et al., 2014; Inagaki et al., 2013; Chagas et al., 2011). Pooled sensitivity was 0.84 (95% CI = 0.68–0.92) and pooled specificity was 0.76 (95% CI = 0.65–0.85).

Descriptive variables (setting, age, proportion female, language) and the individual quality criteria were not identified as sources of heterogeneity in meta-regression analyses for the studies that reported cut-off point 2 ( $p > 0.05$ ).

Fig. 4 uses the pooled sensitivity and specificity at cut-off  $\geq 2$  to estimate the performance of the PHQ-2 at this cut-off point as prevalence varies. The diagonal line in blue represents the prevalence of depression. The probability that a person is depressed according to the gold standard given a positive score is represented by the red line; the probability that a person is depressed given a negative score is represented by the green line.

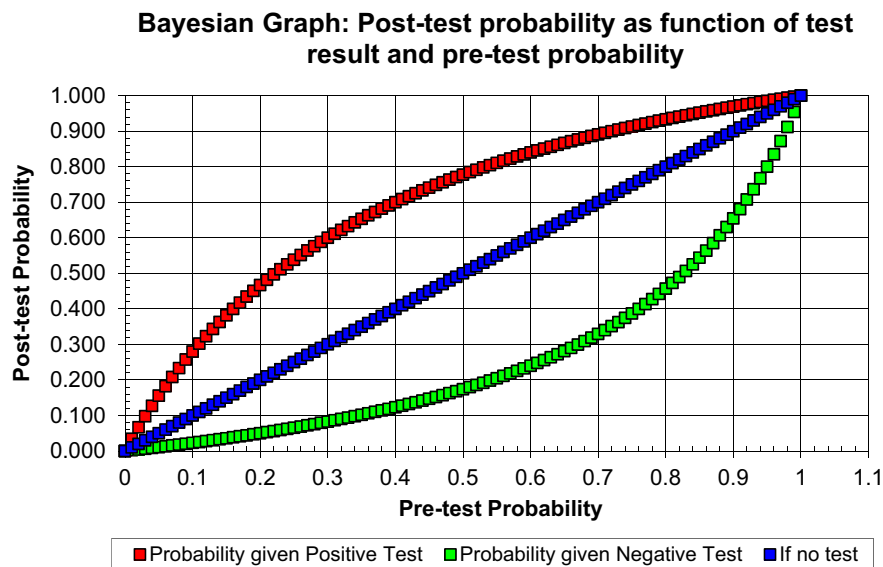
#### 4. Discussion

The original validation study of the PHQ-2 recommended a cut-off point of  $\geq 3$  on the basis of a sensitivity of 0.83 and specificity of 0.90 (Kroenke et al., 2003). This systematic review suggests that the accuracy of the PHQ-2 in identifying major depression is lower than that reported in the original study at this cut-off point. In general, sensitivity was lower than that reported in the original validation study (Kroenke et al., 2003). This, however, was not necessarily linked to the other studies reporting higher specificity, as may be expected given that sensitivity and specificity are inversely related. As a result, for those studies for which a diagnostic odds ratio could be calculated, with the exception of two studies (Inagaki et al., 2013; De Lima Osorio et al., 2009), all had a lower diagnostic odds ratio than the figure of 43.6 (95% CI = 18.8–101) calculated for Kroenke et al. (2003). There was substantial heterogeneity at  $\geq 3$ , which makes difficult the interpretation of pooled sensitivity and specificity. For the primary care studies, the sensitivity was substantially lower than Kroenke et al. (2003) (0.64 compared to 0.83 in the original validation study) and this was paired with broadly comparable levels of specificity. (0.91 compared to 0.90).

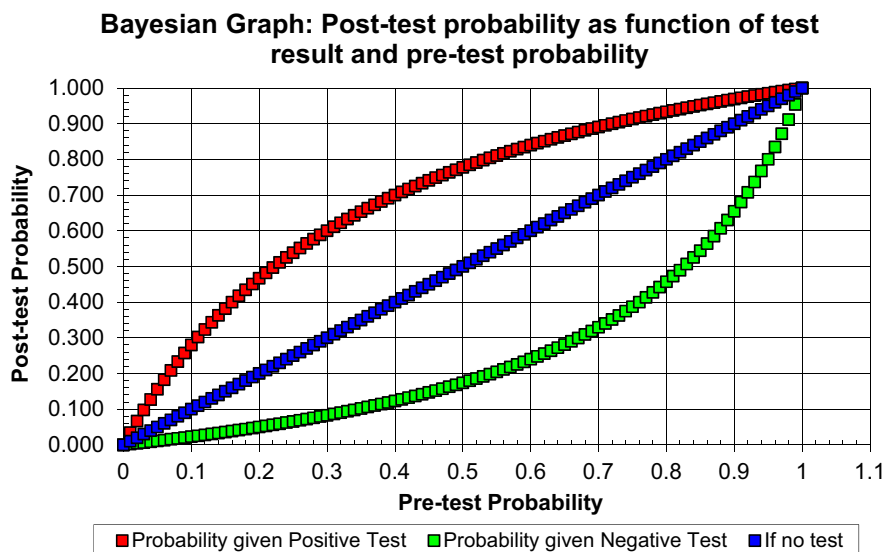
Lowering the cut-off point will increase sensitivity. Pooled sensitivity at the cut-off point of  $\geq 2$  was 0.91 (95% CI = 0.85–0.94), which is higher than the sensitivity reported in the original validation study at cut-off point  $\geq 3$ . This, however, would come at the cost of lowered specificity given its inverse relationship with sensitivity. At a cut-off point of  $\geq 2$  pooled specificity was 0.70 (95% CI = 0.64–0.76). The pooled values for the primary care samples were broadly comparable (pooled sensitivity = 0.84, 95% CI = 0.80–0.88; pooled specificity = 0.76, 95% CI = 0.74–0.79).

While the lowering of the cut-off point may limit the number of people that would be missed by the screen, it is unclear whether the level of false positives generated by this strategy would be acceptable to clinicians. The extent to which this would be a problem depends on the prevalence of depression in which the screen is being used and the cost and availability of strategies to further assess those who score positively on the initial screen.

As prevalence falls, the proportion of people who score positively but who are not depressed will increase. Prevalence estimates from the studies reported here vary substantially, though for some of the higher estimates this is likely to be related to sampling strategies that over-selected people who were likely to be depressed (Richardson et al., 2010a; Williams et al., 2005; Margrove et al., 2011b). Some idea of the value of using a cut-off point of  $\geq 2$  can be gained by using the pooled sensitivity and specificity values to estimate the proportion of people scoring  $\geq 2$  who were in fact depressed according to the reference standard at different prevalence estimates (see Fig. 4). For illustrative



**Fig. 4.** Performance of PHQ-2 at  $\geq 2$  using pooled sensitivity and specificity at different prevalence estimates. (For interpretation of the references to color in this figure legend, the reader is referred to the web version of this article.)



**Fig. 5.** Performance of PHQ-2 at  $\geq 2$  using pooled sensitivity and specificity at different prevalence estimates in primary care studies (Gilbody et al., 2007).

purposes, prevalence values of 5%, 15% and 25% are discussed. On the basis of the pooled sensitivity and specificity values, at a 5% prevalence of depression approximately 14% of people who scored at  $\geq 2$  would be depressed according to the gold standard; at 15% prevalence the value becomes approximately 37% and at 25% prevalence the value would be 51%. The pooled sensitivity and specificity of the primary care studies at this cut-off point gives similar results (5% prevalence: 16%; 15% prevalence: 38%; 25% prevalence 54%) (see Fig. 5). This analysis assumes that no patients are being treated for depression, which is perhaps an unrealistic assumption. About half of patients are recognised without screening and in primary care and a large number are already treated. However the studies do not present sufficiently detailed data to re-run the analyses for people not known to be depressed (Thombs et al., 2011).

At the lower estimates of prevalence, this cut-off point may generate too high a proportion of people scoring positively who are not depressed to make it a useful clinical tool. This suggests that it may be of limited use as a case-finding instrument, in which all people presenting to a service, such as a general practitioner

surgery, are opportunistically screened, because in such a context the prevalence is likely to be low. As the prevalence increases, however, it may become useful. This suggests that the PHQ-2 at a cut-off point of  $\geq 2$  may be of use in screening situations in which a group known to be at high risk of depression is targeted for screening, because of the increased prevalence of depression. There are, however, a number of caveats to this conclusion. First, the studies reviewed here typically used it in a general screening context; evaluation in selective contexts would be needed to confirm its performance in these situations. Secondly, as already mentioned, the studies reviewed do not distinguish between those people who are already known to services to be depressed and those who are depressed but not known. The aim of selective screening would be to identify cases that are not already known to clinical services. The prevalence of previously unknown depression will be lower than the overall depression prevalence, which may again limit the value of any identification tool. It is also unclear how the different context of identifying only previously unidentified depression would affect the diagnostic characteristics of the measure. Thirdly, the value of a screening tool cannot be

assessed solely on the basis of its sensitivity and specificity, but can only be assessed as part of a wider evaluation that examines the effectiveness and cost-effectiveness of not only screening, but the consequences of screening in terms of treatment and the outcome of that treatment (Allaby, 2010).

While this cut-off point may have some limitations in identifying people likely to have depression when there is a low prevalence of depression, given the high false positive rate, the negative likelihood ratios for this cut-off point suggest that those people who are predicted to be not depressed according to this cut-off point are unlikely to be depressed, particularly when the prevalence of depression is low. The PHQ-2 at  $\geq 2$ , therefore, may have value in ruling out depression. Fig. 4 illustrates this for the pooled sensitivity and specificity. If the pooled sensitivity and specificity values are used, at 5% prevalence approximately 99% of people scoring below the cut-off would not be depressed; at 15% the figure is 97% and at 25% the figure is 94%. The corresponding figures based on the primary care pooled estimates of sensitivity and specificity are 99% (5% prevalence), 96% (15% prevalence) and 93% (25% prevalence) (see Fig. 5).

It is important to note that the results of this meta-analysis do not apply to the Whooley questions (also known as the 'yes/no' PHQ-2). The Whooley questions are often confused with, and referred to as, the PHQ-2. However, the relatively poor sensitivity and specificity reported for the PHQ-2 in this study does not apply to the Whooley questions. A recent diagnostic meta-analysis of the Whooley questions has shown that the Whooley questions appear to be more sensitive but less specific (Bosanquet et al., 2015).

#### 4.1. Limitations

Although we sought to review grey literature databases, we cannot rule out the possibility of publication bias. Study selection and data extraction were performed by one author, which may have also introduced bias.

Three studies (Richardson et al., 2010a; Williams et al., 2005; Margrove et al., 2011b) used a design in which participants who were more likely to be depressed were also more likely to be given the reference standard, which may have introduced a partial verification bias. The QUADAS-II assessment identified variability in study quality, with only a small number of studies rated as at low risk of bias across all domains. Variations in study quality, however, did not appear to be related to outcome according to the meta-regression for cut-off point  $\geq 3$ .

There was some lack of detail in the reporting of studies, which made it difficult to assess some of the QUADAS-2 criteria. This was particularly the case for the reporting of whether the reference standard was conducted blind to the PHQ-2. Future studies should make clear statements about the blinding of the reference standard and more generally ensure that the method is reported in sufficient detail to assess the standard QUADAS-2 criteria.

Some studies may have selectively reported cut-off points – the studies that reported the two cut-off points (2 and 3) varied. It is possible that there is a relationship between the observed performance of the PHQ-2 at a particular cut-off point and the likelihood that it is reported for a particular study. Future studies should report the performance of the PHQ-2 at all available cut-off points to protect against the possibility of selective outcome reporting. Some studies reported details of sensitivity and specificity but were excluded because we were unable to identify the additional information required to calculate the 2\*2 tables that permit the calculation of the full range of accuracy statistics. Future studies should also report sufficient information to ensure that a 2\*2 table can be reconstructed from the information reported. As described above, the role of screening is to identify previously unknown cases, yet typically the studies identified in this review do

not differentiate between previously known and previously unknown cases. It is not clear what impact restricting the analysis to previously unknown cases would have on sensitivity and specificity, but such an approach would necessarily reduce the prevalence of depression, which may affect whether the instrument is likely to be useful in a particular clinical context. Future validation studies should seek to report the diagnostic performance of the PHQ-2 in identifying previously unknown cases.

The pooled estimates should be interpreted with caution given the high level of heterogeneity. Although  $I^2$  may exaggerate heterogeneity in DTA studies, there is no clear guidance available on the best way to manage this.

Another interesting finding of this review is the relatively small number of validation studies of the PHQ-2 compared to the number of validation studies of the PHQ-9, which incorporates the PHQ-2. A recent meta-analysis of the PHQ-9 has identified 36 validation studies and most of these do not specifically report the psychometric properties of the PHQ-2.

#### 4.2. Conclusion

In screening situations, reasonably high sensitivity is often required to ensure that the screening process misses few people with the diagnosis. The original validation study of Kroenke et al. (2003) reported sensitivity of 0.83 at a cut-off point of  $\geq 3$ , but a number of subsequent studies have tended to report somewhat lower sensitivity at this cut-off point. If sensitivity comparable to that reported in the original validation study is required in a screening situation, then the lower cut-off point may be needed to ensure sufficiently high sensitivity. However, the associated specificity value at this cut-off point is modest, which may limit the usefulness of the PHQ-2 at this cut-off point to identify people likely to be depressed when the prevalence of depression is low.

#### Conflicts of interest

No authors have any conflicts of interest disclosures.

#### Acknowledgements

We would like to thank the authors of both the included and excluded studies for their help in answering our questions about their studies. Dr Manea was supported by an NIHR Lectureship award. There was no specific funding for this study, and no funders had any role in the study design, in the collection, analysis or interpretation of data, in the writing of the manuscript or in the decision to submit the manuscript for publication.

#### Appendix A. Search terms used in Embase, MEDLINE and PsycINFO

(phq adj5 "2").ti, ab.  
 (phq adj5 abbreviate\$).ti, ab.  
 (phq adj5 brief).ti, ab.  
 (phq adj5 item\$).ti, ab.  
 (phq adj5 short\$).ti, ab.  
 (phq adj5 two).ti, ab.  
 (patient health questionnaire adj5 "2").ti, ab.  
 (patient health questionnaire adj5 abbreviate\$).ti, ab.  
 (patient health questionnaire adj5 brief).ti, ab.  
 (patient health questionnaire adj5 item\$).ti, ab.  
 (patient health questionnaire adj5 short\$).ti, ab.  
 (patient health questionnaire adj5 two).ti, ab.  
 (prime md adj5 "2").ti, ab.

(prime md adj5 abbreviate\$.ti, ab.  
 (prime md adj5 brief).ti, ab.  
 (prime md adj5 item\$.ti, ab.  
 (prime md adj5 short\$.ti, ab.  
 (prime md adj5 two).ti, ab.

## Appendix B. Excluded studies and reasons for exclusion

see Table B1.

**Table B1**  
 Excluded studies and reasons for exclusion.

Study	Reason for exclusion	Further information
Allgaier et al. (2012)	Reference standard not solely major depression	If either of the two questions were scored as positive, the test was considered positive.
Baker-Glenn et al. (2011)	Non-standard PHQ-2 scoring	
Boyle et al. (2011)	Overlap in sample	Overlap with Richardson et al. (2010a, 2010b)
Brody et al. (n.d.)	Not PHQ-2	From description of the measure, it is not clear that it is the PHQ-2
Bunevicius et al. (2013)	Inadequate reference standard	
Celano et al. (2013)	Inadequate reference standard	
Chen et al. (2010)	Insufficient information to calculate 2*2 table	Sensitivity and specificity reported, but other information needed to calculate 2*2 table such as base rate of depression according to gold standard not reported
de Man-van Ginkel et al. (2012)	Inadequate reference standard	
Elderon et al. (2011)	Overlap in sample	Overlap with Thombs et al. (2008)
Gjerdingen et al. (2009)	Non-standard PHQ-2 scoring	PHQ-2 scored as positive if either question scored $\geq 2$
Hahn et al. (2006)	Not PHQ-2	Uses PHQ-9 not PHQ-2
Thapar et al. (2014)	PHQ-9/PHQ-2 used to detect recurrent depression	Included patients already known to have depression
Henkel et al. (2003)	Not PHQ-2	Uses PHQ-9 not PHQ-2
Henkel et al. (2004)	Insufficient information to calculate 2*2 table	Sufficient information reported to calculate 2*2 table for 'any depressive disorder' but not major depression
Henkel et al. (n.d.)	Not PHQ-2	Uses PHQ-9 not PHQ-2
Jiang and Hesser (2011)	Inadequate reference standard	PHQ-8 is treated as the reference standard. (In addition, reference standard is 'any depressive disorder' not major depression.)
Kochhar et al. (2007)	Not PHQ-2	Uses PHQ-9 not PHQ-2 (In addition, reference standard is clinician diagnosis)
Kroenke and Spitzer (2002)	Overlap in sample	Overlap with Kroenke et al. (2003)
Li et al. (2007)	Not PHQ-2	Although called PHQ-2 it uses different questions to standard PHQ-2 items
Löwe et al. (2005)	Overlap in sample	Overlap with Lowe et al. (2005)
McGuire (2011)	Reference standard not solely major depression	Reference standard diagnosis was either major or minor depression
McManus et al. (2005)	Overlap in sample	Overlap with Thombs et al. (2008)
Mitchell et al. (2009)	Not PHQ-2	Items were from the Structured Clinical Interview for DSM-IV
Mitchell et al. (2008)	Non-standard PHQ-2 scoring	PHQ-2 scored as positive if either question was scored as positive
Mitchell et al. (2010)	Non-standard PHQ-2 scoring	PHQ-2 scored as positive if either question was scored as positive

**Table B1** (continued)

Study	Reason for exclusion	Further information
Monahan et al. (2009)	Inadequate reference standard	PHQ-9 used as the reference standard
Pibernik-Okanović et al. (2009)	Reference standard not solely major depression	Reference standard diagnosis combines major depression and dysthymia
Richardson et al.	Overlap in sample	Overlap with Richardson et al. (2010a, 2010b)
Rickels et al. (2009)	Non-standard PHQ-2 scoring	Items are scored yes / no
Robison et al. (2002)	Not PHQ-2	Uses the Whooley questions not the PHQ-2
Rollman et al. (2012)	Non-standard PHQ-2 scoring	PHQ-2 scored as positive if either question was scored as positive.
Ryan et al. (2012)	Not PHQ-2	
Smolderen et al. (2011)	Inadequate reference standard	Uses a variety of case records to determine depression status
Tiffin (2011)	Overlap in sample	A review of Richardson et al. (2010a, 2010b)
Wagner et al. (2013)	Insufficient information	Only abstract available
Watson et al. (2009)	Non-standard PHQ-2 scoring	PHQ-2 scored with yes-no response (In addition, reference standard is 'any depressive disorder' not major depression.)

## References

- Allaby, M., 2010. Screening for Depression: A Report for the UK National Screening Committee (revised report). UK National Screening Committee.
- Allgaier, A.-K., et al., 2012. Screening for depression in adolescents: validity of the patient health questionnaire in pediatric care. *Depress. Anxiety* 29 (10), 906–913 (<http://www.ncbi.nlm.nih.gov/pubmed/22753313>), accessed 25.06.16.
- Arroll, B., Goodyear-Smith, F., Crengle, S., Gunn, J., Kerse, N., et al., 2010. Validation of PHQ-2 and PHQ-9 to screen for major depression in the primary care population. *Ann. Fam. Med.* 8, 348–353.
- Bosanquet, K., Bailey, D., Gilbody, S., Harden, M., Manea, L., et al., 2015. Diagnostic accuracy of the Whooley questions for the identification of depression: a diagnostic meta-analysis. *BMJ Open*, 5.
- Boyle, L.L., et al., 2011. How do the PHQ-2, the PHQ-9 perform in aging services clients with cognitive impairment? *Int. J. Geriatr. Psychiatry* 26 (9), 952–960 (<http://www.ncbi.nlm.nih.gov/pubmed/21845598>), Accessed June 25, 2016.
- Brody, D.S. et al., Identifying patients with depression in the primary care setting: a more efficient method. *Archiv. Int. Med.*, 158(22), 2469–2475. Available at: (<http://www.ncbi.nlm.nih.gov/pubmed/9855385>) (accessed 25.06.16).
- Bunevicius, A., et al., 2013. Screening for psychological distress in neurosurgical brain tumor patients using the Patient Health Questionnaire-2. *Psycho-oncology* 22 (8), 1895–1900 (<http://www.ncbi.nlm.nih.gov/pubmed/23233453>), accessed 25.06.16.
- Celano, C.M., et al., 2013. Feasibility and utility of screening for depression and anxiety disorders in patients with cardiovascular disease. *Circ. Cardiovas. Qual. Outcomes* 6 (4), 498–504 (<http://www.ncbi.nlm.nih.gov/pubmed/23759474>), accessed 25.06.16.
- Centre for Reviews and Dissemination, 2009. Systematic Reviews: CRD's Guidance for Undertaking Reviews in Health Care. University Of York, York.
- Chagas, M.H., Crippa, J.A., Loureiro, S.R., Hallak, J.E., Meneses-Gaya, C., et al., 2011. Validity of the PHQ-2 for the screening of major depression in Parkinson's disease: two questions and one important answer. *Aging Ment. Health* 15, 838–843.
- Chagas, M.H.N., Crippa, J.A.S., Loureiro, S.R., Hallak, J.E.C., de Meneses-Gaya, C., et al., 2011. Validity of the PHQ-2 for the screening of major depression in Parkinson's disease: two questions and one important answer. *Aging Ment. Health* 15, 838–843.
- Chen, S., et al., 2010. Reliability and validity of the PHQ-9 for screening late-life depression in Chinese primary care. *Int. J. Geriatr. Psychiatry* 25 (11), 1127–1133 (<http://www.ncbi.nlm.nih.gov/pubmed/20029795>), accessed 25.06.16.
- De Lima Osorio, F., Vilela Mendes, A., Crippa, J.A., Loureiro, S.R., 2009. 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* 45, 216–227.
- de Man-van Ginkel, J.M., Gooskens, F., Schepers, V.P., Schuurmans, M.J., Lindeman, E., et al., 2012. Screening for poststroke depression using the patient health questionnaire. *Nurs. Res.* 61 (5), 333–341. Available at: (accessed 10.08.15).
- Delgadillo, J., Payne, S., Gilbody, S., Godfrey, C., Gore, S., et al., 2011. How reliable is depression screening in alcohol and drug users? A validation of brief and ultra-



- brief questionnaires. *J. Affect. Disord.* 134, 266–271.
- Elderson, L., et al., 2011. Accuracy and prognostic value of American Heart Association: recommended depression screening in patients with coronary heart disease: data from the Heart and Soul Study. *Circ. Cardiovas. Qual. Outcomes* 4 (5), 533–540 (<http://www.ncbi.nlm.nih.gov/pubmed/21862720>), accessed 25.06.16.
- Fiest, K.M., Patten, S.B., Wiebe, S., Bullock, A.G.M., Maxwell, C.J., et al., 2014. Validating screening tools for depression in epilepsy. *Epilepsia* 55, 1642–1650.
- Gilbody, S., Richards, D., Brealey, S., Hewitt, C., 2007. Screening for depression in medical settings with the Patient Health Questionnaire (PHQ): a diagnostic meta-analysis. *J. Gen. Intern. Med.* 22, 1596–1602.
- Gilbody, S., Sheldon, T., House, A., 2008. Screening and case-finding instruments for depression: a meta-analysis. *Can. Med. Assoc. J.* 178, 997–1003.
- Gjerdingen, D., et al., 2009. Postpartum depression screening at well-child visits: validity of a 2-question screen and the PHQ-9. *Ann. Fam. Med.* 7 (1), 63–70 (<http://www.ncbi.nlm.nih.gov/pubmed/19139451>), Accessed June 25, 2016.
- Hahn, D., Reuter, K., Härter, M., 2006. Screening for affective and anxiety disorders in medical patients - comparison of HADS, GHQ-12 and Brief-PHQ. *Psycho-social Med.* 3, Doc09 (<http://www.ncbi.nlm.nih.gov/pubmed/19742274>), accessed 25.06.16.
- Henkel, V. et al., Use of brief depression screening tools in primary care: consideration of heterogeneity in performance in different patient groups. *General hospital psychiatry*, 26(3) 190–198. Available at: (<http://www.ncbi.nlm.nih.gov/pubmed/15121347>) (accessed 25.06.16).
- Henkel, V., et al., 2003. Identifying depression in primary care: a comparison of different methods in a prospective cohort study. *Br. Med. J. (Clinical research ed.)* 326 (7382), 200–201 (<http://www.ncbi.nlm.nih.gov/pubmed/12543837>), accessed 25.06.16.
- Henkel, V., et al., 2004. Screening for depression in primary care: will one or two items suffice? *Eur. Arch. Psychiatry Clin. Neurosci.* 254 (4), 215–223 (<http://www.ncbi.nlm.nih.gov/pubmed/15309389>), accessed 25.06.16.
- Inagaki, M., Ohtsuki, T., Yonemoto, N., Kawashima, Y., Saitoh, A., et al., 2013. Validity of the Patient Health Questionnaire (PHQ)-9 and PHQ-2 in general internal medicine primary care at a Japanese rural hospital: a cross-sectional study. *Gen. Hosp. Psychiatry* 35, 592–597.
- Jiang, Y., Hesser, J.E., 2011. A comparison of depression and mental distress indicators, Rhode Island Behavioral Risk Factor Surveillance System, 2006. *Prev. Chron. Dis.* 8 (2), A37 (<http://www.ncbi.nlm.nih.gov/pubmed/21324251>), accessed 25.06.16.
- Joffres, M., Jaramillo, A., Dickinson, J., Lewin, G., Pottie, K., et al., 2013. Recommendations on screening for depression in adults. *Can. Med. Assoc. J. (J. de l'Assoc. Med. Can.)* 185, 775–782.
- Kochhar, P., Rajadhyaksha, S., Suvarna, V., 2007. Translation and validation of brief patient health questionnaire against DSM IV as a tool to diagnose major depressive disorder in Indian patients. *J. Postgrad. Med.* 53 (2), 102 (<http://www.jpgonline.com/text.asp?2007/53/2/102/32209>), accessed 25.06.16.
- Kroenke, K., Spitzer, R.L., 2002. The PHQ-9. *Psychiatr. Ann.* 32 (9), 509–515.
- Kroenke, K., Spitzer, R.L., Williams, J.B.W., 2001. The PHQ-9: validity of a brief depression severity measure. *J. Gen. Intern. Med.* 16, 606–613.
- Kroenke, K., Spitzer, R.L., Williams, J.B., 2003. The patient health questionnaire-2: validity of a two-item depression screener. *Med. Care* 41, 1284–1292.
- Kroenke, K., Spitzer, R.L., Williams, J.B.W., Lowe, B., 2010. The patient health questionnaire somatic, anxiety, and depressive symptom scales: a systematic review. *Gen. Hosp. Psychiatry* 32, 345–359.
- Li, C., et al., 2007. Validity of the patient health questionnaire 2 (PHQ-2) in identifying major depression in older people. *J. Am. Geriatr. Soc.* 55 (4), 596–602 (<http://www.ncbi.nlm.nih.gov/pubmed/17397440>), accessed 25.06.16.
- Liu, S.I., Yeh, Z.T., Huang, H.C., Sun, F.J., Tjeng, J.J., et al., 2011. Validation of patient health questionnaire for depression screening among primary care patients in Taiwan. *Compr. Psychiatry* 52, 96–101.
- Lowe, B., Kroenke, K., Grafe, K., 2005. Detecting and monitoring depression with a two-item questionnaire (PHQ-2). *J. Psychosom. Res.* 58, 163–171.
- Löwe, B., Kroenke, K., Grafe, K., 2005. Detecting and monitoring depression with a two-item questionnaire (PHQ-2). *J. Psychosom. Res.* 58 (2), 163–171 (<http://www.ncbi.nlm.nih.gov/pubmed/15820844>), accessed 25.06.16.
- Manea, L., Gilbody, S., McMillan, D., 2012. Optimal cut-off score for diagnosing depression with the Patient Health Questionnaire (PHQ-9): a meta-analysis. *Can. Med. Assoc. J.* 184, E191–E196.
- Mann, R., Hewitt, C.E., Gilbody, S.M., 2009. Assessing the quality of diagnostic studies using psychometric instruments: applying QUADAS. *Soc. Psychiatry Psychiatr. Epidemiol.* 44, 300–307.
- Margrove, K., Mensah, S., Thapar, A., Kerr, M., 2011a. Depression screening for patients with epilepsy in a primary care setting using the patient health questionnaire-2 and the neurological disorders depression inventory for epilepsy. *Epilepsy Behav.* 21, 387–390.
- Margrove, K., Mensah, S., Thapar, A., Kerr, M., 2011b. Depression screening for patients with epilepsy in a primary care setting using the patient health questionnaire-2 and the neurological disorders depression inventory for epilepsy. *Epilepsy Behav.* 21, 387–390.
- McGuire, A.W., 2011. Depression Screening by Nurses in Hospitalized Acute Coronary Syndrome Patients.
- McManus, D., Pipkin, S.S., Whooley, M.A., 2005. Screening for depression in patients with coronary heart disease (data from the Heart and Soul Study). *The American journal of cardiology* 96 (8), 1076–1081 (<http://www.ncbi.nlm.nih.gov/pubmed/16214441>), accessed 25.06.16.
- Mitchell, A.J., Coyne, J.C.J., 2007. Do ultra-short screening instruments accurately detect depression in primary care? A pooled analysis and meta-analysis of 22 studies. *Br. J. Gen. Pract.* 57, 144–151.
- Mitchell, A.J., et al., 2008. Acceptability of common screening methods used to detect distress and related mood disorders—preferences of cancer specialists and non-specialists. *Psycho-Oncology* 17 (3), 226–236. <http://dx.doi.org/10.1002/pon.1228>, accessed 25.06.16.
- Mitchell, A.J., et al., 2009. Accuracy of specific symptoms in the diagnosis of major depressive disorder in psychiatric out-patients: data from the MIDAS project. *Psychol. Med.* 39 (07), 1107 ([http://www.journals.cambridge.org/abstract\\_S0033291708004674](http://www.journals.cambridge.org/abstract_S0033291708004674)), accessed 25.06.16.
- Mitchell, A.J., Rao, S., Vaze, A., 2010. Do primary care physicians have particular difficulty identifying late-life depression? A meta-analysis stratified by age. *Psychother. Psychosomat.* 79 (5), 285–294. <http://dx.doi.org/10.1159/000318295>, accessed 25.06.16.
- Moher, D., Liberati, A., Tetzlaff, J., Altman, D.G., The PRISMA Group, 2009. Preferred reporting items for systematic reviews and meta-analyses: the PRISMA statement. *J. Clin. Epidemiol.* 62, 1006–1012.
- Monahan, P.O., et al., 2009. Validity/reliability of PHQ-9 and PHQ-2 depression scales among adults living with HIV/AIDS in western Kenya. *J. Gen. Int. Med.* 24 (2), 189–197 (<http://www.ncbi.nlm.nih.gov/pubmed/19031037>), accessed 25.06.16.
- Osorio, F., Carvalho, A., Fracalossi, T., Crippa, J., Loureiro, E., 2012. Are two items sufficient to screen for depression within the hospital context. *Int. J. Psychiatry Med.* 44, 141–148.
- Phelan, E., Williams, B., Meeker, K., Bonn, K., Frederick, J., et al., 2010a. A study of the diagnostic accuracy of the PHQ-9 in primary care elderly. *BMC Fam. Pract.* 11, 1–9.
- Phelan, E., Williams, B., Meeker, K., Bonn, K., Frederick, J., et al., 2010b. A study of the diagnostic accuracy of the PHQ-9 in primary care elderly. *BMC Fam. Pract.* 11, 1–9.
- Pibernik-Okanović, M. et al., 2009. Screening performance of a short versus long version of the Patient Health Questionnaire-Depression in outpatients with diabetes.
- Richardson, L.P., McCauley, E., Grossman, D.C., McCarty, C.A., Richards, J., et al., 2010. Evaluation of the patient health questionnaire-9 item for detecting major depression among adolescents. *Pediatrics* 126, 1117–1123.
- Richardson, T.M., He, H., Podgorski, C., Tu, X., Conwell, Y., 2010. Screening depression aging services clients. *Am. J. Geriatr. Psychiatry* 18, 1116–1123.
- Rickels, M.R., et al., 2009. Assessment of anxiety and depression in primary care: value of a four-item questionnaire. *J. Am. Osteopath. Assoc.* 109 (4) 798–219.
- Robison, J., et al., 2002. Screening for depression in middle-aged and older puerto rican primary care patients. *J. Gerontol. Ser. A, Biol. Sci. Med. Sci.* 57 (5), M308–M314 (<http://www.ncbi.nlm.nih.gov/pubmed/11983725>), accessed 25.06.16.
- Rollman, B.L., et al., 2012. A positive 2-item Patient Health Questionnaire depression screen among hospitalized heart failure patients is associated with elevated 12-month mortality. *J. Card. Fail.* 18 (3), 238–245 (<http://www.ncbi.nlm.nih.gov/pubmed/22385945>), accessed 25.06.16.
- Ryan, D.A., et al., 2012. Sensitivity and specificity of the Distress Thermometer and a two-item depression screen (Patient Health Questionnaire-2) with a “help” question for psychological distress and psychiatric morbidity in patients with advanced cancer. *Psycho-oncology* 21 (12), 1275–1284 (<http://www.ncbi.nlm.nih.gov/pubmed/21919118>), accessed 25.06.16.
- Smith, M.V., Gotman, N., Lin, H., Yonkers, K.A., 2010. Do the PHQ-8 and the PHQ-2 accurately screen for depressive disorders in a sample of pregnant women? *Gen. Hosp. Psychiatry* 32, 544–548.
- Smolderen, K.G., et al., 2011. Real-World Lessons From the Implementation of a Depression Screening Protocol in Acute Myocardial Infarction Patients: Implications for the American Heart Association Depression Screening Advisory. *Circ.: Cardiovas. Qual. Outcomes* 4 (3), 283–292. <http://dx.doi.org/10.1161/CIRCOUTCOMES.110.960013>, accessed 25.06.16.
- Thapar, A., et al., 2014. Detecting recurrent major depressive disorder within primary care rapidly and reliably using short questionnaire measures. *The British J. Gen. Pract.: J. R. Coll. Gen. Pract.* 64 (618), e31–e37 (<http://www.ncbi.nlm.nih.gov/pubmed/24567580>), accessed 25.06.16.
- Thombs, B.D., Ziegelstein, R.C., Whooley, M.A., 2008a. Optimizing detection of major depression among patients with coronary artery disease using the patient health questionnaire: data from the heart and soul study. *J. Gen. Intern. Med.* 23, 2014–2017.
- Thombs, B.D., Ziegelstein, R.C., Whooley, M.A., 2008b. Optimizing detection of major depression among patients with coronary artery disease using the patient health questionnaire: data from the heart and soul study. *J. Gen. Intern. Med.* 23, 2014–2017.
- Thombs, B.D., Arthurs, E., El-Baalbaki, G., Meijer, A., Ziegelstein, R.C., et al., 2011. Risk of bias from inclusion of patients who already have diagnosis of or are undergoing treatment for depression in diagnostic accuracy studies of screening tools for depression: systematic review. *Br. Med. J.* 343.
- Thombs, B.D., Coyne, J.C., Cuijpers, P., de Jonge, P., Gilbody, S., et al., 2012. Rethinking recommendations for screening for depression in primary care. *Can. Med. Assoc. J.* 184, 413–418.
- Tiffin, P.A., 2011. The Patient Health Questionnaire 2-item is a rapid, sensitive and specific screening tool for identifying adolescents with major depression. *Evid.-Based Ment. Health* 13 (4). <http://dx.doi.org/10.1136/ebmh1110>, accessed 25.06.16.
- Tsai, F.J., Huang, Y.H., Liu, H.C., Huang, K.Y., Liu, S.I., 2014. Patient health questionnaire for school-based depression screening among Chinese adolescents.

- Pediatrics, 133.
- U.S. Preventive Services Task Force, 2009. Screening for depression in adults: US preventive services task force recommendation statement. *Ann. Intern. Med.* 151, 784–792.
- Wagner, L.L., et al., 2013. Screening for depression in community-based radiation oncology settings: Results from RTOG 0841. *ASCO Meeting Abstracts* 31 (15\_suppl), 9527.
- Watson, L.C., et al., 2009. Practical depression screening in residential care/assisted living: five methods compared with gold standard diagnoses. *Am. J. Geriatric Psychiatry: Off. J. Am. Assoc. Geriatr. Psychiatry* 17 (7), 556–564 (<http://www.ncbi.nlm.nih.gov/pubmed/19554670>), accessed 25.06.16.
- Whiting, P.F., Rutjes, A.W.S., Westwood, M.E., Mallett, S., Deeks, J.J., et al., 2011. QUADAS-2: a revised tool for the quality assessment of diagnostic accuracy studies. *Ann. Intern. Med.* 155, 529–536.
- Williams, L.S., Brizendine, E.J., Plue, L., Bakas, T., Tu, W., et al., 2005. Performance of the PHQ-9 as a screening tool for depression after stroke. *Stroke* 36, 635–638.
- Wittkamp, K.A., Naeije, L., Schene, A.H., Husyer, J., van Weert, H.C., 2007. Diagnostic accuracy of the mood module of the patient health questionnaire: a systematic review. *Gen. Hosp. Psychiatry* 29, 388–395.
- Zhang, Y., Ting, R., Lam, M., Lam, J., Nan, H., et al., 2013. Measuring depressive symptoms using the patient health questionnaire-9 in Hong Kong Chinese subjects with type 2 diabetes. *J. Affect. Disord.* 151, 660–666.
- Zuithoff, N.P., Vergouwe, Y., King, M., Nazareth, I., van Wezep, M.J., et al., 2010a. The patient health questionnaire-9 for detection of major depressive disorder in primary care: consequences of current thresholds in a crosssectional study. *BMC Fam. Pract.* 11.
- Zuithoff, N.P., Vergouwe, Y., King, M., Nazareth, I., van Wezep, M.J., et al., 2010b. The patient health questionnaire-9 for detection of major depressive disorder in primary care: consequences of current thresholds in a crosssectional study. *BMC Fam. Pract.* 11, 1–7.