BMJ
**open**

# EXTENDING THE USE OF PROMS IN THE NHS: USING THE OXFORD KNEE SCORE TO MONITOR THE PROGRESSION OF KNEE OSTEOARTHRITIS. A VALIDATION STUDY

SCHOLARONE™
Manuscripts

# EXTENDING THE USE OF PROMS IN THE NHS: USING THE OXFORD KNEE SCORE TO MONITOR THE PROGRESSION OF KNEE OSTEOARTHRITIS. A VALIDATION STUDY

## ABSTRACT

**Objectives** To assess the validity of the OKS for use in patients undergoing non-operative management for their knee OA within the NHS.

**Design** Observational cohort study.

**Setting** Single orthopaedic centre in England.

**Participants** 134 patients undergoing non operative management for knee OA.

**Main outcome measures** OKS, ICOAP, KOOS-PS, at baseline and three month follow up, transition item of change at three months.

**Results** The OKS summary scale and its pain and functional component subscales demonstrated good test-retest reliability (ICC 0.93, 0.91, 0.92 respectively) and measurement precision, which allows its use with groups of patients with knee OA (research/audit) and with individuals (clinical practice). The results in this study were consistent with *a priori* set hypotheses about the relationship of the OKS with other validated measures (KOOS-PS, ICOAP, SF12), which provided evidence of construct validity and responsiveness of the score and its subscales. Confirmatory Factor Analysis confirmed the structural validity of the OKS. However, there was a lack of satisfactory

evidence of structural validity for the ICOAP and KOOS. Minimal important changes, minimal important differences and the precision of the change score were calculated for the OKS, its subscales, the ICOAP and the KOOS-PS.

**Conclusions** The OKS summary scale, together with its pain and functional component subscales, have excellent measurement properties when used with patients with knee OA, undergoing non-operative treatment. This evidence provides support for the validity of the use of the OKS when used across the spectrum of knee OA disease severity, both in research and clinical practice.

## Article summary

### Article focus
This study examines the measurement properties of the OKS, an instrument designed for patients undergoing knee replacement surgery, when used in patients undergoing non-operative treatment for knee OA.

### Key messages
The OKS summary scale, and its pain and functional component subscales, were found to have acceptable evidence of measurement properties.

The findings of this study support the use of the OKS with groups of patients (in research/audit) and for individuals (in clinical practice) who are undergoing non-operative treatment for knee OA, in addition to patients undergoing total knee replacement.

### Strengths and limitations
Despite the reliability, construct validity and responsiveness of the OKS and its subscales have been proven to be satisfactory when used in patients undergoing non-operative management of knee OA, there might be a need to further verify its content validity in this extended context.

## INTRODUCTION

The Oxford Knee Score (OKS) is a widely used patient reported outcome measure

(PROM), originally developed in1998 to be used in clinical trials for assessing

the patient-perceived outcomes of knee replacement surgery. In this form it

has proven to be reliable, valid and responsive.(1, 2) The remit of the OKS

was extended in 2009, when it was adopted by the NHS PROMs programme

in England and Wales as a primary outcome measure for knee replacement

surgery.(3) Thus, OKS data are now collected on all patients undergoing knee

replacement surgery preoperatively and at 6 months post operation, in order

to monitor and benchmark the performance of health providers.

The increasing popularity of the OKS has also resulted in its being

used for different populations and contexts from that for which it was originally

developed. In particular there has been a growing interest in using the OKS in

clinical practice as a means of standardizing clinical assessment, monitoring

individual's self-reported health state across the spectrum of OA disease, and

using the scores as an aid to clinical decision making. Extending the potential

uses of PROMs in this manner has generally been highlighted as an

opportunity to achieve maximum benefit from these measures, although the

challenges of the application of such systems have also been

recognised.(4,5)

Using the OKS as a single score across the patient pathway, to aid

diagnosis, monitor progression, assist in shared decision making and

measure the outcome of intervention offers great potential for continuity of care and understanding for patients. However robust evidence is required of the score's overall validity (i.e., the consistency of its measurement properties, such as reliability), when applied in these proposed new contexts. Generally, a measure is valid when applied to populations and contexts similar to the context in which the instrument was originally developed and tested, but measurement properties may change when the measure is applied in other contexts. The fact that the OKS was developed and tested to be used in the knee OA context (albeit end stage) is justification for considering its application in people with knee OA 'in general', but evidence has not been presented demonstrating that the OKS remains as reliable, valid and responsive when used with patients who are at earlier stages of their disease management.

The aim of our study was to assess the measurement properties of the OKS when used with patients who are undergoing non operative management for knee OA, by examining its reliability, validity, responsiveness and interpretability when applied in this context.

1 **METHODS**

2     We obtained ethical approval for a prospective cohort study from a local

3 ethics committee (11/SC/005). Informed consent was obtained from all participants

4 in the study.

5

6 Study procedures and assessments

7     This study took place at an orthopaedic centre between June 2011 and

8 August 2012. Patients were eligible for inclusion if they were referred for knee

9 problems, had a confirmed diagnosis of knee OA and were enrolled in the non-

10 operative management pathway for their knee OA (as recommended by the

11 National Institute of Clinical Excellence (NICE)(6)). Treatments for patients were

12 tailored individually, taking into account patients' preferences and needs. As such,

13 they represented standard practice in the NHS. All patients who met these criteria

14 were sent an invitation letter containing information about the study, consent forms

15 and baseline questionnaires. Patients who consented to participate in the study

16 were asked to complete the OKS(2) the Intermittent and Constant Osteoarthritis

17 Pain (ICOAP)(7) the Knee Injury and Osteoarthritis Score-Physical function Short

18 form (KOOS-PS)(8) and SF12(9) patient-reported questionnaires.

19     The OKS is a 12 item questionnaire. It's item content was devised using

20 patient interviews, which addresses pain and functional impairment in relation to

21 their knee, in patients who are undergoing knee replacement surgery.(2) Likert

22 responses are recommended to be scored from 0 to 4, which are summed to

23 produce a summary score of 0 (worst) to 48 (best)(10). More recently, we

24 presented evidence (in the context of joint replacement) that supported the original

25 conceptual basis of the OKS using its composite summary scales, but which also

26  offered an option to perform additional analyses using pain and function

27  subscales.(11) The Pain Component Score (OKS-PCS) consists of items 2, 3, 7,

28  11 and 12 and the Functional Component Score (OKS-FCS) consists of items 1, 4,

29  5, 6, 8, 9 and 10. Subscale raw scores are standardized from 0 (worst) to 100

30  (best). Patients completed the OKS at baseline, 2 and 5 days (for test-retest

31  reliability) and at 3 months.

32      We asked the patients to complete the KOOS-PS and ICOAP at baseline

33  and 3 month follow up. These scores were developed to measure pain and

34  functional disability related to knee OA, and are now a recommended outcome

35  measures by the Osteoarthritis Research Society International (OARSI).

36      The KOOS-PS consists of 7 Likert-response items and was developed from

37  a longer version of the questionnaire (KOOS(12)) using Rasch analysis to measure

38  physical function in patients with various degrees of knee OA. It is scored as the

39  KOOS from 0 (best) to 4 (worst), with a summary raw score ranging from 0 to 28.

40  The score is converted to a true interval score that ranges from 0 (best) to 100

41  (worst). The ICOAP is an 11 item questionnaire whose items were informed from

42  focus groups with patents with hip or knee OA. It has two subscales that measure

43  the intermittent and constant pain with a standardized summary score ranging from

44  0 (best) to 100 (worst).

45      Patients also completed the generic SF-12, a 12-item general health

46  measure with 8 items that have Likert-type response categories and 4 items with

47  dichotomous (yes/no) response categories. The SF-12 is scored as a Physical

48  Component Summary (PCS) and Mental Component Summary (MCS) ranging

49  from 0 (worst) to 100 (best).

50    Lastly, we asked the patients to complete a transition question in regards to

51    the change they experienced from the baseline measurement: "Compared to one

52    week before your clinic visit, please indicate how much your knee problem has

53    changed?" The question had three response options: "1. My knee has got better; 2.

54    My knee has stayed the same; 3. My knee has got worse".

55    We supplemented patient reported outcome data with information on their

56    body mass index (BMI) and the degree of structural changes observed in the knee,

57    which was available from the patients' medical records. An orthopaedic surgeon

58    (LDJ) performed Kellgren-Lawrence (K-L) grading using available knee OA

59    radiographs. The degree of structural changes in the knee was classified using (K-

60    L) grading.(13) In the absence of X-rays, we assessed intra-operative

61    documentation from previous knee arthroscopy or available MRIs to examine the

62    extent of cartilage loss and confirm the diagnosis of osteoarthritis.

63

64    Statistical methods

65    The recommended minimum sample sizes for validation studies (based on

66    optimal numbers for correlations) often range from 50 to 100.(14, 15) For

67    confirmatory factor analysis (CFA) the literature agrees with a minimum sample

68    size of about 100-150 or about 10 subjects per questionnaire item.(16, 17) These

69    sample sizes are required for data analyses and should be adjusted (i.e.

70    increased) for the risk of loss to follow up. In this study we stopped recruiting when

71    the dataset enabled us to perform CFA with at least 10 subjects per item.

72    We analysed the data using SPSS version 20 and LISREL V 8.80. Baseline

73    and 3 month follow up scores were generally non-normally distributed and change

74    scores approximated to normal (except the ICOAP and the OKS-PCS). We used

75  non-parametric statistics, where appropriate. We did not use data imputation and

76  we excluded cases with missing data on analysis by analysis basis (unless

77  mentioned otherwise). We examined the following measurement properties of the

78  OKS:

79

80  **Reliability**

81  Reliability is an estimation of the consistency and stability of a measure. It

82  includes analysis of the extent to which a measure is internally consistent

83  (measured by the inter-correlation of all items) and free from measurement error.

84  We used Cronbach's alpha to assess the internal consistency of the OKS

85  summary scale and its subscales. Alpha values of at least 0.7 are recommended in

86  order to demonstrate internal consistency.(18) We calculated an intraclass

87  correlation coefficient ($ICC_{2,1}$)(19) to assess the test-retest reliability of the OKS

88  and its subscales. Minimum ICC values of 0.7 are normally considered acceptable

89  (18) although higher values are required for the use of the score applied at an

90  individual level. To inform the potential use of the OKS on the individual level, we

91  calculated the precision of individual scores at 90% CI level by multiplying the

92  standard error of measurement (SEM) by the 2-tailed z value at 90%.

93

94  **Construct validity**

95  The validity of a measure is concerned with whether a measure actually

96  measures what it purports to measure.(20, 21) The definition of validity has

97  recently been further refined as: "The degree to which accumulated evidence and

98  theory support specific interpretations of test scores entailed by proposed uses of a

99  test".(22) Construct validity of a measure is supported by the accumulation of

100  evidence obtained by testing hypotheses about the relationship that the measure

101  exhibits with other (validated) measures.(21)

102      We examined the construct validity of the OKS summary scale and its

103  subscales by testing an *a priori* set *of* hypotheses about the expected relationships

104  between the instruments at baseline:

105      (i) the OKS and the physical component summary of the SF12 (PCS-12) are

106  measuring sufficiently similar constructs (SF-PCS measures self-reported physical

107  function and the OKS measures self-reported pain and physical functioning related

108  to the knee), so the correlation between these two instruments' scales should be

109  moderate and in the same direction,

110      (ii) the correlation between the OKS and the mental component summary of

111  the SF12 (MCS-12) should be weaker than the one between the PCS-12 and OKS

112  as these two scale constructs are not considered to be related to such an extent,

113      (iii) the OKS and KOOS-PS are measuring a sufficiently similar construct

114  (the KOOS-PS measures self-reported knee function and the OKS measures self-

115  reported pain and physical functioning related to the knee) that the correlation

116  between these two measures should be strong and negative (as scores go in the

117  opposite direction),

118      (iv) the OKS and the ICOAP are measuring sufficiently similar constructs

119  (the ICOAP measures self-reported knee pain and the OKS measures self-

120  reported pain and physical functioning related to the knee) that the correlation

121  between these two measures should be strong and negative,

122      (v) the OKS-PCS should be correlated more with the ICOAP than with the

123  KOOS-PS and negatively, in each case (the OKS-PCS measures self-reported

124  knee pain as does the ICOAP),

125    (vi) the OKS-FCS should be correlated more with the KOOS-PS that the

126    ICOAP and negatively (the OKS-FCS measures self-reported knee function, as

127    does the KOOS-PS).

128    We classified correlations (r) as: r=0 to 0.29 as none/weak; r= 0.3 to 0.69 as

129    moderate; and r > 0.7 as strong.

130    **Structural validity** is one particular aspect of construct validity; it examines

131    the extent to which the dimensionality of a measure corresponds to the construct

132    (i.e. latent variable) that is supposed to be measured.(21) For instance, if a

133    measure is unidimensional (i.e. it is supposed to measure one construct, such as

134    pain) all of its items will measure the same underlying construct. We examined the

135    structural validity of the OKS by conducting Confirmatory Factor Analysis (CFA)

136    that tested the fit of the one and two factor models of the OKS to the data, using

137    LISREL V8.80 software. In line with the standard CFA testing guidelines, we

138    considered the following indices as satisfactory: a non-significant $\chi^2$ (p>0.05),

139    standardised root mean square residual (SRMR)>0.08, comparative fit index (CFI)

140    >0.95, root mean square error of approximation  (RMSEA): <0.05 close fit,

141    <0.08good fit, <0.1 satisfactory fit; RMSEA p test of close fit>0.05.(23) Additionally,

142    we used the Chi-square  ($\chi$2) difference test and Parsimonious Normed Fit Index

143    (PNFI) to compare the fit between the two models of the OKS and the ICOAP.(24)

144    We calculated the $\chi$2 difference tests by looking at the difference of $\chi$2 of two

145    models along with the difference in their degrees of freedom.  We checked the $\chi$2

146    difference, with its degrees of freedom in the  $\chi$2 distribution table. If this value is

147    statistically significant, then the model with more degrees of freedom is favoured.

148

149    **Responsiveness**

150    The ability of a measure to detect meaningful clinical change (where it has

151    occurred) over time is critical for the use and the application of a measure.(25) This

152    change might occur following an intervention, or just occur 'naturally' during a

153    period of observation. Generally, as with construct validity, responsiveness is

154    assessed by testing *a priori* hypotheses about the relationship of the changes in

155    one measure to the changes in another (validated) measure, or with reference to a

156    change in a gold standard (as with testing criterion validity). Responsiveness can

157    also be tested with reference to a transition item, where the responsiveness is

158    tested only in subjects who have reported that clinical change has occurred.

159    We used a one sample t-test (2 tailed) to assess if the changes at 3 months

160    for the OKS, its subscales (OKS-PCS and OKS-FCS), KOOS-PS and the ICOAP

161    were significantly different from 0. We constructed a Cumulative Distribution

162    Function (CDF) plot for the; (i) OKS, (ii) OKS-PCS and ICOAP, and (iii) OKS-FCS

163    and KOOS-PS to examine the proportion of individual patients who experienced

164    deterioration and improvement beyond the measurement error of the instrument at

165    the individual level and to compare the proportion of change in pain and function

166    detected by the different measures.

167    As with construct validity, we tested the responsiveness by setting *a priori*

168    hypotheses about the direction and magnitude of changes of the validated

169    comparator instruments and the OKS:

170    (i) the change scores in the OKS should correlate strongly with the change

171    scores in the KOOS-PS and ICOAP,

172    (ii) the change scores in the OKS-PCS should correlate more strongly with

173    the change scores in the ICOAP than with the change scores in the KOOS-PS,

174      (iii) change scores for the OKS-FCS should correlate more strongly with

175  change scores for the KOOS-PS than the change scores for the ICOAP.

176      All correlations should be negative.

177

178      There was a concern about the amount of overall change that can be

179  experienced as a result of such a management pathway (which included a wide

180  range of individually tailored treatments administered to a heterogeneous sample),

181  so we additionally defined the construct of change using a patient rated item of

182  change. We then used the responses to this item to calculate anchor based values

183  of minimal important change and difference.

184

185  **Interpretability**

186      Interpretability is defined as the degree to which one can assign qualitative

187  meaning to a quantitative score.(26) In clinical trials, this issue can concern the

188  question of what is considered to be a 'good', 'bad' or 'indifferent' outcome (as

189  measured by a particular criterion or score) and what is considered to be a

190  clinically relevant change. The minimum amount of change that is discerned as

191  meaningful by patients is particularly important as it affects interpretation of study

192  results.

193      We assessed the interpretability by relating the change in the PROMs

194  scores to the patient reported item of change (using an anchor based method) and

195  by relating the observed change in the score to its measurement error at the

196  individual level (using a distribution based method). Average change in the score

197  associated with the group of patients who responded with "My knee has got better"

198  on the transition item was taken as the anchor based minimal important change

199   (MIC). The difference in the change score between the groups of patients who

200   responded with "My knee has stayed the same" and "My knee has got better on

201   the global item of change was taken as the minimal important difference (MID).

202   Finally, the minimum change in the instrument that represents real change (beyond

203   measurement error) was calculated using the Minimum Detectable Change

204   ($MDC_{90}$), which was obtained by multiplying the SEM with the z-value at the 90%

205   level and the square root of two (to account for two measurement occasions).(27,

206   28)

207

208

## RESULTS

**Sample characteristics.** 137 patients were recruited in the study. 21 patients did not complete follow up questionnaires at 3 months, out of which 3 patients were listed for a surgical procedure (2 osteotomies and 1 arthroplasty) before 3 month follow-up, 7 patients no longer wanted to participate in the study and 11 were lost to follow-up.  134 patients were included in the main baseline analysis of whom 67 (50 %) were male and 67 patients were female. The mean age of patients was 59 (SD 11), which is about 10 years less than the average age of the developmental sample of the OKS. 70% of patients had information on Body Mass Index (BMI), out of whom 30% were classified as obese (BMI>30), 41% as overweight (BMI between 25 and 29.9), 29% as normal weight (BMI between 18.5 and 24.9). No one was classified as underweight. All of the patients had a diagnosis of knee osteoarthritis. 2% of the patients had Kellgren-Lawrence (KL) grading of 0 (but evidence of cartilage loss on MRI scan), 8% had K-L of 1, 43% had K-L of 2, 16% had K-L of 3, 4% had K-L of 4.  For 26% of cases, X-ray information was unavailable, of whom, 20% had their diagnosis confirmed on the basis of MRI, while 6% of patients did not have X-rays or MRIs accessible (however, these patients had the diagnosis of OA previously confirmed in the primary care setting, different trust, or in a private clinic). All patients underwent standard non-operative management of knee OA.(29)

116 (87%) out of 134 recruited patients returned the questionnaires at three month follow up.  There was no difference in age or BMI between those patients who did not respond at three months versus those who did, but baseline OKS was different between these groups. The group that did not respond had scored, on

234 average, 7.3 points lower (worse) on the OKS than responders at three months

235 (Independent samples t-test, p<0.05). A summary of the baseline scores is

236 presented in Table 1.

237

238 Table 1. Baseline scores for the OKS, its subscales (OKS-PCS and OKS-FCS),
239 ICOAP, KOOS-PS, and SF-12 physical and mental summaries (PCS-12 and MCS-
240 12).
241
242

| | N | | Mean (SD) | Median | Percentiles | |
|---|---|---|---|---|---|---|
| | Valid | Missing | | | 25 | 75 |
| OKS | 121 | 13 | 29.3 (10) | 30 | 22 | 37 |
| OKS-PCS | 123 | 11 | 57.4 (23) | 57 | 43 | 75 |
| OKS-FCS | 137 | 7 | 66.5 (22) | 70 | 50 | 85 |
| ICOAP | 124 | 10 | 37.8 (26) | 31.8 | 16 | 57 |
| KOOS-PS | 112 | 22 | 40.5 (18) | 38.6 | 32 | 49 |
| PCS-12 | 130 | 4 | 36.7 (10) | 35 | 29 | 45 |
| MCS-12 | 130 | 4 | 51 (12) | 56 | 43 | 60 |

243
244

245 **Reliability**

246 Cronbach's alpha for the 12-item OKS was 0.94, 0.88 for the OKS-FCS and

247 0.90 for the OKS-PCS. For the ICOAP and KOOS-PS, the Cronbach's alpha was

248 0.97 and 0.94 respectively. The alpha value did not change considerably if any of

249 the items were sequentially removed from the total scores.

250 Test retest reliability ICCs were 0.93 (95% CI, 0.91-0.95) for the summary

251 OKS, 0.91 (95% CI, 0.88-0.94) for the OKS-PCS and 0.92 (95% CI, 0.90-0.95) for

252 the OKS-FCS.

253 The standard error of measurement (SEM) for the summary OKS was 2.65

254 and the confidence in individual single score at 90% was ±4.4 OKS points. SEM for

255 the OKS-FCS was 6.2 with ±10.2 90% CI for individual score and the SEM for the

256 OKS-PCS was 6.9 with ±11.3 points as 90% CI for individual score (noting that the

257 OKS-PCS and the OKS-FCS are presented on a different scale than the OKS).

258    The SEM for the ICOAP was 10.1 with ±16.6 points as 90% CI for individual score.

259    We calculated the SEM for the ICOAP by using the test-retest reliability that was

260    reported in the developmental study (0.85).(30) For the KOOS-PS, this information

261    for the English version of the questionnaire was not available, so we used the test-

262    retest reliability value of 0.86 from the validation of the French version of the

263    questionnaire. The SEM for the KOOS-PS was 6.7 with ±11.1 points as 90% CI for

264    individual score.

265

266    **Construct validity**

267          **Construct validity (hypothesis-testing).** All correlations were generally

268    consistent with *a priori* hypotheses concerning the relationships of the OKS with

269    comparator instruments. Spearman's ρ between the baseline OKS, KOOS-PS,

270    ICOAP, SF12-MCS and SF-12-PCS are shown in Table 2. The OKS correlated

271    strongly with the KOOS-PS and ICOAP. The correlation between the SF12-PCS

272    and the OKS was slightly higher than expected.  As expected, the OKS was most

273    poorly related to the SF12-MCS. The OKS-PCS correlated more with ICOAP than

274    with KOOS-PS and the OKS-FCS correlated more with the KOOS-PS that with

275    ICOAP. This evidence supports convergent and divergent validity of the OKS.

276

277    Table 2: Baseline Spearman's correlations between the scores. All correlations
278    were significant at the 0.01 level (2-tailed). The number of cases with complete
279    information that allowed the calculation of the correlation coefficients is in brackets
280    for each correlation.
281

|          | OKS          | OKS-PCS      | OKS-FCS      |
|----------|--------------|--------------|--------------|
| ICOAP    | -.879 (115)  | -.884 (117)  | -.792 (121)  |
| KOOS-PS  | -.849 (106)  | -.779 (107)  | -.867 (111)  |
| PCS-12   | .648 (121)   | /            | /            |
| MCS-12   | .370 (121)   | /            | /            |

282
283

284    **Structural validity.** 122 pre-operative OKSs, 125 pre-operative ICOAP and

285    113 pre-operative KOOS-PS were available for the CFA. Fit indices of one and two

286    factor models for the OKS are presented in Table 3. Neither of the one and two

287    factor models was rejected. Fit indices favoured the 2 factor model and the

288    reduction in $\chi^2$ in the two factor model was significant ( $\chi 2 diff > 7.879$, with df=1, at

289    the a=0.005 level).

290

291    Table 3. Fit indices of one and two-factor model of the OKS.

| Factors | χ2 (p value) | df | RMSEA | 90% CI RMSEA | RMSEA p test | CFI | SRMR | PNFI |
|---------|--------------|-----|-------|--------------|--------------|-----|------|------|
| 1 | 71.32 (p=0.06) | 54 | 0.052 | 0.00-0.08 | 0.44 | 0.99 | 0.043 | 0.80 |
| 2 | 56.64 (p=0.34) | 53 | 0.024 | 0.0-0.06 | 0.83 | 1 | 0.039 | 0.79 |

292    Note. F=number of factors; 2 =chi-square; df=degrees of freedom; RMSEA=root mean square of
293    approximation; CI=confidence intervals; p-value for test of close fit (RMSEA<.05);
294    SRMR=standardized root mean square residual; CFI-comparative t index; PNFI=parsimonious
295    normed fit index.
296

297    CFA revealed that a one-factor KOOS-PS model was rejected by the $\chi 2$ test and

298    its RMSEA was above the highest acceptable threshold of an acceptable fit (0.1)

299    (Table 4). The SRMR was acceptable and CFI was on the threshold of a good fit.

300    Both one and two factor ICOAP models were rejected by the $\chi 2$ test and both

301    models had RMSEA values far above the lowest threshold of an acceptable fit.

302    However, SRMR and CFI were acceptable for both scores. There was no

303    significant reduction (at the 0.05 level) in $\chi 2$ for the 2 factor model of the ICOAP

304    ($\chi 2 diff < 3.84$, with df=1).

305

306

307

308

309

310 Table 4. Fit indices of one and two-factor model of the ICOAP and KOOS-PS.

| | $\chi^2$ (p value) | df | RMSEA | 90% CI RMSEA | RMSEA p test | CFI | SRMR | PNFI |
|---|---|---|---|---|---|---|---|---|
| ICOAP (1F) | 242.31 (p=0.00) | 44 | 0.19 | 0.17-0.22 | 0.00 | 0.95 | 0.064 | 0.75 |
| ICOAP (2F) | 228.19 (p=0.00) | 43 | 0.19 | 0.16-0.21 | 0.00 | 0.96 | 0.057 | 0.74 |
| KOOS-PS (1F) | 40.88 (p=0.00) | 14 | 0.13 | 0.09-0.18 | 0.00 | 0.98 | 0.046 | / |

311 Note. F=number of factors; 2 =chi-square; df=degrees of freedom; RMSEA=root mean square of
312 approximation; CI=confidence intervals; p-value for test of close fit (RMSEA<.05);
313 SRMR=standardized root mean square residual; CFI-comparative t index; PNFI=parsimonious
314 normed fit index.
315

316 **Responsiveness**

317 Figure 1 shows the CDF plot for the OKS. The plot demonstrates that,

318 based on the OKS summary score, approximately 20% of patients in the study

319 experienced deterioration in health state, at three month follow up, that was

320 beyond the $MDC_{90}$ of 4 points, approximately 40% of patients experienced

321 improvement and 40% of patients did not experience change beyond this value.

322 Also, about 25% of the patients experienced improvement that was beyond the

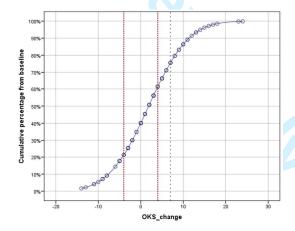323 MIC of 7 points on the OKS.



324

325 Figure 1. Cumulative percentage of patients experiencing the change on the OKS
326 from baseline less or equal to the value on the x-axis. Red line marks the minimum
327 detectable change beyond the measurement error of the score ($MDC_{90}$ of 4
328 points).
329

330        Table 5 shows the mean baseline, three month follow-up change scores,

331    and p values for the significance of 3 month change and ES for the OKS, OKS-

332    PCS, OKS-FCS, KOOS-PS and ICAOP for the overall cohort. All mean changes

333    were significant at the 0.01 level (2-tailed t-test) except the OKS-FCS.

334

335    Table 5: Significance of change in OKS, its subscales (OKS-PCS and OKS-FCS),
336    ICOAP and KOOS-PS scores at three months (one sample t-test).
337

| | N | Baseline (SD) | 3 months (SD) | Change (SD) | p-value | ES |
|---|---|---|---|---|---|---|
| OKS | 104 | 30.29 (10) | 32.15 (11) | 1.87 (7) | 0.01 | 0.19 |
| OKS-PCS | 107 | 59.36 (22) | 65.13 (24) | 5.77 (17) | <0.01 | 0.26 |
| OKS-FCS | 108 | 67.22 (21) | 68.66 (23) | 1.44 (16) | 0.4 | 0.07 |
| ICOAP[a] | 104 | 37.19 (25) | 31.53 (25) | -5.66 (19) | <0.01 | 0.23 |
| KOOS-PS[a] | 92 | 39.42 (18) | 34.88 (20) | -4.5 (14) | <0.01 | 0.25 |

338    Note. N=number of complete cases available for calculation of 3 month follow up; SD=standard
339    deviation; ES=effect size;[a] The ICOAP and the KOOS-PS represent severity of the disease in the
340    opposite direction from the OKS and its subscales.
341


342

343        The correlations between the changes in the OKS and changes in the

344    KOOS-PS and the ICOAP were somewhat less than anticipated (0.67 and 0.62

345    respectively). As hypothesized, the changes in the OKS-PCS correlated more with

346    the changes in ICOAP (also assessing knee pain) than KOOS-PS, and the

347    changes in the OKS-FCS correlated more strongly with the changes in the KOOS-

348    PS (also assessing knee function) than with the changes in the ICOAP (Table 6).

349

350    Table 6: Spearman's correlations between the 3 month changes in the OKS and its
351    subscales (OKS-PCS and OKS-FCS), ICOAP and KOOS-PS.
352

| | ICOAP | KOOS-PS |
|---|---|---|
| OKS | -.674 (96) | -.617 (87) |
| OKS-PCS | -.669 (99) | -.551 (88) |
| OKS-FCS | -.598 (100) | -.622 (90) |

353    Note. All correlations are significant at the 0.01 level (2-tailed). The number of cases with complete
354    information that allowed the calculation of the correlation coefficients is in brackets for each
355    correlation.

356
357
358 **Interpretability**

359      Tables 7 and 8 present the percentage of responses for different response

360 categories, effect sizes and mean score changes by response category. We

361 conducted independent sample t-tests for the equality of means between the mean

362 scores for groups of patients who responded 'better' and 'the same' on the

363 transition item. Only the OKS, OKS-PCS, and OKS-FCS had registered significant

364 differences between the means (2 tailed, $p < 0.05$) of groups who responded that

365 they were better/the same. Here, the OKS and OKS-PCS mean differences were

366 close to (and generally just above) scale MDC/MID values and thus likely beyond

367 measurement error, while the OKS-FCS mean differences were just less than the

368 subscale's MDC/MID values. All OKS scales' mean differences were greater than

369 the scales' relevant SEM values.

370      Table 9 presents the summary of interpretability indices.

371

372 Table 7: Number (N) and percentage of responses for different response
373 categories with effect sizes (ES), mean score changes by response category and
374 ANOVA tests for linear trend for the OKS and its subscales (OKS-PCS and OKS-
375 FCS).
376

|         |                        | Better      | Same      | Worse       |
|---------|------------------------|-------------|-----------|-------------|
| OKS     | N (% of responses)     | 30 (33)     | 26 (28)   | 36 (39)     |
|         | Mean change (SD)       | 7.1 (8)     | 0.7 (6)   | -1.88 (5)   |
|         | ES                     | .7          | .1        | -.2         |
|         | P-value for linear trend | <.001     | <.001     | <.001       |
| OKS-PCS | N (% of responses)     | 31 (33)     | 28 (30)   | 38 (35)     |
|         | Mean change (SD)       | 17.27 (19)  | 2.93 (14) | -2.68 (11)  |
|         | ES                     | .8          | .2        | -.1         |
|         | P-value for linear trend | <.001     | <.001     | <.001       |
| OKS-FCS | N (% of responses)     | 28 (33)     | 26 (31)   | 30 (36)     |
|         | Mean change (SD)       | 10.63 (14)  | 1.11 (16) | -6.35 (14)  |
|         | ES                     | .5          | .1        | -.3         |
|         | P-value for linear trend | <.001     | <.001     | <.001       |

377

378

379 Table 8: Number (N) and percentage of responses for different response
380 categories with effect sizes (ES), mean score changes by response category and
381 ANOVA tests for linear trend for the ICOAP and the KOOS-PS.
382

|  |  | Better | Same | Worse |
|---|---|---|---|---|
| ICOAP | N (% of responses) | 32 (34) | 27 (29) | 35 (37) |
|  | Mean change (SD) | -13.42 (23) | -5.64 (17) | 2.73 (16) |
|  | ES | -.6 | -.3 | .1 |
|  | P-value for linear trend | <.003 | <.003 | <.003 |
| KOOS-PS | N (% of responses) | 25 (31) | 27 (33) | 30 (37) |
|  | Mean change (SD) | -11.98 (15) | -4.22 (12) | 1.61 (12) |
|  | ES | -.8 | -.3 | .1 |
|  | P-value for linear trend | <.001 | <.001 | <.001 |

383

384 Table 9: Anchor based and distribution based MIC/MID values for the OKS, its
385 subscales, ICOAP and KOOS-PS.
386

|  | Distribution based | Anchor based | |
|---|---|---|---|
|  | $MDC_{90}$ | MID | MIC |
| OKS | ±4 | 6.4 | 7.1 |
| OKS-PCS | ±11 | 14.3 | 17.3 |
| OKS-FCS | ±10 | 9.5 | 10.6 |
| ICOAP | ±17 | 7.8 | 13.4 |
| KOOS-PS | ±11 | 7.8 | 12.0 |

387 Note. $MDC_{90}$=minimum detectable change; MID=minimum important difference; MIC=minimum
388 important change.
389

## DISCUSSION

The OKS summary scale and its pain and functional component subscales were each found to have acceptable evidence of their measurement properties to support their use with groups of patients (research/audit) and for individuals (clinical practice) who are undergoing non-operative treatment for knee OA. The OKS summary scale and its subscales were validated against the KOOS-PS, the ICOAP (measures developed for use in patients with knee OA) and the SF-12 by testing logical *a priori* hypotheses regarding the construct validity and responsiveness of the OKS and its subscales in comparison to these other (validated) measures. Thus, CFA demonstrated excellent fit and confirmed the structural validity of the OKS and both subscales. Furthermore, assessment of test-retest reliability demonstrated that the OKS and its subscales could all be used both at group and individual levels (clinical practice).(31)

The OKS subscales can be used to specifically target the improvement or deterioration in pain or function, whether in research (as an endpoint or for sample size calculations) or in clinical practice. Anchor based MIC of ≈7 for the OKS, ≈17 for the OKS-PCS, and ≈11 for the OKS-FCS can be used in cohort studies to assess if the change in the OKS (from baseline) is clinically relevant. Anchor based MID of ≈6 for the OKS, ≈14 for the OKS-PCS, and ≈10 for the OKS-FCS can be used in clinical trials to assess if the difference in change between two arms of treatment is clinically relevant. Finally, changes in individual patient scores beyond the $MDC_{90}$ (≈4 points for the OKS, ≈11 points for the OKS-PCS, and ≈10 points for the OKS-FCS) can be used as a benchmark of improvement or deterioration that is beyond the

measurement error of the score. These values are likely to be different if the OKS is used in a different population of patients (i.e. patients undergoing knee replacement surgery).

## Limitations

Even though the reliability, construct validity and responsiveness of the OKS and its subscales have been proven to be satisfactory when used in patients undergoing non-operative management for their knee OA, there might be a need to further verify its content validity in this extended context.(32) The items for the OKS were originally devised using a representative sample of patients with end stage disease, who were undergoing knee replacement surgery. It could be argued that the measure in its current form might not fully represent the concerns of this slightly different population of patients whose knee OA is generally at an earlier stage. If a measure is used in a different context or with different type of patients than that which was used  in its design/development, then the content validity may be suspect (in relation to the new/different usage).(33) On the other hand it may be assumed to be appropriate if the context is considered to be 'similar enough'.(20) Another argument is that it is unrealistic to have a new/different measure (and a new study conducted to design and test one) for every possible sub category of patient or type of treatment within all diseases or conditions. In such cases a researcher should make a judgement about the best available/closest measure, but as a minimum should  check that the measurement properties are still otherwise maintained. Any further

examination of the content validity of the OKS in this extended context would necessitate a new study (based on qualitative interviews) being undertaken.

**Comparative performance of the OKS and its subscales versus the ICOAP and the KOOS-PS in this study**

Even though the ICOAP and the KOOS-PS are currently widely used as outcome measures for knee OA, the OKS performed better in this study on several counts.

The 11-item ICOAP had a Cronbach's alpha of 0.97 (compared to the alpha of the OKS-PCS of 0.9) and the alpha was 0.94 for the KOOS-PS (compared to the alpha of 0.87 for the OKS-FCS). A high alpha value can mean that some of the items on a scale are redundant and this seems to be more of a concern for the ICOAP and KOOS-PS than for the OKS subscales. Furthermore, the reliability and precision of the score was better for the OKS and its subscales than for the KOOS-PS and ICOAP, which makes it more suitable to be used in clinical practice.

There was evidence to support both one and two factor models of the OKS, but no acceptable evidence of structural validity was found for the KOOS-PS or the ICOAP. The KOOS-PS and the one and two-factor ICOAP models were rejected by the $\chi^2$ test. Furthermore, RMSEAs were unacceptably high for both scales. The exploration of the sources of poor fit of these measures is beyond the scope of this study and future studies should investigate this problem further (perhaps also using exploratory factor analysis).

We have some concerns about the interpretability of the ICOAP and

KOOS-PS. It seems that these measures performed less well than the OKS in

this regard. First, due to the fact that the ICOAP has low precision at the

individual level (the $MDC_{90}$ is 4 points larger than the MIC) this makes it less

suitable to interpret change scores in individual patients. Second, although

around one third of the patients in our sample reported being better following

3 months of non-operative management for knee OA, neither the ICOAP or

the KOOS-PS obtained statistically significant differences in the change score

between the groups of patients who reported themselves to be better or the

same (in contrast with the OKS and its subscales). This could indicate

problems with the sensitivity of these scores to change. Third, whilst there

was some lack of symmetry between the mean change in the OKS score and

its subscales in relation to the patient rated item of change (patients who

claim they had not experienced change on the global transition item, actually

experienced change as measured by the PROM), this lack of symmetry

seems to be more pronounced for the KOOS-PS and ICOAP.


**Implications for clinicians and policymakers**

In this study, we obtained evidence that supports the use of the OKS

and its pain and functional subscales in patients who are undergoing non-

operative management for their knee. When used with patients in this context,

the OKS has demonstrated evidence of validity, reliability, and

responsiveness in measuring the health state of individuals. The measure

could be used in clinical practice to monitor disease progression in individual

patients undergoing non-operative management for their knee OA, or for

hospital audit where the information from groups of patients is analysed to assess the effectiveness of current patient management pathways for treating OA in terms of health gain/deterioration.

The use of a single valid score across a patient pathway is a compelling goal when considering how to develop standardisation of patient care in the NHS. Our new evidence suggests extending the use of the OKS in the patient pathway for managing knee OA may be possible. However the practicalities and feasibility of widespread score administration need further exploration focusing on appropriate timing, frequency and method of score administration.(34) Most importantly, more work is required to understand how results of the OKS, if adopted earlier in the pathway, should be interpreted to support patients in shared decision making regarding treatment options and the influence that such routine use of the OKS might have on the quality of care that patients receive (i.e. the effect on the quality of service and influence on patients' clinical outcomes)(35).

**ACKNOWLEDGMENTS**

A copy of the OHS and OKS questionnaires and permission to use this measure can be acquired from Isis Innovation Ltd, the technology transfer company of the University of Oxford via website:

http://www.isis-innovation.com/outcomes/index.html or email:

healthoutcomes@isis.ox.ac.uk

**COMPETING INTEREST DECLARATION**

All authors have completed the Unified Competing Interest form at www.icmje.org/coi_disclosure.pdf (available on request from the corresponding author) and declare that KKH, LDJ, AJP, DJB have no financial interests that may be relevant to the submitted work. JD is one of the original inventors of the OHS and OKS. She has received consultancy payments, via Isis Innovation, in relation to work involving both questionnaires.

**CONTRIBUTIONS**

Conception and design: JD, DJB, AJP

Acquisition of data: KKH, LDJ

Analysis and interpretation of data: KKH, JD, DJB, AJP

Drafting of the article and revision it critically for important intellectual content:

KKH, JD, LDJ, DJB, AJP

Final approval of the article: KKH, JD, LDJ, DJB, AJP

All authors, external and internal, had full access to all of the data (including statistical reports and tables) in the study and can take responsibility for the integrity of the data and the accuracy of the data analysis.

**ETHICS APPROVAL**

This study obtained ethics approval from the Oxfordshire Research Ethics Committee B (11/SC/005). Informed consent was obtained from all participants in the study.

**ROLE OF THE FUNDING SOURCE**

The research was supported under the general programme of research undertaken by the Nuffield Department of Orthopaedics, Rheumatology & Musculoskeletal Sciences as a NIHR Biomedical Research Unit.

**DATA SHARING STATEMENT**

Anonymised data and statistical codes are available from the corresponding author.

## REFERENCES

1.      Dawson JF, R. Churchman, D. Verjee-Lorenz, A. Clayson, D. Oxford Knee Score (OKS) User Manual. 2010.

2.      Dawson J, Fitzpatrick R, Murray D, Carr A. Questionnaire on the perceptions of patients about total knee replacement. Journal of Bone and Joint Surgery British Volume. 1998 Jan;80(1):63-9. PubMed PMID: 9460955. Epub 1998/02/14. eng.

3.      Department of Health. Guidance of the Routine Collection of Patient Reported Outcome Measures (PROMs). In: Department of Health, editor. London2008.

4.      Devlin NJ, Appleby J. GettinG the most out of proms. King's Fund, Office of Health Economics. 2010.

5.      Valderas J, Alonso J. Patient reported outcome measures: a model-based classification system for research and clinical practice. Qual Life Res. 2008;17(9):1125-35.

6.      Conaghan PG, Dickson J, Grant RL. Guidelines: Care and management of osteoarthritis in adults: summary of NICE guidance. BMJ: British Medical Journal. 2008;336(7642):502.

7.      Hawker GA, Davis AM, French MR, Cibere J, Jordan JM, March L, et al. Development and preliminary psychometric testing of a new OA pain measure--an OARSI/OMERACT initiative. Osteoarthritis and Cartilage. 2008 Apr;16(4):409-14. PubMed PMID: 18381179. Epub 2008/04/03. eng.

8.      Perruccio AV, Stefan Lohmander L, Canizares M, Tennant A, Hawker GA, Conaghan PG, et al. The development of a short measure of physical function for knee OA KOOS-Physical Function Shortform (KOOS-PS) - an

OARSI/OMERACT initiative. Osteoarthritis and Cartilage. 2008

May;16(5):542-50. PubMed PMID: 18294869. Epub 2008/02/26. eng.

9.      Ware J, Jr., Kosinski M, Keller SD. A 12-Item Short-Form Health

Survey: construction of scales and preliminary tests of reliability and validity.

Medical Care. 1996 Mar;34(3):220-33. PubMed PMID: 8628042. Epub

1996/03/01. eng.

10.     Murray D, Fitzpatrick R, Rogers K, Pandit H, Beard D, Carr A, et al.

The use of the Oxford hip and knee scores. The Journal of bone and joint

surgery British volume. 2007;89(8):1010.

11.     Harris K, Dawson J, Doll H, Field RE, Murray DW, Fitzpatrick R, et al.

Can pain and function be distinguished in the Oxford Knee Score in a

meaningful way? An exploratory and confirmatory factor analysis. Quality of

Life Research. 2013:1-8.

12.     Roos EM, Roos HP, Lohmander LS, Ekdahl C, Beynnon BD. Knee

Injury and Osteoarthritis Outcome Score (KOOS)--development of a self-

administered outcome measure. Journal of Orthopaedic and Sports Physical

Therapy. 1998 Aug;28(2):88-96. PubMed PMID: 9699158. Epub 1998/08/12.

eng.

13.     Kellgren J, Lawrence J. Radiological assessment of osteo-arthrosis.

Ann Rheum Dis. 1957;16(4):494-502.

14.     De Vet HCW, Terwee CB, Mokkink LB, Knol DL. Measurement in

Medicine a Practical Guide Cambridge: Cambridge University Press; 2011.

Available from:

http://public.eblib.com/EBLPublic/PublicView.do?ptiID=802925.

15.     Kline P. An easy guide to factor analysis. 1993.

16.     Ding L, Velicer WF, Harlow LL. Effects of estimation methods, number of indicators per factor, and improper solutions on structural equation modeling fit indices. Structural Equation Modeling: A Multidisciplinary Journal. 1995 1995/01/01;2(2):119-43.

17.     Bentler PM, Chou CP. Practical issues in structural modeling. Sociological Methods & Research. 1987;16(1):78-117.

18.     Nunnally JC, Bernstein IH. Psychometric theory. New York: McGraw-Hill; 1994.

19.     Shrout PE, Fleiss JL. Intraclass correlations: uses in assessing rater reliability. Psychol Bull. 1979;86(2):420.

20.     Streiner DL, Norman GR. Health measurement scales: a practical guide to their development and use: Oxford University Press, USA; 2008.

21.     Mokkink LB, Terwee CB, Patrick DL, Alonso J, Stratford PW, Knol DL, et al. The COSMIN checklist for assessing the methodological quality of studies on measurement properties of health status measurement instruments: an international Delphi study. Quality of Life Research. 2010;19(4):539-49.

22.     Association AER, Association AP, Education NCoMi, Educational JCoSf, Testing P. Standards for educational and psychological testing: Amer Educational Research Assn; 1999.

23.     Browne MW, Cudeck R. Alternative ways of assessing model fit. Testing structural equation models. 1993;154:136–62.

24.     Schumacker RE, Lomax RG. A beginner's guide to structural equation modeling. Mahwah, N.J.: Lawrence Erlbaum Associates; 2004.

25.     Fitzpatrick R, Davey C, Buxton M, Jones D. Evaluating patient-based outcome measures for use in clinical trials. 1998  Contract No.: 14.

26.     Mokkink LB, Terwee CB, Patrick DL, Alonso J, Stratford PW, Knol DL, et al. The COSMIN study reached international consensus on taxonomy, terminology, and definitions of measurement properties for health-related patient-reported outcomes. Journal of clinical epidemiology. 2010;63(7):737.

27.     Beckerman H, Roebroeck M, Lankhorst G, Becher J, Bezemer P, Verbeek A. Smallest real difference, a link between reproducibility and responsiveness. Quality of Life Research. 2001;10(7):571-8.

28.     De Vet HC, Terwee CB, Ostelo RW, Beckerman H, Knol DL, Bouter LM. Minimal changes in health status questionnaires: distinction between minimally detectable change and minimally important change. Health and Quality of Life Outcomes. 2006;4(1):54.

29.     National Institute for Health and Clinical Excellence (NICE). Osteoarthritis. National clinical guideline for care and management in adults. London: Royal College of Physicans; 2008.

30.     Hawker G, Davis A, French M, Cibere J, Jordan J, March L, et al. Development and preliminary psychometric testing of a new OA pain measure–an OARSI/OMERACT initiative. Osteoarthritis and Cartilage. 2008;16(4):409-14.

31.     Charter RA, Feldt LS. Confidence intervals for true scores: Is there a correct approach? Journal of Psychoeducational Assessment. 2001;19(4):350-64.

32.     Rothman M, Burke L, Erickson P, Leidy NK, Patrick DL, Petrie CD. Use of Existing Patient-Reported Outcome (PRO) Instruments and Their

Modification: The ISPOR Good Research Practices for Evaluating and

Documenting Content Validity for the Use of Existing Instruments and Their

Modification PRO Task Force Report. Value in Health. 2009;12(8):1075-83.

33.     Nunnally, JC. Psychometric theory. New York McGraw-Hill. 1994.

34.     Dawson J, Doll H, Fitzpatrick R, Jenkinson C, Carr AJ. The routine use

of patient reported outcome measures in healthcare settings. BMJ (Clinical

research ed). 2010;340:c186.

35.     Snyder CF, Aaronson NK, Choucair AK, Elliott TE, Greenhalgh J,

Halyard MY, et al. Implementing patient-reported outcomes assessment in

clinical practice: a review of the options and considerations. Quality of Life

Research. 2012;21(8):1305-14.

STROBE Statement—Checklist of items that should be included in reports of *cohort studies*

| | Item No | Recommendation |
|---|---|---|
| **Title and abstract** **p.1 & 2** | 1 | (*a*) Indicate the study's design with a commonly used term in the title or the abstract |
| | | (*b*) Provide in the abstract an informative and balanced summary of what was done and what was found |
| **Introduction** | | |
| Background/rationale **p.3&4** | 2 | Explain the scientific background and rationale for the investigation being reported |
| Objectives **p.4** | 3 | State specific objectives, including any prespecified hypotheses |
| **Methods** | | |
| Study design **p.5** | 4 | Present key elements of study design early in the paper |
| Setting **p. 5& 6** | 5 | Describe the setting, locations, and relevant dates, including periods of recruitment, exposure, follow-up, and data collection |
| Participants **p. 5** | 6 | (*a*) Give the eligibility criteria, and the sources and methods of selection of participants. Describe methods of follow-up |
| | | (*b*) For matched studies, give matching criteria and number of exposed and unexposed |
| Variables **n/a** | 7 | Clearly define all outcomes, exposures, predictors, potential confounders, and effect modifiers. Give diagnostic criteria, if applicable |
| Data sources/ measurement **p. 5&6** | 8* | For each variable of interest, give sources of data and details of methods of assessment (measurement). Describe comparability of assessment methods if there is more than one group |
| Bias **p.12** | 9 | Describe any efforts to address potential sources of bias |
| Study size **p.7** | 10 | Explain how the study size was arrived at |
| Quantitative variables **p. 7-8** | 11 | Explain how quantitative variables were handled in the analyses. If applicable, describe which groupings were chosen and why |
| Statistical methods | 12 | **p.7-13** (*a*) Describe all statistical methods, including those used to control for confounding |
| | | **n/a** (*b*) Describe any methods used to examine subgroups and interactions |
| | | **p.8** (*c*) Explain how missing data were addressed |
| | | **p. 14/15** (*d*) If applicable, explain how loss to follow-up was addressed |
| | | **n/a** (*e*) Describe any sensitivity analyses |
| **Results** | | |
| Participants | 13* | **p.14** (a) Report numbers of individuals at each stage of study—eg numbers potentially eligible, examined for eligibility, confirmed eligible, included in the study, completing follow-up, and analysed |
| | | **p.14** (b) Give reasons for non-participation at each stage |
| | | **n/a** (c) Consider use of a flow diagram |
| Descriptive data | 14* | **p.14** (a) Give characteristics of study participants (eg demographic, clinical, social) and information on exposures and potential confounders |
| | | **p.14/15** (b) Indicate number of participants with missing data for each variable of interest |
| | | **n/a** (c) Summarise follow-up time (eg, average and total amount) |
| Outcome data | 15* | Report numbers of outcome events or summary measures over time |
| Main results | 16 | (*a*) Give unadjusted estimates and, if applicable, confounder-adjusted estimates and |

| **n/a** | | their precision (eg, 95% confidence interval). Make clear which confounders were adjusted for and why they were included |
| | | (*b*) Report category boundaries when continuous variables were categorized |
| | | (*c*) If relevant, consider translating estimates of relative risk into absolute risk for a meaningful time period |
| Other analyses **n/a** | 17 | Report other analyses done—eg analyses of subgroups and interactions, and sensitivity analyses |
| **Discussion** | | |
| Key results **p.22/23** | 18 | Summarise key results with reference to study objectives |
| Limitations **p.23** | 19 | Discuss limitations of the study, taking into account sources of potential bias or imprecision. Discuss both direction and magnitude of any potential bias |
| Interpretation **p.25/26** | 20 | Give a cautious overall interpretation of results considering objectives, limitations, multiplicity of analyses, results from similar studies, and other relevant evidence |
| Generalisability **p.22/23, 25/26** | 21 | Discuss the generalisability (external validity) of the study results |
| **Other information** | | |
| Funding **p.28** | 22 | Give the source of funding and the role of the funders for the present study and, if applicable, for the original study on which the present article is based |

*Give information separately for exposed and unexposed groups.

**Note:** An Explanation and Elaboration article discusses each checklist item and gives methodological background and published examples of transparent reporting. The STROBE checklist is best used in conjunction with this article (freely available on the Web sites of PLoS Medicine at http://www.plosmedicine.org/, Annals of Internal Medicine at http://www.annals.org/, and Epidemiology at http://www.epidem.com/). Information on the STROBE Initiative is available at http://www.strobe-statement.org.

# BMJ open

## EXTENDING THE USE OF PROMS IN THE NHS: USING THE OXFORD KNEE SCORE IN PATIENTS UNDERGOING NON-OPERATIVE MANAGEMENT FOR KNEE OSTEOARTHRITIS. A VALIDATION STUDY

**SCHOLAR**ONE™
Manuscripts

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

# EXTENDING THE USE OF PROMS IN THE NHS: USING THE OXFORD KNEE SCORE IN PATIENTS UNDERGOING NON-OPERATIVE MANAGEMENT FOR KNEE OSTEOARTHRITIS. A VALIDATION STUDY

## ABSTRACT

**Objectives** To assess the validity of the OKS for use in patients undergoing non-operative management for their knee OA within the NHS.

**Design** Observational cohort study.

**Setting** Single orthopaedic centre in England.

**Participants** 134 patients undergoing non operative management for knee OA.

**Main outcome measures** OKS, ICOAP, KOOS-PS, at baseline and three month follow up, transition item of change at three months.

**Results** The OKS summary scale and its pain and functional component subscales demonstrated good test-retest reliability (ICC 0.93, 0.91, 0.92 respectively) and measurement precision which allows its use with groups of patients with knee OA (research/audit) and with individuals (clinical practice). The results in this study were consistent with *a priori* set hypotheses about the relationship of the OKS with other validated measures (KOOS-PS, ICOAP, SF12), which provided evidence of its construct validity and responsiveness.. Confirmatory Factor Analysis confirmed the structural validity of the OKS.

However, there was a lack of satisfactory evidence of structural validity for the ICOAP and KOOS. The minimum detectable change ($MDC_{90}$) was ±6 for the OKS (±16 for the OKS-PCS, and ±15 for the OKS-FCS), Minimal important changes were $\approx$ 7 for the OKS ($\approx$17 for the OKS–PCS and $\approx$11 for the OKS–FCS) and minimal important differences were $\approx$6 for the OKS ($\approx$14 for the OKS–PCS and $\approx$10 for the OKS–FCS). These values were also calculated for the ICOAP and the KOOS-PS.

**Conclusions** The OKS summary scale, together with its pain and functional component subscales, have excellent measurement properties when used with patients with knee OA, undergoing non-operative treatment and is superior to the ICOAP and the KOOS-PS for this purpose. This evidence provides support for the validity of the use of the OKS when used across the spectrum of knee OA disease severity, both in research and clinical practice.

## Article focus

- The OKS is a widely used patient reported outcome measure that was originally developed to measure the outcomes of knee replacement surgery.

- There is a growing interest to use the OKS in clinical practice, across the spectrum of OA disease.

- The aim of this study was to assess the measurement properties of the OKS when used with (individuals and groups of) patients who are undergoing non operative management for knee OA and compare it with most commonly used measures in this population of patients.

## Key Messages

- The OKS, and its pain and functional component subscales, have acceptable evidence of its measurement properties when used in patients (individual and groups) undergoing non-operative treatment for knee OA.

- The OKS performed better than the ICOAP and the KOOS-PS (widely used outcome measures for knee OA) on several counts.

## Strength and limitations of this study

- This study has conducted a comprehensive examination of scores' measurement properties.

- There might be a need to additionally re-evaluate evidence on some of the measurement properties presented here (such as interpretability or content validity), using different methods.

- The impact of the routine use of such scores in clinical practice should also be evaluated.

programme in England and Wales as a primary outcome measure for knee replacement surgery.(3) Thus, OKS data are now collected on all patients undergoing knee replacement surgery preoperatively and at 6 months post operation, in order to monitor and benchmark the performance of health providers.

The increasing popularity of the OKS has also resulted in its being used for different populations and contexts from that for which it was originally developed. In particular there has been a growing interest in using the OKS in clinical practice as a means of standardizing clinical assessment, monitoring individual's self-reported health state across the spectrum of OA disease, and using the scores as an aid to clinical decision making. Extending the potential uses of PROMs in this manner has generally been highlighted as an opportunity to achieve maximum benefit from these measures, although the challenges of the application of such systems have also been recognised.(4, 5)

Using the OKS as a single score across the patient pathway, to aid diagnosis, monitor progression, assist in shared decision making and measure the outcome of intervention offers great potential for continuity of care and understanding for patients. However robust evidence is required of the score's overall validity (i.e., the consistency of its measurement properties, such as reliability), when applied in these proposed new contexts. Generally, a measure is valid when applied to populations and contexts similar to the context in which the instrument was originally developed and tested, but measurement properties may change when the measure is applied in other contexts. The fact that the OKS was developed and tested to be used in the

knee OA context (albeit end stage) is justification for considering its

application in people with knee OA 'in general', but evidence has not been

presented demonstrating that the OKS remains as reliable (both on an

individual and a group level), valid and responsive when used with patients

who are at earlier stages of their disease management.

The principal aim of our study was to assess the measurement

properties of the OKS when used with (individuals and groups of) patients

who are undergoing non operative management for knee OA, by examining

its reliability, validity, responsiveness, and interpretability when applied in this

context. Furthermore, we examined some of the measurement properties of

the two most commonly used measures in this population; the Intermittent and

Constant Osteoarthritis Pain (ICOAP)(6) and the Knee Injury and

Osteoarthritis Score-Physical function Short form (KOOS-PS)(7).

1 **METHODS**

2      We obtained ethical approval for a prospective cohort study from a local

3 ethics committee (11/SC/005). Informed consent was obtained from all participants

4 in the study.

5

6 Study procedures and assessments

7      This study took place at an orthopaedic centre between June 2011 and

8 August 2012. Patients were eligible for inclusion if they were referred for knee

9 problems, had a confirmed diagnosis of knee OA and were enrolled in the non-

10 operative management pathway for their knee OA (as recommended by the

11 National Institute of Clinical Excellence (NICE)(8)). Treatments for patients were

12 tailored individually, taking into account patients' preferences and needs. As such,

13 they represented standard practice in the NHS. All patients who met these criteria

14 were sent an invitation letter containing information about the study, consent forms

15 and baseline questionnaires. Patients who consented to participate in the study

16 were asked to complete the OKS(2) the ICOAP(6), the KOOS-PS(7), and the

17 SF12(9) patient-reported questionnaires.

18      The OKS is a 12 item questionnaire. It's item content was devised using

19 patient interviews, which addresses pain and functional impairment in relation to

20 their knee, in patients who are undergoing knee replacement surgery.(2) Likert

21 responses are recommended to be scored from 0 to 4, which are summed to

22 produce a summary score of 0 (worst) to 48 (best)(10). More recently, we

23 presented evidence (in the context of joint replacement) that supported the original

24 conceptual basis of the OKS using its composite summary scales, but which also

25 offered an option to perform additional analyses using pain and function

26  subscales.(11) The Pain Component Score (OKS-PCS) consists of items 2, 3, 7,

27  11 and 12 and the Functional Component Score (OKS-FCS) consists of items 1, 4,

28  5, 6, 8, 9 and 10. Subscale raw scores are standardized from 0 (worst) to 100

29  (best). Patients completed the OKS at baseline, 2 and 5 days (for test-retest

30  reliability) and at 3 months.

31      We asked the patients to complete the KOOS-PS and ICOAP at baseline

32  and 3 month follow up. These scores were developed to measure pain and

33  functional disabilities related to knee OA, and are now a recommended outcome

34  measures by the Osteoarthritis Research Society International (OARSI).

35      The KOOS-PS consists of 7 Likert-response items and was developed from

36  a longer version of the questionnaire (KOOS(12)) using Rasch analysis to measure

37  physical function in patients with various degrees of knee OA. It is scored as the

38  KOOS from 0 (best) to 4 (worst), with a summary raw score ranging from 0 to 28.

39  The score is converted to a true interval score that ranges from 0 (best) to 100

40  (worst). The ICOAP is an 11 item questionnaire whose items were informed from

41  focus groups with patents with hip or knee OA. It has two subscales that measure

42  the intermittent and constant pain with a standardized summary score ranging from

43  0 (best) to 100 (worst).

44      Patients also completed the generic SF-12, a 12-item general health

45  measure with 8 items that have Likert-type response categories and 4 items with

46  dichotomous (yes/no) response categories. The SF-12 is scored as a Physical

47  Component Summary (PCS) and Mental Component Summary (MCS) ranging

48  from 0 (worst) to 100 (best).

49      Lastly, we asked the patients to complete a transition question in regards to

50  the change they experienced from the baseline measurement: "Compared to one

51  week before your clinic visit, please indicate how much your knee problem has

52  changed?" The question had three response options: "1. My knee has got better; 2.

53  My knee has stayed the same; 3. My knee has got worse".

54  We supplemented patient reported outcome data with information on their

55  body mass index (BMI) and the degree of structural changes observed in the knee,

56  which was available from the patients' medical records. An orthopaedic surgeon

57  (LDJ) performed Kellgren-Lawrence (K-L) grading using available knee OA

58  radiographs. (13) The degree of structural changes in the knee was classified

59  using (K-L) grading. In the absence of X-rays, we assessed intra-operative

60  documentation from previous knee arthroscopy or available MRIs to examine the

61  extent of cartilage loss and confirm the diagnosis of osteoarthritis.

62

63  Statistical methods

64  The recommended minimum sample sizes for validation studies (based on

65  optimal numbers for correlations) often range from 50 to 100.(14, 15) For

66  confirmatory factor analysis (CFA) the literature agrees with a minimum sample

67  size of about 100-150 or about 10 subjects per questionnaire item.(16, 17) These

68  sample sizes are required for data analyses and should be adjusted (i.e.

69  increased) for the risk of loss to follow up. In this study we stopped recruiting when

70  the dataset enabled us to perform CFA with at least 10 subjects per item.

71  We analysed the data using SPSS version 20 and LISREL V 8.80. Baseline

72  and 3 month follow up scores were generally non-normally distributed and change

73  scores approximated to normal (except the ICOAP and the OKS-PCS). We used

74  non-parametric statistics, where appropriate. We did not use data imputation and

75  we excluded cases with missing data on analysis by analysis basis (unless

76  mentioned otherwise). We examined the following measurement properties of the

77  OKS:

78

79  **Reliability**

80      Reliability is an estimation of the consistency and stability of a measure. It

81  includes analysis of the extent to which a measure is internally consistent

82  (measured by the inter-correlation of all items) and free from measurement error.

83  We used Cronbach's alpha to assess the internal consistency of the OKS

84  summary scale and its subscales. Alpha values of at least 0.7 are recommended in

85  order to demonstrate internal consistency. (18) We calculated an intraclass

86  correlation coefficient ($ICC_{2,1}$)(19) to assess the test-retest reliability of the OKS

87  and its subscales. Minimum ICC values of 0.7 are normally considered acceptable

88  (18) although higher values are required for the use of the score applied at an

89  individual level. To inform the potential use the OKS on the individual level, we

90  calculated the precision of individual scores at 90% CI level by multiplying the

91  standard error of measurement (SEM) by the 2-tailed z value at 90%.

92

93  **Construct validity**

94      The validity of a measure is concerned with whether a measure actually

95  measures what it purports to measure.(20, 21) The definition of validity has

96  recently been further refined as: "The degree to which accumulated evidence and

97  theory support specific interpretations of test scores entailed by proposed uses of a

98  test".(22) Construct validity of a measure is supported by the accumulation of

99  evidence obtained by testing hypotheses about the relationship that the measure

100  exhibits with other (validated) measures.(20)

101    We examined the construct validity of the OKS summary scale and its

102    subscales by testing an *a priori* set *of* hypotheses about the expected relationships

103    between the instruments at baseline:

104    (i) the OKS and the physical component summary of the SF12 (PCS-12) are

105    measuring sufficiently similar constructs (SF-PCS measures self-reported physical

106    function and the OKS measures self-reported pain and physical functioning related

107    to the knee), so the correlation between these two instruments' scales should be

108    moderate and in the same direction,

109    (ii) the correlation between the OKS and the mental component summary of

110    the SF12 (MCS-12) should be weaker than the one between the PCS-12 and OKS

111    as these two scale constructs are not considered to be related to such an extent,

112    (iii) the OKS and KOOS-PS are measuring a sufficiently similar construct

113    (the KOOS-PS measures self-reported knee function and the OKS measures self-

114    reported pain and physical functioning related to the knee) that the correlation

115    between these two measures should be strong and negative (as scores go in the

116    opposite direction),

117    (iv) the OKS and the ICOAP are measuring sufficiently similar constructs

118    (the ICOAP measures self-reported knee pain and the OKS measures self-

119    reported pain and physical functioning related to the knee) that the correlation

120    between these two measures should be strong and negative,

121    (v) the OKS-PCS should be correlated more with the ICOAP than with the

122    KOOS-PS and negatively, in each case (the OKS-PCS measures self-reported

123    knee pain as does the ICOAP),

124     (vi) the OKS-FCS should be correlated more with the KOOS-PS that the

125     ICOAP and negatively (the OKS-FCS measures self-reported knee function, as

126     does the KOOS-PS).

127     We classified correlations (r) as: r=0 to 0.29 as none/weak; r= 0.3 to 0.69 as

128     moderate; and r > 0.7 as strong.

129     **Structural validity** is one particular aspect of construct validity; it examines

130     the extent to which the dimensionality of a measure corresponds to the construct

131     (i.e. latent variable) that is supposed to be measured.(20) For instance, if a

132     measure is unidimensional (i.e. it is supposed to measure one construct, such as

133     pain) all of its items will measure the same underlying construct. We examined the

134     structural validity of the OKS by conducting Confirmatory Factor Analysis (CFA)

135     that tested the fit of the one and two factor models of the OKS to the data, using

136     LISREL V8.80 software. In line with the standard CFA testing guidelines, we

137     considered the following indices as satisfactory: a non-significant $\chi^2$ (p>0.05),

138     standardised root mean square residual (SRMR)>0.08, comparative fit index (CFI)

139     >0.95, root mean square error of approximation  (RMSEA): <0.05 close fit,

140     <0.08good fit, <0.1 satisfactory fit; RMSEA p test of close fit>0.05.(23) Additionally,

141     we used the Chi-square ($\chi2$) difference test and Parsimonious Normed Fit Index

142     (PNFI) to compare the fit between the two models of the OKS and the ICOAP. (24)

143     We calculated the $\chi2$ difference tests by looking at the difference of $\chi2$ of two

144     models along with the difference in their degrees of freedom.

145

146     **Responsiveness**

147     The ability of a measure to detect meaningful clinical change (where it has

148     occurred) over time is critical for the use and the application of a measure. (25)

149 This change might occur following an intervention, or just occur 'naturally' during a

150 period of observation. Generally, as with construct validity, responsiveness is

151 assessed by testing *a priori* hypotheses about the relationship of the changes in

152 one measure to the changes in another (validated) measure, or with reference to a

153 change in a gold standard (as with testing criterion validity). Responsiveness can

154 also be tested with reference to a transition item, where the responsiveness is

155 tested only in subjects who have reported that clinical change has occurred.

156        We used a one sample t-test (2 tailed) to assess if the changes at 3 months

157 for the OKS, its subscales (OKS-PCS and OKS-FCS), KOOS-PS and the ICOAP

158 were significantly different from 0. We constructed a Cumulative Distribution

159 Function (CDF) plot for the; (i) OKS, (ii) OKS-PCS and ICOAP, and (iii) OKS-FCS

160 and KOOS-PS to examine the proportion of individual patients who experienced

161 deterioration and improvement beyond the measurement error of the instrument at

162 the individual level and to compare the proportion of change in pain and function

163 detected by the different measures.

164        As with construct validity, we tested the responsiveness by setting *a priori*

165 hypotheses about the direction and magnitude of changes of the validated

166 comparator instruments and the OKS:

167        (i) the change scores in the OKS should correlate strongly with the change

168 scores in the KOOS-PS and ICOAP,

169        (ii) the change scores in the OKS-PCS should correlate more strongly with

170 the change scores in the ICOAP than with the change scores in the KOOS-PS,

171        (iii) change scores for the OKS-FCS should correlate more strongly with

172 change scores for the KOOS-PS than the change scores for the ICOAP.

173        All correlations should be negative.

174

175      There was a concern about the amount of overall change that can be

176      experienced as a result of such a management pathway (which included a wide

177      range of individually tailored treatments administered to a heterogeneous sample),

178      so we additionally defined the construct of change using a patient rated item of

179      change. We then used the responses to this item to calculate anchor based values

180      of minimal important change and difference.

181

182      **Interpretability**

183      Interpretability is defined as the degree to which one can assign qualitative

184      meaning to a quantitative score.[20] In clinical trials, this issue can concern the

185      question of what is considered to be a 'good', 'bad' or 'indifferent' outcome (as

186      measured by a particular criterion or score) and what is considered to be a

187      clinically relevant change. The minimum amount of change that is discerned as

188      meaningful by patients is particularly important as it affects interpretation of study

189      results.

190      We assessed the interpretability by relating the change in the PROMs

191      scores to the patient reported item of change (using an anchor based method) and

192      by relating the observed change in the score to its measurement error at the

193      individual level (using a distribution based method). Average change in the score

194      associated with the group of patients who responded with "My knee has got better"

195      on the transition item was taken as the anchor based minimal important change

196      (MIC). The difference in the change score between the groups of patients who

197      responded with "My knee has stayed the same" and "My knee has got better on

198      the global item of change was taken as the minimal important difference (MID).

199 Finally, the minimum change in the instrument that represents real change (beyond

200 measurement error) was calculated using the Minimum Detectable Change

201 ($MDC_{90}$)(26, 27))

202

203

204 **RESULTS**

205

206     **Sample characteristics.** 137 patients were recruited in the study. 21

207 patients did not complete follow up questionnaires at 3 months, out of which 3

208 patients were listed for a surgical procedure (2 osteotomies and 1 arthroplasty)

209 before 3 month follow-up, 7 patients no longer wanted to participate in the study

210 and 11 were lost to follow-up.  134 patients were included in the main baseline

211 analysis of whom 67 (50 %) were male and 67 patients were female. The mean

212 age of patients was 59 (SD 11). 70% of patients had information on Body Mass

213 Index (BMI), out of whom 30% were classified as obese (BMI>30), 41% as

214 overweight (BMI between 25 and 29.9), 29% as normal weight (BMI between 18.5

215 and 24.9). No one was classified as underweight. All of the patients had a

216 diagnosis of knee osteoarthritis. 2% of the patients had Kellgren-Lawrence (KL)

217 grading of 0 (but evidence of cartilage loss on MRI scan), 8% had K-L of 1, 43%

218 had K-L of 2, 16% had K-L of 3, 4% had K-L of 4.  For 26% of cases, X-ray

219 information was unavailable, of whom, 20% had their diagnosis confirmed on the

220 basis of MRI, while 6% of patients did not have X-rays or MRIs accessible

221 (however, these patients had the diagnosis of OA previously confirmed in the

222 primary care setting, different trust, or in a private clinic). All patients underwent

223 standard non-operative management of knee OA.(8)

224     116 (87%) out of 134 recruited patients returned the questionnaires at three

225 month follow up.  There was no difference in age or BMI between those patients

226 who did not respond at three months versus those who did, but baseline OKS was

227 different between these groups. The group that did not respond had scored, on

228 average, 7.3 points lower (worse) on the OKS than responders at three months

229 (Independent samples t-test, p<0.05). A summary of the baseline scores is

230 presented in Table 1.

231

232 Table 1. Baseline scores for the OKS, its subscales (OKS-PCS and OKS-FCS),
233 ICOAP, KOOS-PS, and SF-12 physical and mental summaries (PCS-12 and MCS-
234 12).

235
236

|  | N | | Mean (SD) | Median | Percentiles | |
|---|---|---|---|---|---|---|
|  | Valid | Missing |  |  | 25 | 75 |
| OKS | 121 | 13 | 29.3 (10) | 30 | 22 | 37 |
| OKS-PCS | 123 | 11 | 57.4 (23) | 57 | 43 | 75 |
| OKS-FCS | 137 | 7 | 66.5 (22) | 70 | 50 | 85 |
| ICOAP | 124 | 10 | 37.8 (25) | 31.8 | 16 | 57 |
| KOOS-PS | 112 | 22 | 40.5 (18) | 38.6 | 32 | 49 |
| PCS-12 | 130 | 4 | 36.7 (10) | 35 | 29 | 45 |
| MCS-12 | 130 | 4 | 51 (12) | 56 | 43 | 60 |

237
238

239

240 For comparison, in the developmental study of the OKS, the median age of

241 patients undergoing knee replacement was 73 and in this study the median age

242 was 58 (mean 59). (2) There was also considerable difference in self-reported pain

243 and functional disability between the patients in the two studies. The mean

244 baseline OKS in this sample was 29, compared to the mean preoperative OKS of

245 in the developmental study sample of 17 (when transformed to the 0-48 scoring

246 system).

247

248 **Reliability**

249

250 Cronbach's alpha for the 12-item OKS was 0.94, 0.88 for the OKS-FCS and

251 0.90 for the OKS-PCS. For the ICOAP and KOOS-PS, the Cronbach's alpha was

252 0.97 and 0.94 respectively. The alpha value did not change considerably if any of

253 the items were sequentially removed from the total scores.

254     Test retest reliability ICCs were 0.93 (95% CI, 0.91-0.95) for the summary

255 OKS, 0.91 (95% CI, 0.88-0.94) for the OKS-PCS and 0.92 (95% CI, 0.90-0.95) for

256 the OKS-FCS.

257     The standard error of measurement (SEM) for the summary OKS was 2.65

258 and the confidence in individual single score at 90% was ±4.4 OKS points. SEM for

259 the OKS-FCS was 6.2 with ±10.2 90% CI for individual score and the SEM for the

260 OKS-PCS was 6.9 with ±11.3 points as 90% CI for individual score (noting that the

261 OKS-PCS and the OKS-FCS are presented on a different scale than the OKS).

262 The SEM for the ICOAP was 9.68 with ±15.9 points as 90% CI for individual score.

263 We calculated the SEM for the ICOAP by using the test-retest reliability that was

264 reported in the developmental study (0.85)(6). For the KOOS-PS, this information

265 for the English version of the questionnaire was not available, so we used the test-

266 retest reliability value of 0.86 from the validation of the French version of the

267 questionnaire. (28) The SEM for the KOOS-PS was 6.7 with ±11.1 points as 90%

268 CI for individual score.

269

270 **Construct validity**

271     **Construct validity (hypothesis-testing).** All correlations were generally

272 consistent with *a priori* hypotheses concerning the relationships of the OKS with

273 comparator instruments. Spearman's ρ between the baseline OKS, KOOS-PS,

274 ICOAP, SF12-MCS and SF-12-PCS are shown in Table 2. The OKS correlated

275 strongly with the KOOS-PS and ICOAP. The correlation between the SF12-PCS

276 and the OKS was slightly higher than expected.  As expected, the OKS was most

277 poorly related to the SF12-MCS. The OKS-PCS correlated more with ICOAP than

278 with KOOS-PS and the OKS-FCS correlated more with the KOOS-PS that with

279 ICOAP. This evidence supports convergent and divergent validity of the OKS.

280

281 Table 2: Baseline Spearman's correlations between the scores. All correlations
282 were significant at the 0.01 level (2-tailed). The number of cases with complete
283 information that allowed the calculation of the correlation coefficients is in brackets
284 for each correlation.

285

|  | OKS | OKS-PCS | OKS-FCS |
|---|---|---|---|
| ICOAP | -.879 (115) | -.884 (117) | -.792 (121) |
| KOOS-PS | -.849 (106) | -.779 (107) | -.867 (111) |
| PCS-12 | .648 (121) | / | / |
| MCS-12 | .370 (121) | / | / |

286
287

288 **Structural validity.** 122 pre-operative OKSs, 125 pre-operative ICOAP and

289 113 pre-operative KOOS-PS were available for the CFA. Fit indices of one and two

290 factor models for the OKS are presented in Table 3. Neither of the one and two

291 factor models was rejected. Fit indices favoured the 2 factor model and the

292 reduction in $\chi^2$ in the two factor model was significant ( χ2diff>7.879, with df=1, at

293 the a=0.005 level).

294

295 Table 3. Fit indices of one and two-factor model of the OKS.

| Factors | χ2 (p value) | df | RMSEA | 90% CI RMSEA | RMSEA p test | CFI | SRMR | PNFI |
|---|---|---|---|---|---|---|---|---|
| 1 | 71.32 (p=0.06) | 54 | 0.052 | 0.00-0.08 | 0.44 | 0.99 | 0.043 | 0.80 |
| 2 | 56.64 (p=0.34) | 53 | 0.024 | 0.0-0.06 | 0.83 | 1 | 0.039 | 0.79 |

296 Note. F=number of factors; 2 =chi-square; df=degrees of freedom; RMSEA=root mean square of
297 approximation; CI=confidence intervals; p-value for test of close fit (RMSEA<.05);
298 SRMR=standardized root mean square residual; CFI-comparative t index; PNFI=parsimonious
299 normed fit index.
300

301 CFA revealed that a one-factor KOOS-PS model was rejected by the  χ2 test and

302 its RMSEA was above the highest acceptable threshold of an acceptable fit (0.1)

303 (Table 4). The SRMR was acceptable and CFI was on the threshold of a good fit.

304 Both one and two factor ICOAP models were rejected by the χ2 test and both

305 models had RMSEA values far above the lowest threshold of an acceptable fit.

306 However, SRMR and CFI were acceptable for both scores. There was no

307 significant reduction (at the 0.05 level) in χ2 for the 2 factor model of the ICOAP

308 (χ2diff< 3.84, with df=1).

309

310

311

312

313

314 Table 4. Fit indices of one and two-factor model of the ICOAP and KOOS-PS.

| | $\chi^2$ (p value) | df | RMSEA | 90% CI RMSEA | RMSEA p test | CFI | SRMR | PNFI |
|---|---|---|---|---|---|---|---|---|
| ICOAP (1F) | 242.31 (p=0.00) | 44 | 0.19 | 0.17-0.22 | 0.00 | 0.95 | 0.064 | 0.75 |
| ICOAP (2F) | 228.19 (p=0.00) | 43 | 0.19 | 0.16-0.21 | 0.00 | 0.96 | 0.057 | 0.74 |
| KOOS-PS (1F) | 40.88 (p=0.00) | 14 | 0.13 | 0.09-0.18 | 0.00 | 0.98 | 0.046 | / |

315 Note. F=number of factors; 2 =chi-square; df=degrees of freedom; RMSEA=root mean square of
316 approximation; CI=confidence intervals; p-value for test of close fit (RMSEA<.05);
317 SRMR=standardized root mean square residual; CFI-comparative t index; PNFI=parsimonious
318 normed fit index.
319

320 **Responsiveness**

321 Figure 1 shows the CDF plot for the OKS. The plot demonstrates that,

322 based on the OKS summary score, approximately 15% of patients in the study

323 experienced deterioration in health state, at three month follow up, that was

324 beyond the $MDC_{90}$ of 6 points, approximately 30% of patients experienced

325 improvement and 55% of patients did not experience change beyond this value.

326 Also, slightly less than 30% of the patients experienced improvement that was

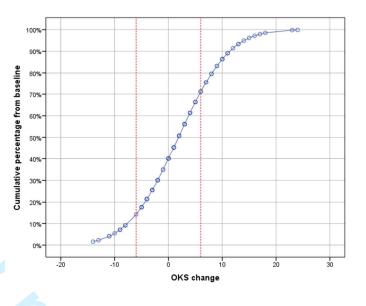327 beyond the MIC of 7 points on the OKS.

328

329 Figure 1. Cumulative percentage of patients experiencing the change on the OKS
330 from baseline less or equal to the value on the x-axis. Red line marks the minimum
331 detectable change beyond the measurement error of the score (MDC$_{90}$ of 6
332 points).
333
334 Table 5 shows the mean baseline, three month follow-up change scores,

335 and p values for the significance of 3 month change and ES for the OKS, OKS-

336 PCS, OKS-FCS, KOOS-PS and ICAOP for the overall cohort. All mean changes

337 were significant at the 0.01 level (2-tailed t-test) except the OKS-FCS.

338

339 Table 5: Significance of change in OKS, its subscales (OKS-PCS and OKS-FCS),
340 ICAOP and KOOS-PS scores at three months (one sample t-test).
341

|  | N | Baseline (SD) | 3 months (SD) | Change (SD) | p-value | ES |
|---|---|---|---|---|---|---|
| OKS | 104 | 30.29 (10) | 32.15 (11) | 1.87 (7) | 0.01 | 0.19 |
| OKS-PCS | 107 | 59.36 (22) | 65.13 (24) | 5.77 (17) | <0.01 | 0.26 |
| OKS-FCS | 108 | 67.22 (21) | 68.66 (23) | 1.44 (16) | 0.4 | 0.07 |
| ICOAP[a] | 104 | 37.19 (25) | 31.53 (25) | -5.66 (19) | <0.01 | 0.23 |
| KOOS-PS[a] | 92 | 39.42 (18) | 34.88 (20) | -4.5 (14) | <0.01 | 0.25 |

342 Note. N=number of complete cases available for calculation of 3 month follow up; SD=standard
343 deviation; ES=effect size;[a] The ICOAP and the KOOS-PS represent severity of the disease in the
344 opposite direction from the OKS and its subscales.
345

346

347    The correlations between the changes in the OKS and changes in the

348    KOOS-PS and the ICOAP were somewhat less than anticipated (0.67 and 0.62

349    respectively). As hypothesized, the changes in the OKS-PCS correlated more with

350    the changes in ICOAP (also assessing knee pain) than KOOS-PS, and the

351    changes in the OKS-FCS correlated more strongly with the changes in the KOOS-

352    PS (also assessing knee function) than with the changes in the ICOAP (Table 6).

353

354    Table 6: Spearman's correlations between the 3 month changes in the OKS and its
355    subscales (OKS-PCS and OKS-FCS), ICOAP and KOOS-PS.
356

|         | ICOAP      | KOOS-PS   |
|---------|------------|-----------|
| OKS     | -.674 (96) | -.617 (87) |
| OKS-PCS | -.669 (99) | -.551 (88) |
| OKS-FCS | -.598 (100)| -.622 (90) |

357    Note. All correlations are significant at the 0.01 level (2-tailed). The number of cases with complete
358    information that allowed the calculation of the correlation coefficients is in brackets for each
359    correlation.
360
361
362    **Interpretability**

363    Tables 7 and 8 present the percentage of responses for different response

364    categories, effect sizes and mean score changes by response category. We

365    conducted independent sample t-tests for the equality of means between the mean

366    scores for groups of patients who responded 'better' and 'the same' on the

367    transition item. Only the OKS, OKS-PCS, and OKS-FCS had registered significant

368    differences between the means (2 tailed, $p<0.05$) of groups who responded that

369    they were better/the same. Table 9 presents the summary of interpretability

370    indices.

371

372    Table 7: Number (N) and percentage of responses for different response
373    categories with effect sizes (ES), mean score changes by response category and
374    ANOVA tests for linear trend for the mean score across the three response
375    categories for the OKS and its subscales (OKS-PCS and OKS-FCS).
376
377

|  |  | Better | Same | Worse |
|---|---|---|---|---|
| OKS | N (% of responses) | 30 (33) | 26 (28) | 36 (39) |
|  | Mean change (SD) | 7.1 (8) | 0.7 (6) | -1.88 (5) |
|  | ES | .7 | .1 | -.2 |
|  | P-value for linear trend | <.001 | <.001 | <.001 |
| OKS-PCS | N (% of responses) | 31 (33) | 28 (30) | 38 (35) |
|  | Mean change (SD) | 17.27 (19) | 2.93 (14) | -2.68 (11) |
|  | ES | .8 | .2 | -.1 |
|  | P-value for linear trend | <.001 | <.001 | <.001 |
| OKS-FCS | N (% of responses) | 28 (33) | 26 (31) | 30 (36) |
|  | Mean change (SD) | 10.63 (14) | 1.11 (16) | -6.35 (14) |
|  | ES | .5 | .1 | -.3 |
|  | P-value for linear trend | <.001 | <.001 | <.001 |

378

379

380 Table 8: Number (N) and percentage of responses for different response
381 categories with effect sizes (ES), mean score changes by response category and
382 ANOVA tests for linear trend for the mean score across the three response
383 categories for the ICOAP and the KOOS-PS.
384

|  |  | Better | Same | Worse |
|---|---|---|---|---|
| ICOAP | N (% of responses) | 32 (34) | 27 (29) | 35 (37) |
|  | Mean change (SD) | -13.42 (23) | -5.64 (17) | 2.73 (16) |
|  | ES | -.6 | -.3 | .1 |
|  | P-value for linear trend | <.003 | <.003 | <.003 |
| KOOS-PS | N (% of responses) | 25 (31) | 27 (33) | 30 (37) |
|  | Mean change (SD) | -11.98 (15) | -4.22 (12) | 1.61 (12) |
|  | ES | -.8 | -.3 | .1 |
|  | P-value for linear trend | <.001 | <.001 | <.001 |

385

386 Table 9: Anchor based and distribution based MIC/MID values for the OKS, its
387 subscales, ICOAP and KOOS-PS.
388

|  | Distribution based | Anchor based | |
|---|---|---|---|
|  | $MDC_{90}$ | MID | MIC |
| OKS | ±6 | 6.4 | 7.1 |
| OKS-PCS | ±16 | 14.3 | 17.3 |
| OKS-FCS | ±15 | 9.5 | 10.6 |
| ICOAP | ±23 | 7.8 | 13.4 |
| KOOS-PS | ±16 | 7.8 | 12.0 |

389 Note. $MDC_{90}$=minimum detectable change; MID=minimum important difference; MIC=minimum
390 important change.
391

## DISCUSSION

The OKS summary scale and its pain and functional component subscales were each found to have acceptable evidence of their measurement properties to support their use with groups of patients (research/audit) and for individuals (clinical practice) who are undergoing non-operative treatment for knee OA. The OKS summary scale and its subscales were validated against the KOOS-PS, the ICOAP (measures developed for use in patients with knee OA) and the SF-12 by testing logical *a priori* hypotheses regarding the construct validity and responsiveness of the OKS and its subscales in comparison to these other (validated) measures. Thus, CFA demonstrated excellent fit and confirmed the structural validity of the OKS and both subscales. Furthermore, assessment of test-retest reliability demonstrated that the OKS and its subscales could all be used both at group and individual levels (clinical practice)(29).

The OKS subscales can be used to specifically target the improvement or deterioration in pain or function, whether in research (as an endpoint or for sample size calculations) or in clinical practice. Anchor based MIC of ≈7 for the OKS, ≈17 for the OKS-PCS, and ≈11 for the OKS-FCS can be used in cohort studies to assess if the change in the OKS (from baseline) is clinically relevant. Anchor based MID of ≈6 for the OKS, ≈14 for the OKS-PCS, and ≈10 for the OKS-FCS can be used in clinical trials to assess if the difference in change between two arms of treatment is clinically relevant. Finally, changes in individual patient scores beyond the $MDC_{90}$ (≈6 points for the OKS, ≈16 points for the OKS-PCS, and ≈15 points for the OKS-FCS) can be used as a benchmark of improvement or deterioration that is beyond the

measurement error of the score. These values are likely to be different if the OKS is used in a different population of patients (i.e. patients undergoing knee replacement surgery).

## Limitations

Even though the reliability, construct validity and responsiveness of the OKS and its subscales have been proven to be satisfactory when used in patients undergoing non-operative management for their knee OA, there might be a need to further verify its content validity in this extended context.(30) The items for the OKS were originally devised using a representative sample of patients with end stage disease, who were undergoing knee replacement surgery. It could be argued that the measure in its current form might not fully represent the concerns of this slightly different population of patients whose knee OA is generally at an earlier stage. If a measure is used in a different context or with different type of patients than that which was used  in its design/development, then the content validity may be suspect (in relation to the new/different usage).(18) A counterargument is that it is unrealistic to have a new/different measure (and a new study conducted to design and test one) for every possible sub category of patient or type of treatment within all diseases or conditions. In such cases a researcher should make a judgement about the best available/closest measure (21), but as a minimum should check that the measurement properties are still otherwise maintained. Any further examination of the content validity of the OKS in this extended context would necessitate a new study (based on qualitative interviews) being undertaken.

One of the limitations concerns the use of the transition question with three response levels (better, the same, worse). MIC/MID values depend on the number of response categories on the transition question. If, for instance, a response category 'a little better' was used instead of 'better' the final MIC value would have probably been smaller. Indeed, the methods of MIC/MID estimation have been a subject of debate within the scientific community and we would recommend that any application of the MIC/MID values presented in this paper is done with awareness of its caveats. However, regardless of the shortcomings of the transition item, the same was used in the comparative analysis of interpretability between the OKS, its subscales, the KOOS-PS and the ICOAP and in terms of drawing conclusions about the comparative performance between the scores, this is not such a source of concern.

**Comparative performance of the OKS and its subscales versus the ICOAP and the KOOS-PS in this study**

Even though the ICOAP and the KOOS-PS are currently widely used as outcome measures for knee OA, the OKS performed better in this study on several counts.

The 11-item ICOAP had a Cronbach's alpha of 0.97 (compared to the alpha of the OKS-PCS of 0.9) and the alpha was 0.94 for the KOOS-PS (compared to the alpha of 0.87 for the OKS-FCS). A high alpha value can mean that some of the items on a scale are redundant and this seems to be more of a concern for the ICOAP and KOOS-PS than for the OKS subscales. Furthermore, the reliability and precision of the score was better for the OKS

and its subscales than for the KOOS-PS and ICOAP, which makes it more suitable to be used in clinical practice.

There was evidence to support both one and two factor models of the OKS, but no acceptable evidence of structural validity was found for the KOOS-PS or the ICOAP. The KOOS-PS and the one and two-factor ICOAP models were rejected by the $\chi^2$ test. Furthermore, RMSEAs were unacceptably high for both scales. The exploration of the sources of poor fit of these measures is beyond the scope of this study and future studies should investigate this problem further (perhaps also using exploratory factor analysis).

We have some concerns about the interpretability of the ICOAP and KOOS-PS. It seems that these measures performed less well than the OKS in this regard. First, due to the fact that the ICOAP has low precision at the individual level (the $MDC_{90}$ is almost 10 points larger than the MIC) this makes it less suitable to interpret change scores in individual patients. Second, although around one third of the patients in our sample reported being better following 3 months of non-operative management for knee OA, neither the ICOAP or the KOOS-PS obtained statistically significant differences in the change score between the groups of patients who reported themselves to be better or the same (in contrast with the OKS and its subscales). This could indicate problems with the sensitivity of these scores to change. Third, whilst there was some lack of symmetry between the mean change in the OKS score and its subscales in relation to the patient rated item of change (patients who claim they had not experienced change on the global transition item,

actually experienced change as measured by the PROM), this lack of

symmetry seems to be more pronounced for the KOOS-PS and ICOAP.

**Implications for clinicians and policymakers**

In this study, we obtained evidence that supports the use of the OKS

and its pain and functional subscales in patients who are undergoing non-

operative management for their knee. When used with patients in this context,

the OKS has demonstrated evidence of validity, reliability, and

responsiveness in measuring the health state of individuals. The measure

could be used in clinical practice to monitor disease progression in individual

patients undergoing non-operative management for their knee OA, or for

hospital audit where the information from groups of patients is analysed to

assess the effectiveness of current patient management pathways for treating

OA in terms of health gain/deterioration.

Although this study was conducted on a sample of patients with knee

OA presenting themselves in the secondary care setting, we consider that the

findings presented here may be generalizable to the primary care setting.

Studies have shown no significant differences in the pain severity and function

between the groups of patients with knee OA who get referred to secondary

care and who do not. (31, 32) Other factors, such as the chronicity of the

disease, or complex interaction of psychological and social factors, are more

associated with secondary care referral. However, further research, involving

larger sample sizes, is needed to confirm these findings.

The use of a single valid score across a patient pathway is a

compelling goal when considering how to develop standardisation of patient

care in the NHS. Our new evidence suggests extending the use of the OKS in the patient pathway for managing knee OA may be possible. However the practicalities and feasibility of widespread score administration need further exploration focusing on appropriate timing, frequency and method of score administration. (33) Most importantly, more work is required to understand how results of the OKS, if adopted earlier in the pathway, should be interpreted to support patients in shared decision making regarding treatment options and the influence that such routine use of the OKS might have on the quality of care that patients receive (i.e. the effect on the quality of service and influence on patients' clinical outcomes). (34)

**ACKNOWLEDGMENTS**

A copy of the OHS and OKS questionnaires and permission to use this

measure can be acquired from Isis Innovation Ltd, the technology transfer

company of the University of Oxford via website:

http://www.isis-innovation.com/outcomes/index.html or email:

healthoutcomes@isis.ox.ac.uk

**COMPETING INTEREST DECLARATION**

All authors have completed the Unified Competing Interest form at

www.icmje.org/coi_disclosure.pdf (available on request from the

corresponding author) and declare that KKH, LDJ, AJP, DJB have no financial

interests that may be relevant to the submitted work. JD is one of the original

inventors of the OHS and OKS. She has received consultancy payments, via

Isis Innovation, in relation to work involving both questionnaires.

**CONTRIBUTIONS**

Conception and design: JD, DJB, AJP

Acquisition of data: KKH, LDJ

Analysis and interpretation of data: KKH, JD, DJB, AJP

Drafting of the article and revision it critically for important intellectual content:

KKH, JD, LDJ, DJB, AJP

Final approval of the article: KKH, JD, LDJ, DJB, AJP

All authors, external and internal, had full access to all of the data (including statistical reports and tables) in the study and can take responsibility for the integrity of the data and the accuracy of the data analysis.

**ETHICS APPROVAL**

This study obtained ethics approval from the Oxfordshire Research Ethics Committee B (11/SC/005). Informed consent was obtained from all participants in the study.

**ROLE OF THE FUNDING SOURCE**

The research was supported under the general programme of research undertaken by the Nuffield Department of Orthopaedics, Rheumatology & Musculoskeletal Sciences as a NIHR Biomedical Research Unit.

**DATA SHARING STATEMENT**

Anonymised data and statistical codes are available from the corresponding author.

**For peer review only**

# REFERENCES

1.      Dawson J, Fitzpatrick M, Churchman D, Verjee-Lorenz A, Claysonm D. User Manual for the Oxford Knee Score (OKS). 2010.

2.      Dawson J, Fitzpatrick R, Murray D, Carr A. Questionnaire on the perceptions of patients about total knee replacement. Journal of Bone and Joint Surgery British Volume. 1998 Jan;80(1):63-9. PubMed PMID: 9460955. Epub 1998/02/14. eng.

3.      Department of Health. Guidance of the Routine Collection of Patient Reported Outcome Measures (PROMs). In: Department of Health, editor. London2008.

4.      Devlin NJ, Appleby J. Getting the most out of PROMS. The Kings Fund Office of health economics. 2010.

5.      Valderas J, Alonso J. Patient reported outcome measures: a model-based classification system for research and clinical practice. Quality of Life Research. 2008;17(9):1125-35.

6.      Hawker GA, Davis AM, French MR, Cibere J, Jordan JM, March L, et al. Development and preliminary psychometric testing of a new OA pain measure--an OARSI/OMERACT initiative. Osteoarthritis and Cartilage. 2008 Apr;16(4):409-14. PubMed PMID: 18381179. Epub 2008/04/03. eng.

7.      Perruccio AV, Stefan Lohmander L, Canizares M, Tennant A, Hawker GA, Conaghan PG, et al. The development of a short measure of physical function for knee OA KOOS-Physical Function Shortform (KOOS-PS) - an OARSI/OMERACT initiative. Osteoarthritis and Cartilage. 2008 May;16(5):542-50. PubMed PMID: 18294869. Epub 2008/02/26. eng.

8.      National Institute for Health and Clinical Excellence (NICE). Osteoarthritis. National clinical guideline for care and management in adults. London: Royal College of Physicans; 2008.

9.      Ware J, Jr., Kosinski M, Keller SD. A 12-Item Short-Form Health Survey: construction of scales and preliminary tests of reliability and validity. Medical Care. 1996 Mar;34(3):220-33. PubMed PMID: 8628042. Epub 1996/03/01. eng.

10.     Murray DW, Fitzpatrick R, Rogers K, Pandit H, Beard DJ, Carr AJ, et al. The use of the Oxford hip and knee scores. Journal of Bone and Joint Surgery British Volume. 2007 August 1, 2007;89-B(8):1010-4.

11.     Harris K, Dawson J, Doll H, Field R, Murray D, Fitzpatrick R, et al. Can pain and function be distinguished in the Oxford Knee Score in a meaningful way? An exploratory and confirmatory factor analysis. Quality of Life Research. 2013 2013/03/23:1-8. English.

12. Roos EM, Roos HP, Lohmander LS, Ekdahl C, Beynnon BD. Knee Injury and Osteoarthritis Outcome Score (KOOS)--development of a self-administered outcome measure. Journal of Orthopaedic and Sports Physical Therapy. 1998 Aug;28(2):88-96. PubMed PMID: 9699158. Epub 1998/08/12. eng.

13. Kellgren J, Lawrence J. Radiological assessment of osteo-arthrosis. Annals of the Rheumatic Diseases. 1957;16(4):494-502.

14. De Vet HCW, Terwee CB, Mokkink LB, Knol DL. Measurement in Medicine a Practical Guide Cambridge: Cambridge University Press; 2011. Available from: http://public.eblib.com/EBLPublic/PublicView.do?ptiID=802925.

15. Kline P. An easy guide to factor analysis. 1993.

16. Bentler PM, Chou CP. Practical issues in structural modeling. Sociological Methods & Research. 1987;16(1):78-117.

17. Ding L, Velicer WF, Harlow LL. Effects of estimation methods, number of indicators per factor, and improper solutions on structural equation modeling fit indices. Structural Equation Modeling: A Multidisciplinary Journal. 1995 1995/01/01;2(2):119-43.

18. Nunnally JC, Bernstein IH. Psychometric theory. New York: McGraw-Hill; 1994.

19. Shrout PE, Fleiss JL. Intraclass correlations: uses in assessing rater reliability. Psychol Bull. 1979;86(2):420-8.

20. Mokkink LB, Terwee CB, Patrick DL, Alonso J, Stratford PW, Knol DL, et al. The COSMIN study reached international consensus on taxonomy, terminology, and definitions of measurement properties for health-related patient-reported outcomes. Journal of Clinical Epidemiology. 2010;63(7):737-45.

21.    Streiner DL, Norman GR. Health measurement scales: a practical guide to

their development and use: Oxford University Press, USA; 2008.

22.    AERA, APA, NCME. Standards for educational and psychological testing.

Washington: American Educational Research Association; 1999. p. 194.

23.    Browne MW, Cudeck R. Alternative ways of assessing model fit. Testing

structural equation models. 1993;154:136–62.

24.    Schumacker RE, Lomax RG. A beginner's guide to structural equation

modeling. Mahwah, N.J.: Lawrence Erlbaum Associates; 2004.

25.    Fitzpatrick R, Davey C, Buxton M, Jones D. Evaluating patient-based

outcome measures for use in clinical trials. Health Technology Assessment

1998;2(14):1-74.

26.    Beckerman H, Roebroeck M, Lankhorst G, Becher J, Bezemer P, Verbeek A.

Smallest real difference, a link between reproducibility and responsiveness. Quality of

Life Research. 2001;10(7):571-8.

27.    De Vet HC, Terwee CB, Ostelo RW, Beckerman H, Knol DL, Bouter LM.

Minimal changes in health status questionnaires: distinction between minimally

detectable change and minimally important change. Health and Quality of Life

Outcomes. 2006;4(1):54.

28.    Ornetti P, Perruccio A, Roos E, Lohmander L, Davis A, Maillefert J.

Psychometric properties of the French translation of the reduced KOOS and HOOS

(KOOS-PS and HOOS-PS). Osteoarthritis and Cartilage. 2009;17(12):1604-8.

29.    Charter RA, Feldt LS. Confidence intervals for true scores: Is there a correct

approach? Journal of Psychoeducational Assessment. 2001;19(4):350-64.

30.    Rothman M, Burke L, Erickson P, Leidy NK, Patrick DL, Petrie CD. Use of

Existing Patient-Reported Outcome (PRO) Instruments and Their Modification: The

ISPOR Good Research Practices for Evaluating and Documenting Content Validity for the Use of Existing Instruments and Their Modification PRO Task Force Report. Value in Health. 2009;12(8):1075-83.

31.    Mitchell H, Carr A, Scott D. The management of knee pain in primary care: factors associated with consulting the GP and referrals to secondary care. Rheumatology. 2006;45(6):771-6.

32.    Hopman-Rock M, De Bock GH, Bijlsma JW, Springer MP, Hofman A, Kraaimaat FW. The pattern of health care utilization of elderly people with arthritic pain in the hip or knee. International Journal for Quality in Health Care. 1997;9(2):129-37.

33.    Dawson J, Doll H, Fitzpatrick R, Jenkinson C, Carr AJ. The routine use of patient reported outcome measures in healthcare settings. BMJ (Clinical research ed). 2010;340:c186.

34.    Snyder CF, Aaronson NK, Choucair AK, Elliott TE, Greenhalgh J, Halyard MY, et al. Implementing patient-reported outcomes assessment in clinical practice: a review of the options and considerations. Quality of Life Research. 2012;21(8):1305-14.

STROBE Statement—Checklist of items that should be included in reports of *cohort studies*

| | Item No | Recommendation |
|---|---|---|
| **Title and abstract** **p.1 & 2** | 1 | (*a*) Indicate the study's design with a commonly used term in the title or the abstract |
| | | (*b*) Provide in the abstract an informative and balanced summary of what was done and what was found |
| **Introduction** | | |
| Background/rationale **p.3&4** | 2 | Explain the scientific background and rationale for the investigation being reported |
| Objectives **p.4** | 3 | State specific objectives, including any prespecified hypotheses |
| **Methods** | | |
| Study design **p.5** | 4 | Present key elements of study design early in the paper |
| Setting **p. 5& 6** | 5 | Describe the setting, locations, and relevant dates, including periods of recruitment, exposure, follow-up, and data collection |
| Participants **p. 5** | 6 | (*a*) Give the eligibility criteria, and the sources and methods of selection of participants. Describe methods of follow-up |
| | | (*b*) For matched studies, give matching criteria and number of exposed and unexposed |
| Variables **n/a** | 7 | Clearly define all outcomes, exposures, predictors, potential confounders, and effect modifiers. Give diagnostic criteria, if applicable |
| Data sources/ measurement **p. 5&6** | 8* | For each variable of interest, give sources of data and details of methods of assessment (measurement). Describe comparability of assessment methods if there is more than one group |
| Bias **p.12** | 9 | Describe any efforts to address potential sources of bias |
| Study size **p.7** | 10 | Explain how the study size was arrived at |
| Quantitative variables **p. 7-8** | 11 | Explain how quantitative variables were handled in the analyses. If applicable, describe which groupings were chosen and why |
| Statistical methods | 12 | **p.7-13** (*a*) Describe all statistical methods, including those used to control for confounding |
| | | **n/a** (*b*) Describe any methods used to examine subgroups and interactions |
| | | **p.8** (*c*) Explain how missing data were addressed |
| | | **p. 14/15** (*d*) If applicable, explain how loss to follow-up was addressed |
| | | **n/a** (*e*) Describe any sensitivity analyses |
| **Results** | | |
| Participants | 13* | **p.14** (a) Report numbers of individuals at each stage of study—eg numbers potentially eligible, examined for eligibility, confirmed eligible, included in the study, completing follow-up, and analysed |
| | | **p.14** (b) Give reasons for non-participation at each stage |
| | | **n/a** (c) Consider use of a flow diagram |
| Descriptive data | 14* | **p.14** (a) Give characteristics of study participants (eg demographic, clinical, social) and information on exposures and potential confounders |
| | | **p.14/15** (b) Indicate number of participants with missing data for each variable of interest |
| | | **n/a** (c) Summarise follow-up time (eg, average and total amount) |
| Outcome data | 15* | Report numbers of outcome events or summary measures over time |
| Main results | 16 | (*a*) Give unadjusted estimates and, if applicable, confounder-adjusted estimates and |

| | | | |
|---|---|---|---|
| **n/a** | | | their precision (eg, 95% confidence interval). Make clear which confounders were adjusted for and why they were included |
| | | | (*b*) Report category boundaries when continuous variables were categorized |
| | | | (*c*) If relevant, consider translating estimates of relative risk into absolute risk for a meaningful time period |
| Other analyses **n/a** | | 17 | Report other analyses done—eg analyses of subgroups and interactions, and sensitivity analyses |
| **Discussion** | | | |
| Key results **p.22/23** | | 18 | Summarise key results with reference to study objectives |
| Limitations **p.23** | | 19 | Discuss limitations of the study, taking into account sources of potential bias or imprecision. Discuss both direction and magnitude of any potential bias |
| Interpretation **p.25/26** | | 20 | Give a cautious overall interpretation of results considering objectives, limitations, multiplicity of analyses, results from similar studies, and other relevant evidence |
| Generalisability **p.22/23, 25/26** | | 21 | Discuss the generalisability (external validity) of the study results |
| **Other information** | | | |
| Funding **p.28** | | 22 | Give the source of funding and the role of the funders for the present study and, if applicable, for the original study on which the present article is based |

\*Give information separately for exposed and unexposed groups.

**Note:** An Explanation and Elaboration article discusses each checklist item and gives methodological background and published examples of transparent reporting. The STROBE checklist is best used in conjunction with this article (freely available on the Web sites of PLoS Medicine at http://www.plosmedicine.org/, Annals of Internal Medicine at http://www.annals.org/, and Epidemiology at http://www.epidem.com/). Information on the STROBE Initiative is available at http://www.strobe-statement.org.

# EXTENDING THE USE OF PROMS IN THE NHS: USING THE OXFORD KNEE SCORE ~~TO MONITOR THE PROGRESSION OF KNEE OSTEOARTHRITIS~~IN PATIENTS UNDERGOING NON-OPERATIVE MANAGEMENT FOR KNEE OSTEOARTHRITIS. A VALIDATION STUDY

## ABSTRACT

**Objectives** To assess the validity of the OKS for use in patients undergoing non-operative management for their knee OA within the NHS.

**Design** Observational cohort study.

**Setting** Single orthopaedic centre in England.

**Participants** 134 patients undergoing non operative management for knee OA.

**Main outcome measures** OKS, ICOAP, KOOS-PS, at baseline and three month follow up, transition item of change at three months.

**Results** The OKS summary scale and its pain and functional component subscales demonstrated good test-retest reliability (ICC 0.93, 0.91, 0.92 respectively) and measurement precision which allows its use with groups of patients with knee OA (research/audit) and with individuals (clinical practice). The results in this study were consistent with *a priori* set hypotheses about the

relationship of the OKS with other validated measures (KOOS-PS, ICOAP,

SF12), which provided evidence of its construct validity and responsiveness. ~~of the score and its subscales~~. Confirmatory Factor Analysis confirmed the

structural validity of the OKS. However, there was a lack of satisfactory

evidence of structural validity for the ICOAP and KOOS. The minimum

detectable change ($MDC_{90}$) was ±6 for the OKS (±16 for the OKS-PCS, and

±15 for the OKS-FCS), Minimal important changes were ≈ 7 for the OKS

(≈17 for the OKS–PCS and ≈11 for the OKS–FCS) and minimal

important differences were ≈6 for the OKS (≈14 for the OKS–PCS and

≈10 for the OKS–FCS). ~~and the , minimal important differences and the~~

~~precision of the change score were~~ also ~~calculated for the OKS, its subscales,~~

These ~~These~~ values were also calculated for the ICOAP and the KOOS-PS.

**Conclusions** The OKS summary scale, together with its pain and functional

component subscales, have excellent measurement properties when used

with patients with knee OA, undergoing non-operative treatment and is

superior to the ICOAP and the KOOS-PS for this purpose. This evidence

provides support for the validity of the use of the OKS when used across the

spectrum of knee OA disease severity, both in research and clinical practice.

## Article focus

- The OKS is a widely used patient reported outcome measure that was originally developed to measure the outcomes of knee replacement surgery.

- There is a growing interest to use the OKS in clinical practice, across the spectrum of OA disease.

- The aim of this study was to assess the measurement properties of the OKS when used with (individuals and groups of) patients who are undergoing non operative management for knee OA and compare it with most commonly used measures in this population of patients.

## Key Messages

- The OKS, and its pain and functional component subscales, have acceptable evidence of its measurement properties when used in patients (individual and groups) undergoing non-operative treatment for knee OA.

- The OKS performed better than the ICOAP and the KOOS-PS (widely used outcome measures for knee OA) on several counts.

## Strength and limitations of this study

- This study has conducted a comprehensive examination of scores' measurement properties.

- There might be a need to additionally re-evaluate evidence on some of the measurement properties presented here (such as interpretability or content validity), using different methods.

- The impact of the routine use of such scores in clinical practice should also be evaluated.

## INTRODUCTION

The Oxford Knee Score (OKS) is a widely used patient reported outcome measure (PROM), originally developed in1998 to be used in clinical trials for assessing the patient-perceived outcomes of knee replacement surgery. In this form it has proven to be reliable, valid and responsive.(1, 2) The remit of the OKS was extended in 2009, when it was adopted by the NHS PROMs programme in England and Wales as a primary outcome measure for knee replacement surgery.(3) Thus, OKS data are now collected on all patients undergoing knee replacement surgery preoperatively and at 6 months post operation, in order to monitor and benchmark the performance of health providers.

The increasing popularity of the OKS has also resulted in its being used for different populations and contexts from that for which it was originally developed. In particular there has been a growing interest in using the OKS in clinical practice as a means of standardizing clinical assessment, monitoring individual's self-reported health state across the spectrum of OA disease, and using the scores as an aid to clinical decision making. Extending the potential uses of PROMs in this manner has generally been highlighted as an opportunity to achieve maximum benefit from these measures, although the challenges of the application of such systems have also been recognised.(4, 5)

Using the OKS as a single score across the patient pathway, to aid diagnosis, monitor progression, assist in shared decision making and measure the outcome of intervention offers great potential for continuity of

care and understanding for patients. However robust evidence is required of the score's overall validity (i.e., the consistency of its measurement properties, such as reliability), when applied in these proposed new contexts. Generally, a measure is valid when applied to populations and contexts similar to the context in which the instrument was originally developed and tested, but measurement properties may change when the measure is applied in other contexts. The fact that the OKS was developed and tested to be used in the knee OA context (albeit end stage) is justification for considering its application in people with knee OA 'in general', but evidence has not been presented demonstrating that the OKS remains as reliable (both on an individual and a group level), valid and responsive when used with patients who are at earlier stages of their disease management.

The principal aim of our study was to assess the measurement properties of the OKS when used with (individuals and groups of) patients who are undergoing non operative management for knee OA, by examining its reliability, validity, responsiveness, and interpretability when applied in this context. Furthermore, we examined some of the measurement properties of the two most commonly used measures in this population; the Intermittent and Constant Osteoarthritis Pain (ICOAP)(6) and the Knee Injury and Osteoarthritis Score-Physical function Short form (KOOS-PS)(7).

1    **METHODS**

2        We obtained ethical approval for a prospective cohort study from a local

3    ethics committee (11/SC/005). Informed consent was obtained from all participants

4    in the study.

5

6    Study procedures and assessments

7        This study took place at an orthopaedic centre between June 2011 and

8    August 2012. Patients were eligible for inclusion if they were referred for knee

9    problems, had a confirmed diagnosis of knee OA and were enrolled in the non-

10   operative management pathway for their knee OA (as recommended by the

11   National Institute of Clinical Excellence (NICE)(8)). Treatments for patients were

12   tailored individually, taking into account patients' preferences and needs. As such,

13   they represented standard practice in the NHS. All patients who met these criteria

14   were sent an invitation letter containing information about the study, consent forms

15   and baseline questionnaires. Patients who consented to participate in the study

16   were asked to complete the OKS(2) the ~~Intermittent and Constant Osteoarthritis~~

17   ~~Pain (~~ICOAP(6)~~)~~, the ~~Knee Injury and Osteoarthritis Score-Physical function Short~~

18   ~~form (~~KOOS-PS~~)~~(7)~~,~~ and the SF12(9) patient-reported questionnaires.

19       The OKS is a 12 item questionnaire. It's item content was devised using

20   patient interviews, which addresses pain and functional impairment in relation to

21   their knee, in patients who are undergoing knee replacement surgery.(2) Likert

22   responses are recommended to be scored from 0 to 4, which are summed to

23   produce a summary score of 0 (worst) to 48 (best)(10). More recently, we

24   presented evidence (in the context of joint replacement) that supported the original

25   conceptual basis of the OKS using its composite summary scales, but which also

26  offered an option to perform additional analyses using pain and function

27  subscales.(11) The Pain Component Score (OKS-PCS) consists of items 2, 3, 7,

28  11 and 12 and the Functional Component Score (OKS-FCS) consists of items 1, 4,

29  5, 6, 8, 9 and 10. Subscale raw scores are standardized from 0 (worst) to 100

30  (best). Patients completed the OKS at baseline, 2 and 5 days (for test-retest

31  reliability) and at 3 months.

32          We asked the patients to complete the KOOS-PS and ICOAP at baseline

33  and 3 month follow up. These scores were developed to measure pain and

34  functional disabilities related to knee OA, and are now a recommended outcome

35  measures by the Osteoarthritis Research Society International (OARSI).

36          The KOOS-PS consists of 7 Likert-response items and was developed from

37  a longer version of the questionnaire (KOOS(12)) using Rasch analysis to measure

38  physical function in patients with various degrees of knee OA. It is scored as the

39  KOOS from 0 (best) to 4 (worst), with a summary raw score ranging from 0 to 28.

40  The score is converted to a true interval score that ranges from 0 (best) to 100

41  (worst). The ICOAP is an 11 item questionnaire whose items were informed from

42  focus groups with patents with hip or knee OA. It has two subscales that measure

43  the intermittent and constant pain with a standardized summary score ranging from

44  0 (best) to 100 (worst).

45          Patients also completed the generic SF-12, a 12-item general health

46  measure with 8 items that have Likert-type response categories and 4 items with

47  dichotomous (yes/no) response categories. The SF-12 is scored as a Physical

48  Component Summary (PCS) and Mental Component Summary (MCS) ranging

49  from 0 (worst) to 100 (best).

50     Lastly, we asked the patients to complete a transition question in regards to

51 the change they experienced from the baseline measurement: "Compared to one

52 week before your clinic visit, please indicate how much your knee problem has

53 changed?" The question had three response options: "1. My knee has got better; 2.

54 My knee has stayed the same; 3. My knee has got worse".

55     We supplemented patient reported outcome data with information on their

56 body mass index (BMI) and the degree of structural changes observed in the knee,

57 which was available from the patients' medical records. An orthopaedic surgeon

58 (LDJ) performed Kellgren-Lawrence (K-L) grading using available knee OA

59 radiographs. (13) The degree of structural changes in the knee was classified

60 using (K-L) grading. In the absence of X-rays, we assessed intra-operative

61 documentation from previous knee arthroscopy or available MRIs to examine the

62 extent of cartilage loss and confirm the diagnosis of osteoarthritis.

63

64 <u>Statistical methods</u>

65     The recommended minimum sample sizes for validation studies (based on

66 optimal numbers for correlations) often range from 50 to 100.(14, 15) For

67 confirmatory factor analysis (CFA) the literature agrees with a minimum sample

68 size of about 100-150 or about 10 subjects per questionnaire item.(16, 17) These

69 sample sizes are required for data analyses and should be adjusted (i.e.

70 increased) for the risk of loss to follow up. In this study we stopped recruiting when

71 the dataset enabled us to perform CFA with at least 10 subjects per item.

72     We analysed the data using SPSS version 20 and LISREL V 8.80. Baseline

73 and 3 month follow up scores were generally non-normally distributed and change

74 scores approximated to normal (except the ICOAP and the OKS-PCS). We used

75  non-parametric statistics, where appropriate. We did not use data imputation and

76  we excluded cases with missing data on analysis by analysis basis (unless

77  mentioned otherwise). We examined the following measurement properties of the

78  OKS:

79

80  **Reliability**

81  Reliability is an estimation of the consistency and stability of a measure. It

82  includes analysis of the extent to which a measure is internally consistent

83  (measured by the inter-correlation of all items) and free from measurement error.

84  We used Cronbach's alpha to assess the internal consistency of the OKS

85  summary scale and its subscales. Alpha values of at least 0.7 are recommended in

86  order to demonstrate internal consistency. (18) We calculated an intraclass

87  correlation coefficient ($ICC_{2,1}$)(19) to assess the test-retest reliability of the OKS

88  and its subscales. Minimum ICC values of 0.7 are normally considered acceptable

89  (18) although higher values are required for the use of the score applied at an

90  individual level. To inform the potential use of the OKS on the individual level, we

91  calculated the precision of individual scores at 90% CI level by multiplying the

92  standard error of measurement (SEM) by the 2-tailed z value at 90%.

93

94  **Construct validity**

95  The validity of a measure is concerned with whether a measure actually

96  measures what it purports to measure.(20, 21) The definition of validity has

97  recently been further refined as: "The degree to which accumulated evidence and

98  theory support specific interpretations of test scores entailed by proposed uses of a

99  test".(22) Construct validity of a measure is supported by the accumulation of

100 evidence obtained by testing hypotheses about the relationship that the measure

101 exhibits with other (validated) measures.(20)

102      We examined the construct validity of the OKS summary scale and its

103 subscales by testing an *a priori* set *of* hypotheses about the expected relationships

104 between the instruments at baseline:

105      (i) the OKS and the physical component summary of the SF12 (PCS-12) are

106 measuring sufficiently similar constructs (SF-PCS measures self-reported physical

107 function and the OKS measures self-reported pain and physical functioning related

108 to the knee), so the correlation between these two instruments' scales should be

109 moderate and in the same direction,

110      (ii) the correlation between the OKS and the mental component summary of

111 the SF12 (MCS-12) should be weaker than the one between the PCS-12 and OKS

112 as these two scale constructs are not considered to be related to such an extent,

113      (iii) the OKS and KOOS-PS are measuring a sufficiently similar construct

114 (the KOOS-PS measures self-reported knee function and the OKS measures self-

115 reported pain and physical functioning related to the knee) that the correlation

116 between these two measures should be strong and negative (as scores go in the

117 opposite direction),

118      (iv) the OKS and the ICOAP are measuring sufficiently similar constructs

119 (the ICOAP measures self-reported knee pain and the OKS measures self-

120 reported pain and physical functioning related to the knee) that the correlation

121 between these two measures should be strong and negative,

122      (v) the OKS-PCS should be correlated more with the ICOAP than with the

123 KOOS-PS and negatively, in each case (the OKS-PCS measures self-reported

124 knee pain as does the ICOAP),

125    (vi) the OKS-FCS should be correlated more with the KOOS-PS that the

126    ICOAP and negatively (the OKS-FCS measures self-reported knee function, as

127    does the KOOS-PS).

128    We classified correlations (r) as: r=0 to 0.29 as none/weak; r= 0.3 to 0.69 as

129    moderate; and r > 0.7 as strong.

130    **Structural validity** is one particular aspect of construct validity; it examines

131    the extent to which the dimensionality of a measure corresponds to the construct

132    (i.e. latent variable) that is supposed to be measured.(20) For instance, if a

133    measure is unidimensional (i.e. it is supposed to measure one construct, such as

134    pain) all of its items will measure the same underlying construct. We examined the

135    structural validity of the OKS by conducting Confirmatory Factor Analysis (CFA)

136    that tested the fit of the one and two factor models of the OKS to the data, using

137    LISREL V8.80 software. In line with the standard CFA testing guidelines, we

138    considered the following indices as satisfactory: a non-significant $\chi^2$ (p>0.05),

139    standardised root mean square residual (SRMR)>0.08, comparative fit index (CFI)

140    >0.95, root mean square error of approximation  (RMSEA): <0.05 close fit,

141    <0.08good fit, <0.1 satisfactory fit; RMSEA p test of close fit>0.05.(23) Additionally,

142    we used the Chi-square ($\chi$2) difference test and Parsimonious Normed Fit Index

143    (PNFI) to compare the fit between the two models of the OKS and the ICOAP. (24)

144    We calculated the $\chi$2 difference tests by looking at the difference of $\chi$2 of two

145    models along with the difference in their degrees of freedom.  ~~We checked the $\chi$2~~

146    ~~difference, with its degrees of freedom in the  $\chi$2 distribution table. If this value is~~

147    ~~statistically significant, then the model with more degrees of freedom is favoured.~~

148

149    **Responsiveness**

150    The ability of a measure to detect meaningful clinical change (where it has

151    occurred) over time is critical for the use and the application of a measure. (25)

152    This change might occur following an intervention, or just occur 'naturally' during a

153    period of observation. Generally, as with construct validity, responsiveness is

154    assessed by testing *a priori* hypotheses about the relationship of the changes in

155    one measure to the changes in another (validated) measure, or with reference to a

156    change in a gold standard (as with testing criterion validity). Responsiveness can

157    also be tested with reference to a transition item, where the responsiveness is

158    tested only in subjects who have reported that clinical change has occurred.

159    We used a one sample t-test (2 tailed) to assess if the changes at 3 months

160    for the OKS, its subscales (OKS-PCS and OKS-FCS), KOOS-PS and the ICOAP

161    were significantly different from 0. We constructed a Cumulative Distribution

162    Function (CDF) plot for the; (i) OKS, (ii) OKS-PCS and ICOAP, and (iii) OKS-FCS

163    and KOOS-PS to examine the proportion of individual patients who experienced

164    deterioration and improvement beyond the measurement error of the instrument at

165    the individual level and to compare the proportion of change in pain and function

166    detected by the different measures.

167    As with construct validity, we tested the responsiveness by setting *a priori*

168    hypotheses about the direction and magnitude of changes of the validated

169    comparator instruments and the OKS:

170    (i) the change scores in the OKS should correlate strongly with the change

171    scores in the KOOS-PS and ICOAP,

172    (ii) the change scores in the OKS-PCS should correlate more strongly with

173    the change scores in the ICOAP than with the change scores in the KOOS-PS,

174    (iii) change scores for the OKS-FCS should correlate more strongly with

175 change scores for the KOOS-PS than the change scores for the ICOAP.

176    All correlations should be negative.

177

178    There was a concern about the amount of overall change that can be

179 experienced as a result of such a management pathway (which included a wide

180 range of individually tailored treatments administered to a heterogeneous sample),

181 so we additionally defined the construct of change using a patient rated item of

182 change. We then used the responses to this item to calculate anchor based values

183 of minimal important change and difference.

184

185 **Interpretability**

186    Interpretability is defined as the degree to which one can assign qualitative

187 meaning to a quantitative score.(20) In clinical trials, this issue can concern the

188 question of what is considered to be a 'good', 'bad' or 'indifferent' outcome (as

189 measured by a particular criterion or score) and what is considered to be a

190 clinically relevant change. The minimum amount of change that is discerned as

191 meaningful by patients is particularly important as it affects interpretation of study

192 results.

193    We assessed the interpretability by relating the change in the PROMs

194 scores to the patient reported item of change (using an anchor based method) and

195 by relating the observed change in the score to its measurement error at the

196 individual level (using a distribution based method). Average change in the score

197 associated with the group of patients who responded with "My knee has got better"

198 on the transition item was taken as the anchor based minimal important change

199  (MIC). The difference in the change score between the groups of patients who

200  responded with "My knee has stayed the same" and "My knee has got better on

201  the global item of change was taken as the minimal important difference (MID).

202  Finally, the minimum change in the instrument that represents real change (beyond

203  measurement error) was calculated using the Minimum Detectable Change

204  (MDC$_{90}$), which was obtained by multiplying the SEM with the z-value at the 90%

205  level and the square root of two (to account for two measurement occasions).((26,

206  27))

207

208

209    **RESULTS**

210

211         **Sample characteristics.** 137 patients were recruited in the study. 21

212    patients did not complete follow up questionnaires at 3 months, out of which 3

213    patients were listed for a surgical procedure (2 osteotomies and 1 arthroplasty)

214    before 3 month follow-up, 7 patients no longer wanted to participate in the study

215    and 11 were lost to follow-up.  134 patients were included in the main baseline

216    analysis of whom 67 (50 %) were male and 67 patients were female. The mean

217    age of patients was 59 (SD 11)~~, which is about 10 years less than the average age~~

218    ~~of the developmental sample of the OKS~~. 70% of patients had information on Body

219    Mass Index (BMI), out of whom 30% were classified as obese (BMI>30), 41% as

220    overweight (BMI between 25 and 29.9), 29% as normal weight (BMI between 18.5

221    and 24.9). No one was classified as underweight. All of the patients had a

222    diagnosis of knee osteoarthritis. 2% of the patients had Kellgren-Lawrence (KL)

223    grading of 0 (but evidence of cartilage loss on MRI scan), 8% had K-L of 1, 43%

224    had K-L of 2, 16% had K-L of 3, 4% had K-L of 4.  For 26% of cases, X-ray

225    information was unavailable, of whom, 20% had their diagnosis confirmed on the

226    basis of MRI, while 6% of patients did not have X-rays or MRIs accessible

227    (however, these patients had the diagnosis of OA previously confirmed in the

228    primary care setting, different trust, or in a private clinic). All patients underwent

229    standard non-operative management of knee OA.(8)

230         116 (87%) out of 134 recruited patients returned the questionnaires at three

231    month follow up.  There was no difference in age or BMI between those patients

232    who did not respond at three months versus those who did, but baseline OKS was

233    different between these groups. The group that did not respond had scored, on

234  average, 7.3 points lower (worse) on the OKS than responders at three months

235  (Independent samples t-test, p<0.05). A summary of the baseline scores is

236  presented in Table 1.

237

238  Table 1. Baseline scores for the OKS, its subscales (OKS-PCS and OKS-FCS),
239  ICOAP, KOOS-PS, and SF-12 physical and mental summaries (PCS-12 and MCS-
240  12).
241
242

| | N | | Mean (SD) | Median | Percentiles | |
|---|---|---|---|---|---|---|
| | Valid | Missing | | | 25 | 75 |
| OKS | 121 | 13 | 29.3 (10) | 30 | 22 | 37 |
| OKS-PCS | 123 | 11 | 57.4 (23) | 57 | 43 | 75 |
| OKS-FCS | 137 | 7 | 66.5 (22) | 70 | 50 | 85 |
| ICOAP | 124 | 10 | 37.8 (25) | 31.8 | 16 | 57 |
| KOOS-PS | 112 | 22 | 40.5 (18) | 38.6 | 32 | 49 |
| PCS-12 | 130 | 4 | 36.7 (10) | 35 | 29 | 45 |
| MCS-12 | 130 | 4 | 51 (12) | 56 | 43 | 60 |

243
244

245

246  For comparison, in the developmental study of the OKS, the median age of

247  patients undergoing knee replacement was 73 and in this study the median age

248  was 58 (mean 59). (2) There was also considerable difference in self-reported pain

249  and functional disability between the patients in the two studies. The mean

250  baseline OKS in this sample was 29, compared to the mean preoperative OKS of

251  in the developmental study sample of 17 (when transformed to the 0-48 scoring

252  system).

253

254  **Reliability**

255

256  Cronbach's alpha for the 12-item OKS was 0.94, 0.88 for the OKS-FCS and

257  0.90 for the OKS-PCS. For the ICOAP and KOOS-PS, the Cronbach's alpha was

258  0.97 and 0.94 respectively. The alpha value did not change considerably if any of

259  the items were sequentially removed from the total scores.

260  Test retest reliability ICCs were 0.93 (95% CI, 0.91-0.95) for the summary

261  OKS, 0.91 (95% CI, 0.88-0.94) for the OKS-PCS and 0.92 (95% CI, 0.90-0.95) for

262  the OKS-FCS.

263  The standard error of measurement (SEM) for the summary OKS was 2.65

264  and the confidence in individual single score at 90% was ±4.4 OKS points. SEM for

265  the OKS-FCS was 6.2 with ±10.2 90% CI for individual score and the SEM for the

266  OKS-PCS was 6.9 with ±11.3 points as 90% CI for individual score (noting that the

267  OKS-PCS and the OKS-FCS are presented on a different scale than the OKS).

268  The SEM for the ICOAP was 9.68 with ±15.9 points as 90% CI for individual score.

269  We calculated the SEM for the ICOAP by using the test-retest reliability that was

270  reported in the developmental study (0.85)(6). For the KOOS-PS, this information

271  for the English version of the questionnaire was not available, so we used the test-

272  retest reliability value of 0.86 from the validation of the French version of the

273  questionnaire. (28) The SEM for the KOOS-PS was 6.7 with ±11.1 points as 90%

274  CI for individual score.

275

276  **Construct validity**

277  **Construct validity (hypothesis-testing).** All correlations were generally

278  consistent with *a priori* hypotheses concerning the relationships of the OKS with

279  comparator instruments. Spearman's ρ between the baseline OKS, KOOS-PS,

280  ICOAP, SF12-MCS and SF-12-PCS are shown in Table 2. The OKS correlated

281  strongly with the KOOS-PS and ICOAP. The correlation between the SF12-PCS

282  and the OKS was slightly higher than expected.  As expected, the OKS was most

283 poorly related to the SF12-MCS. The OKS-PCS correlated more with ICOAP than

284 with KOOS-PS and the OKS-FCS correlated more with the KOOS-PS that with

285 ICOAP. This evidence supports convergent and divergent validity of the OKS.

286

287 Table 2: Baseline Spearman's correlations between the scores. All correlations
288 were significant at the 0.01 level (2-tailed). The number of cases with complete
289 information that allowed the calculation of the correlation coefficients is in brackets
290 for each correlation.
291

|  | OKS | OKS-PCS | OKS-FCS |
|---|---|---|---|
| ICOAP | -.879 (115) | -.884 (117) | -.792 (121) |
| KOOS-PS | -.849 (106) | -.779 (107) | -.867 (111) |
| PCS-12 | .648 (121) | / | / |
| MCS-12 | .370 (121) | / | / |

292
293

294 **Structural validity.** 122 pre-operative OKSs, 125 pre-operative ICOAP and

295 113 pre-operative KOOS-PS were available for the CFA. Fit indices of one and two

296 factor models for the OKS are presented in Table 3. Neither of the one and two

297 factor models was rejected. Fit indices favoured the 2 factor model and the

298 reduction in $\chi^2$ in the two factor model was significant ( $\chi 2 diff > 7.879$, with df=1, at

299 the a=0.005 level).

300

301 Table 3. Fit indices of one and two-factor model of the OKS.

| Factors | χ2 (p value) | df | RMSEA | 90% CI RMSEA | RMSEA p test | CFI | SRMR | PNFI |
|---|---|---|---|---|---|---|---|---|
| 1 | 71.32 (p=0.06) | 54 | 0.052 | 0.00-0.08 | 0.44 | 0.99 | 0.043 | 0.80 |
| 2 | 56.64 (p=0.34) | 53 | 0.024 | 0.0-0.06 | 0.83 | 1 | 0.039 | 0.79 |

302 Note. F=number of factors; 2 =chi-square; df=degrees of freedom; RMSEA=root mean square of
303 approximation; CI=confidence intervals; p-value for test of close fit (RMSEA<.05);
304 SRMR=standardized root mean square residual; CFI-comparative t index; PNFI=parsimonious
305 normed fit index.
306

307 CFA revealed that a one-factor KOOS-PS model was rejected by the $\chi 2$ test and

308 its RMSEA was above the highest acceptable threshold of an acceptable fit (0.1)

309 (Table 4). The SRMR was acceptable and CFI was on the threshold of a good fit.

310  Both one and two factor ICOAP models were rejected by the χ2 test and both

311  models had RMSEA values far above the lowest threshold of an acceptable fit.

312  However, SRMR and CFI were acceptable for both scores. There was no

313  significant reduction (at the 0.05 level) in χ2 for the 2 factor model of the ICOAP

314  (χ2diff< 3.84, with df=1).

315

316

317

318

319

320  Table 4. Fit indices of one and two-factor model of the ICOAP and KOOS-PS.

| | $\chi^2$ (p value) | df | RMSEA | 90% CI RMSEA | RMSEA p test | CFI | SRMR | PNFI |
|---|---|---|---|---|---|---|---|---|
| ICOAP (1F) | 242.31 (p=0.00) | 44 | 0.19 | 0.17-0.22 | 0.00 | 0.95 | 0.064 | 0.75 |
| ICOAP (2F) | 228.19 (p=0.00) | 43 | 0.19 | 0.16-0.21 | 0.00 | 0.96 | 0.057 | 0.74 |
| KOOS-PS (1F) | 40.88 (p=0.00) | 14 | 0.13 | 0.09-0.18 | 0.00 | 0.98 | 0.046 | / |

321  Note. F=number of factors; 2 =chi-square; df=degrees of freedom; RMSEA=root mean square of
322  approximation; CI=confidence intervals; p-value for test of close fit (RMSEA<.05);
323  SRMR=standardized root mean square residual; CFI-comparative t index; PNFI=parsimonious
324  normed fit index.
325

326  **Responsiveness**

327  Figure 1 shows the CDF plot for the OKS. The plot demonstrates that,

328  based on the OKS summary score, approximately 15% of patients in the study

329  experienced deterioration in health state, at three month follow up, that was

330  beyond the MDC$_{90}$ of 6 points, approximately 30% of patients experienced

331  improvement and 55% of patients did not experience change beyond this value.

332  Also, slightly less than 30% of the patients experienced improvement that was

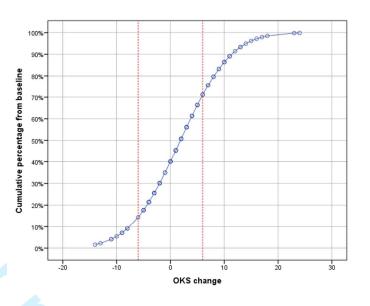333  beyond the MIC of 7 points on the OKS.

334



335 Figure 1. Cumulative percentage of patients experiencing the change on the OKS
336 from baseline less or equal to the value on the x-axis. Red line marks the minimum
337 detectable change beyond the measurement error of the score (MDC$_{90}$ of 6
338 points).
339
340      Table 5 shows the mean baseline, three month follow-up change scores,

341 and p values for the significance of 3 month change and ES for the OKS, OKS-

342 PCS, OKS-FCS, KOOS-PS and ICAOP for the overall cohort. All mean changes

343 were significant at the 0.01 level (2-tailed t-test) except the OKS-FCS.

344

345 Table 5: Significance of change in OKS, its subscales (OKS-PCS and OKS-FCS),
346 ICOAP and KOOS-PS scores at three months (one sample t-test).
347

|  | N | Baseline (SD) | 3 months (SD) | Change (SD) | p-value | ES |
|---|---|---|---|---|---|---|
| OKS | 104 | 30.29 (10) | 32.15 (11) | 1.87 (7) | 0.01 | 0.19 |
| OKS-PCS | 107 | 59.36 (22) | 65.13 (24) | 5.77 (17) | <0.01 | 0.26 |
| OKS-FCS | 108 | 67.22 (21) | 68.66 (23) | 1.44 (16) | 0.4 | 0.07 |
| ICOAP[a] | 104 | 37.19 (25) | 31.53 (25) | -5.66 (19) | <0.01 | 0.23 |
| KOOS-PS[a] | 92 | 39.42 (18) | 34.88 (20) | -4.5 (14) | <0.01 | 0.25 |

348 Note. N=number of complete cases available for calculation of 3 month follow up; SD=standard
349 deviation; ES=effect size;[a] The ICOAP and the KOOS-PS represent severity of the disease in the
350 opposite direction from the OKS and its subscales.
351

352

353    The correlations between the changes in the OKS and changes in the

354    KOOS-PS and the ICOAP were somewhat less than anticipated (0.67 and 0.62

355    respectively). As hypothesized, the changes in the OKS-PCS correlated more with

356    the changes in ICOAP (also assessing knee pain) than KOOS-PS, and the

357    changes in the OKS-FCS correlated more strongly with the changes in the KOOS-

358    PS (also assessing knee function) than with the changes in the ICOAP (Table 6).

359

360    Table 6: Spearman's correlations between the 3 month changes in the OKS and its
361    subscales (OKS-PCS and OKS-FCS), ICOAP and KOOS-PS.
362

|  | ICOAP | KOOS-PS |
|---|---|---|
| OKS | -.674 (96) | -.617 (87) |
| OKS-PCS | -.669 (99) | -.551 (88) |
| OKS-FCS | -.598 (100) | -.622 (90) |

363    Note. All correlations are significant at the 0.01 level (2-tailed). The number of cases with complete
364    information that allowed the calculation of the correlation coefficients is in brackets for each
365    correlation.
366
367
368    **Interpretability**

369    Tables 7 and 8 present the percentage of responses for different response

370    categories, effect sizes and mean score changes by response category. We

371    conducted independent sample t-tests for the equality of means between the mean

372    scores for groups of patients who responded 'better' and 'the same' on the

373    transition item. Only the OKS, OKS-PCS, and OKS-FCS had registered significant

374    differences between the means (2 tailed, $p<0.05$) of groups who responded that

375    they were better/the same. ~~Here, the OKS and OKS-PCS mean differences were~~

376    ~~close to (and generally just above) scale MDC/MID values and thus likely beyond~~

377    ~~measurement error, while the OKS-FCS mean differences were just less than the~~

378    ~~subscale's MDC/MID values. All OKS scales' mean differences were greater than~~

379    ~~the scales' relevant SEM values.~~

380    Table 9 presents the summary of interpretability indices.

381

Table 7: Number (N) and percentage of responses for different response categories with effect sizes (ES), mean score changes by response category and ANOVA tests for linear trend for the mean score across the three response categories for the OKS and its subscales (OKS-PCS and OKS-FCS).

386
387

|  |  | Better | Same | Worse |
|---|---|---|---|---|
| OKS | N (% of responses) | 30 (33) | 26 (28) | 36 (39) |
|  | Mean change (SD) | 7.1 (8) | 0.7 (6) | -1.88 (5) |
|  | ES | .7 | .1 | -.2 |
|  | P-value for linear trend | <.001 | <.001 | <.001 |
| OKS-PCS | N (% of responses) | 31 (33) | 28 (30) | 38 (35) |
|  | Mean change (SD) | 17.27 (19) | 2.93 (14) | -2.68 (11) |
|  | ES | .8 | .2 | -.1 |
|  | P-value for linear trend | <.001 | <.001 | <.001 |
| OKS-FCS | N (% of responses) | 28 (33) | 26 (31) | 30 (36) |
|  | Mean change (SD) | 10.63 (14) | 1.11 (16) | -6.35 (14) |
|  | ES | .5 | .1 | -.3 |
|  | P-value for linear trend | <.001 | <.001 | <.001 |

388

389

Table 8: Number (N) and percentage of responses for different response categories with effect sizes (ES), mean score changes by response category and ANOVA tests for linear trend for the mean score across the three response categories for the ICOAP and the KOOS-PS.

394

|  |  | Better | Same | Worse |
|---|---|---|---|---|
| ICOAP | N (% of responses) | 32 (34) | 27 (29) | 35 (37) |
|  | Mean change (SD) | -13.42 (23) | -5.64 (17) | 2.73 (16) |
|  | ES | -.6 | -.3 | .1 |
|  | P-value for linear trend | <.003 | <.003 | <.003 |
| KOOS-PS | N (% of responses) | 25 (31) | 27 (33) | 30 (37) |
|  | Mean change (SD) | -11.98 (15) | -4.22 (12) | 1.61 (12) |
|  | ES | -.8 | -.3 | .1 |
|  | P-value for linear trend | <.001 | <.001 | <.001 |

395

Table 9: Anchor based and distribution based MIC/MID values for the OKS, its subscales, ICOAP and KOOS-PS.

398

|  | Distribution based | Anchor based | |
|---|---|---|---|
|  | $MDC_{90}$ | MID | MIC |
| OKS | ±6 | 6.4 | 7.1 |
| OKS-PCS | ±16 | 14.3 | 17.3 |
| OKS-FCS | ±15 | 9.5 | 10.6 |
| ICOAP | ±23 | 7.8 | 13.4 |
| KOOS-PS | ±16 | 7.8 | 12.0 |

Note. $MDC_{90}$=minimum detectable change; MID=minimum important difference; MIC=minimum important change.

401

## DISCUSSION

The OKS summary scale and its pain and functional component subscales were each found to have acceptable evidence of their measurement properties to support their use with groups of patients (research/audit) and for individuals (clinical practice) who are undergoing non-operative treatment for knee OA. The OKS summary scale and its subscales were validated against the KOOS-PS, the ICOAP (measures developed for use in patients with knee OA) and the SF-12 by testing logical *a priori* hypotheses regarding the construct validity and responsiveness of the OKS and its subscales in comparison to these other (validated) measures. Thus, CFA demonstrated excellent fit and confirmed the structural validity of the OKS and both subscales. Furthermore, assessment of test-retest reliability demonstrated that the OKS and its subscales could all be used both at group and individual levels (clinical practice)(29).

The OKS subscales can be used to specifically target the improvement or deterioration in pain or function, whether in research (as an endpoint or for sample size calculations) or in clinical practice. Anchor based MIC of ≈7 for the OKS, ≈17 for the OKS-PCS, and ≈11 for the OKS-FCS can be used in cohort studies to assess if the change in the OKS (from baseline) is clinically relevant. Anchor based MID of ≈6 for the OKS, ≈14 for the OKS-PCS, and ≈10 for the OKS-FCS can be used in clinical trials to assess if the difference in change between two arms of treatment is clinically relevant. Finally, changes in individual patient scores beyond the $MDC_{90}$ (≈6 points for the OKS, ≈16 points for the OKS-PCS, and ≈15 points for the OKS-FCS) can be used as a benchmark of improvement or deterioration that is beyond the

measurement error of the score. These values are likely to be different if the OKS is used in a different population of patients (i.e. patients undergoing knee replacement surgery).

## Limitations

Even though the reliability, construct validity and responsiveness of the OKS and its subscales have been proven to be satisfactory when used in patients undergoing non-operative management for their knee OA, there might be a need to further verify its content validity in this extended context.(30) The items for the OKS were originally devised using a representative sample of patients with end stage disease, who were undergoing knee replacement surgery. It could be argued that the measure in its current form might not fully represent the concerns of this slightly different population of patients whose knee OA is generally at an earlier stage. If a measure is used in a different context or with different type of patients than that which was used  in its design/development, then the content validity may be suspect (in relation to the new/different usage).(18) ~~On the other hand it may be assumed to be appropriate if the context is considered to be 'similar enough'.(20)~~ A counter~~nother~~ argument is that it is unrealistic to have a new/different measure (and a new study conducted to design and test one) for every possible sub category of patient or type of treatment within all diseases or conditions. In such cases a researcher should make a judgement about the best available/closest measure (21), but as a minimum should check that the measurement properties are still otherwise maintained. Any further

examination of the content validity of the OKS in this extended context would necessitate a new study (based on qualitative interviews) being undertaken.

One of the limitations concerns the use of the transition question with three response levels (better, the same, worse). MIC/MID values depend on the number of response categories on the transition question. If, for instance, a response category 'a little better' was used instead of 'better' the final MIC value would have probably been smaller. Indeed, the methods of MIC/MID estimation have been a subject of debate within the scientific community and we would recommend that any application of the MIC/MID values presented in this paper is done with awareness of its caveats. However, regardless of the shortcomings of the transition item, the same was used in the comparative analysis of interpretability between the OKS, its subscales, the KOOS-PS and the ICOAP and in terms of drawing conclusions about the comparative performance between the scores, this is not such a source of concern.

## Comparative performance of the OKS and its subscales versus the ICOAP and the KOOS-PS in this study

Even though the ICOAP and the KOOS-PS are currently widely used as outcome measures for knee OA, the OKS performed better in this study on several counts.

The 11-item ICOAP had a Cronbach's alpha of 0.97 (compared to the alpha of the OKS-PCS of 0.9) and the alpha was 0.94 for the KOOS-PS (compared to the alpha of 0.87 for the OKS-FCS). A high alpha value can mean that some of the items on a scale are redundant and this seems to be

more of a concern for the ICOAP and KOOS-PS than for the OKS subscales. Furthermore, the reliability and precision of the score was better for the OKS and its subscales than for the KOOS-PS and ICOAP, which makes it more suitable to be used in clinical practice.

There was evidence to support both one and two factor models of the OKS, but no acceptable evidence of structural validity was found for the KOOS-PS or the ICOAP. The KOOS-PS and the one and two-factor ICOAP models were rejected by the $\chi 2$ test. Furthermore, RMSEAs were unacceptably high for both scales. The exploration of the sources of poor fit of these measures is beyond the scope of this study and future studies should investigate this problem further (perhaps also using exploratory factor analysis).

We have some concerns about the interpretability of the ICOAP and KOOS-PS. It seems that these measures performed less well than the OKS in this regard. First, due to the fact that the ICOAP has low precision at the individual level (the $MDC_{90}$ is almost 10 points larger than the MIC) this makes it less suitable to interpret change scores in individual patients. Second, although around one third of the patients in our sample reported being better following 3 months of non-operative management for knee OA, neither the ICOAP or the KOOS-PS obtained statistically significant differences in the change score between the groups of patients who reported themselves to be better or the same (in contrast with the OKS and its subscales). This could indicate problems with the sensitivity of these scores to change. Third, whilst there was some lack of symmetry between the mean change in the OKS score and its subscales in relation to the patient rated item of change (patients

who claim they had not experienced change on the global transition item, actually experienced change as measured by the PROM), this lack of symmetry seems to be more pronounced for the KOOS-PS and ICOAP.

**Implications for clinicians and policymakers**

In this study, we obtained evidence that supports the use of the OKS and its pain and functional subscales in patients who are undergoing non-operative management for their knee. When used with patients in this context, the OKS has demonstrated evidence of validity, reliability, and responsiveness in measuring the health state of individuals. The measure could be used in clinical practice to monitor disease progression in individual patients undergoing non-operative management for their knee OA, or for hospital audit where the information from groups of patients is analysed to assess the effectiveness of current patient management pathways for treating OA in terms of health gain/deterioration.

Although this study was conducted on a sample of patients with knee OA presenting themselves in the secondary care setting, we consider that the findings presented here may be generalizable to the primary care setting. Studies have shown no significant differences in the pain severity and function between the groups of patients with knee OA who get referred to secondary care and who do not. (31, 32) Other factors, such as the chronicity of the disease, or complex interaction of psychological and social factors, are more associated with secondary care referral. However, further research, involving larger sample sizes, is needed to confirm these findings.

The use of a single valid score across a patient pathway is a compelling goal when considering how to develop standardisation of patient care in the NHS. Our new evidence suggests extending the use of the OKS in the patient pathway for managing knee OA may be possible. However the practicalities and feasibility of widespread score administration need further exploration focusing on appropriate timing, frequency and method of score administration. (33) Most importantly, more work is required to understand how results of the OKS, if adopted earlier in the pathway, should be interpreted to support patients in shared decision making regarding treatment options and the influence that such routine use of the OKS might have on the quality of care that patients receive (i.e. the effect on the quality of service and influence on patients' clinical outcomes). (34)

For peer review only

## ACKNOWLEDGMENTS

A copy of the OHS and OKS questionnaires and permission to use this

measure can be acquired from Isis Innovation Ltd, the technology transfer

company of the University of Oxford via website:

http://www.isis-innovation.com/outcomes/index.html or email:

healthoutcomes@isis.ox.ac.uk

## COMPETING INTEREST DECLARATION

All authors have completed the Unified Competing Interest form at

www.icmje.org/coi_disclosure.pdf (available on request from the

corresponding author) and declare that KKH, LDJ, AJP, DJB have no financial

interests that may be relevant to the submitted work. JD is one of the original

inventors of the OHS and OKS. She has received consultancy payments, via

Isis Innovation, in relation to work involving both questionnaires.

## CONTRIBUTIONS

Conception and design: JD, DJB, AJP

Acquisition of data: KKH, LDJ

Analysis and interpretation of data: KKH, JD, DJB, AJP

Drafting of the article and revision it critically for important intellectual content:

KKH, JD, LDJ, DJB, AJP

Final approval of the article: KKH, JD, LDJ, DJB, AJP

All authors, external and internal, had full access to all of the data (including

statistical reports and tables) in the study and can take responsibility for the

integrity of the data and the accuracy of the data analysis.

## ETHICS APPROVAL

This study obtained ethics approval from the Oxfordshire Research Ethics

Committee B (11/SC/005). Informed consent was obtained from all

participants in the study.

## ROLE OF THE FUNDING SOURCE

**DATA SHARING STATEMENT**

Anonymised data and statistical codes are available from the corresponding

author.

# REFERENCES

1.      Dawson J, Fitzpatrick M, Churchman D, Verjee-Lorenz A, Claysonm D. User Manual for the Oxford Knee Score (OKS). 2010.

2.      Dawson J, Fitzpatrick R, Murray D, Carr A. Questionnaire on the perceptions of patients about total knee replacement. Journal of Bone and Joint Surgery British Volume. 1998 Jan;80(1):63-9. PubMed PMID: 9460955. Epub 1998/02/14. eng.

3.      Department of Health. Guidance of the Routine Collection of Patient Reported Outcome Measures (PROMs). In: Department of Health, editor. London2008.

4.      Devlin NJ, Appleby J. Getting the most out of PROMS. The Kings Fund Office of health economics. 2010.

5.      Valderas J, Alonso J. Patient reported outcome measures: a model-based classification system for research and clinical practice. Quality of Life Research. 2008;17(9):1125-35.

6.      Hawker GA, Davis AM, French MR, Cibere J, Jordan JM, March L, et al. Development and preliminary psychometric testing of a new OA pain measure--an OARSI/OMERACT initiative. Osteoarthritis and Cartilage. 2008 Apr;16(4):409-14. PubMed PMID: 18381179. Epub 2008/04/03. eng.

7.      Perruccio AV, Stefan Lohmander L, Canizares M, Tennant A, Hawker GA, Conaghan PG, et al. The development of a short measure of physical function for knee OA KOOS-Physical Function Shortform (KOOS-PS) - an OARSI/OMERACT initiative. Osteoarthritis and Cartilage. 2008 May;16(5):542-50. PubMed PMID: 18294869. Epub 2008/02/26. eng.

8.      National Institute for Health and Clinical Excellence (NICE). Osteoarthritis. National clinical guideline for care and management in adults. London: Royal College of Physicans; 2008.

9.      Ware J, Jr., Kosinski M, Keller SD. A 12-Item Short-Form Health Survey: construction of scales and preliminary tests of reliability and validity. Medical Care. 1996 Mar;34(3):220-33. PubMed PMID: 8628042. Epub 1996/03/01. eng.

10.      Murray DW, Fitzpatrick R, Rogers K, Pandit H, Beard DJ, Carr AJ, et al. The use of the Oxford hip and knee scores. Journal of Bone and Joint Surgery British Volume. 2007 August 1, 2007;89-B(8):1010-4.

11.      Harris K, Dawson J, Doll H, Field R, Murray D, Fitzpatrick R, et al. Can pain and function be distinguished in the Oxford Knee Score in a meaningful way? An exploratory and confirmatory factor analysis. Quality of Life Research. 2013 2013/03/23:1-8. English.

12.      Roos EM, Roos HP, Lohmander LS, Ekdahl C, Beynnon BD. Knee Injury and Osteoarthritis Outcome Score (KOOS)--development of a self-administered outcome measure. Journal of Orthopaedic and Sports Physical Therapy. 1998 Aug;28(2):88-96. PubMed PMID: 9699158. Epub 1998/08/12. eng.

13.      Kellgren J, Lawrence J. Radiological assessment of osteo-arthrosis. Annals of the Rheumatic Diseases. 1957;16(4):494-502.

14.      De Vet HCW, Terwee CB, Mokkink LB, Knol DL. Measurement in Medicine a Practical Guide Cambridge: Cambridge University Press; 2011. Available from: http://public.eblib.com/EBLPublic/PublicView.do?ptiID=802925.

15.      Kline P. An easy guide to factor analysis. 1993.

16.      Bentler PM, Chou CP. Practical issues in structural modeling. Sociological Methods & Research. 1987;16(1):78-117.

17.     Ding L, Velicer WF, Harlow LL. Effects of estimation methods, number of

indicators per factor, and improper solutions on structural equation modeling fit

indices. Structural Equation Modeling: A Multidisciplinary Journal. 1995

1995/01/01;2(2):119-43.

18.     Nunnally JC, Bernstein IH. Psychometric theory. New York: McGraw-Hill;

1994.

19.     Shrout PE, Fleiss JL. Intraclass correlations: uses in assessing rater reliability.

Psychol Bull. 1979;86(2):420-8.

20.     Mokkink LB, Terwee CB, Patrick DL, Alonso J, Stratford PW, Knol DL, et al.

The COSMIN study reached international consensus on taxonomy, terminology, and

definitions of measurement properties for health-related patient-reported outcomes.

Journal of Clinical Epidemiology. 2010;63(7):737-45.

21.     Streiner DL, Norman GR. Health measurement scales: a practical guide to

their development and use: Oxford University Press, USA; 2008.

22.     AERA, APA, NCME. Standards for educational and psychological testing.

Washington: American Educational Research Association; 1999. p. 194.

23.     Browne MW, Cudeck R. Alternative ways of assessing model fit. Testing

structural equation models. 1993;154:136–62.

24.     Schumacker RE, Lomax RG. A beginner's guide to structural equation

modeling. Mahwah, N.J.: Lawrence Erlbaum Associates; 2004.

25.     Fitzpatrick R, Davey C, Buxton M, Jones D. Evaluating patient-based

outcome measures for use in clinical trials. Health Technology Assessment

1998;2(14):1-74.

26.     Beckerman H, Roebroeck M, Lankhorst G, Becher J, Bezemer P, Verbeek A. Smallest real difference, a link between reproducibility and responsiveness. Quality of Life Research. 2001;10(7):571-8.

27.     De Vet HC, Terwee CB, Ostelo RW, Beckerman H, Knol DL, Bouter LM. Minimal changes in health status questionnaires: distinction between minimally detectable change and minimally important change. Health and Quality of Life Outcomes. 2006;4(1):54.

28.     Ornetti P, Perruccio A, Roos E, Lohmander L, Davis A, Maillefert J. Psychometric properties of the French translation of the reduced KOOS and HOOS (KOOS-PS and HOOS-PS). Osteoarthritis and Cartilage. 2009;17(12):1604-8.

29.     Charter RA, Feldt LS. Confidence intervals for true scores: Is there a correct approach? Journal of Psychoeducational Assessment. 2001;19(4):350-64.

30.     Rothman M, Burke L, Erickson P, Leidy NK, Patrick DL, Petrie CD. Use of Existing Patient-Reported Outcome (PRO) Instruments and Their Modification: The ISPOR Good Research Practices for Evaluating and Documenting Content Validity for the Use of Existing Instruments and Their Modification PRO Task Force Report. Value in Health. 2009;12(8):1075-83.

31.     Mitchell H, Carr A, Scott D. The management of knee pain in primary care: factors associated with consulting the GP and referrals to secondary care. Rheumatology. 2006;45(6):771-6.

32.     Hopman-Rock M, De Bock GH, Bijlsma JW, Springer MP, Hofman A, Kraaimaat FW. The pattern of health care utilization of elderly people with arthritic pain in the hip or knee. International Journal for Quality in Health Care. 1997;9(2):129-37.

33.    Dawson J, Doll H, Fitzpatrick R, Jenkinson C, Carr AJ. The routine use of

patient reported outcome measures in healthcare settings. BMJ (Clinical research ed).

2010;340:c186.

34.    Snyder CF, Aaronson NK, Choucair AK, Elliott TE, Greenhalgh J, Halyard

MY, et al. Implementing patient-reported outcomes assessment in clinical practice: a

review of the options and considerations. Quality of Life Research. 2012;21(8):1305-

14.