

1 **The Oswestry Disability Index, confirmatory factor analysis in a sample of** 2 **35,263 verifies a one-factor structure but practicality issues remain**

3 Charles Philip Gabel^a, Antonio Cuesta-Vargas^b, Meihua Qian^c, Rok Vengust^d, Ulrich Berlemann^e, Emin
4 Aghayev^f, Markus Melloh^{g,h,i}

5
6 a. Coolum Physiotherapy Sunshine Coast, Coolum Beach, Queensland, Australia

7 b. Department of Physiotherapy, Faculty of Medicine, Malaga University, Malaga, Spain

8 c. Department of Education and Human Development, College of Education, Clemson
9 University, 410 Tillman Hall, Clemson, SC 29634, USA

10 d. Department of Orthopaedic Surgery, University Medical Centre Ljubljana, Ljubljana, Slovenia

11 e. dasRückenzentrum, Standort Salem-Spital, Bern, Switzerland

12 f. Swiss RDL, Institute of Social and Preventive Medicine, University of Bern, Bern, Switzerland

13 g. Institute for Health Sciences, School of Health Professions, Zurich University of Applied
14 Sciences, Winterthur, Switzerland

15 h. Faculty of Health Sciences, Curtin University, Perth, Australia

16 i. Centre for Medical Research, University of Western Australia, Nedlands, Australia

17 18 **E-mail addresses**

19 • Meihua Qian - mqian@g.clemson.edu

20 • Charles Philip Gabel - cp.gabel@bigpond.com

21 • Antonio Cuesta-Vargas - acuesta@uma.es

22 • Rok Vengust - rok.vengust@kclj.si

23 • Ulrich Berlemann - uberlemann@hotmail.com

24 • Emin Aghayev - emin.aghayev@ispm.unibe.ch

25 • Markus Melloh - markus.melloh@zhaw.ch

26 **4 Tables and 1 Figure**

27 **Table 1:** Percentiles of Oswestry Disability Index references values (ODI-RV) classified into five
28 categories

29 **Table 2:** Summary of the one-factor solution with or without error covariance using CFA

30 **Table 3:** Sub-group comparisons of CFA outputs—male vs. female participants

31 **Table 4:** Factor loadings from sub-group analyses

32 **Figure 1:** The second model with correlated errors.

33

34 **ABSTRACT**

35 **Purpose**

36 To analyze the factor structure of the Oswestry Disability Index (ODI) in a large symptomatic low back
37 pain (LBP) population using exploratory (EFA) and confirmatory factor analysis (CFA).

38

39 **Methods**

40 Analysis of pooled baseline ODI LBP patient data from the international Spine Tango registry of
41 EUROSPINE, the Spine Society of Europe. The sample, with $n = 35,263$ (55.2% female; age 15–99,
42 median 59 years), included 76.1% of patients with a degenerative disease, and 23.9% of the patients
43 with various other spinal conditions. The initial EFA provided a hypothetical construct for
44 consideration. Subsequent CFA was considered in three scenarios: the full sample and separate
45 genders. Models were compared empirically for best fit.

46

47 **Results**

48 The EFA indicated a one-factor solution accounting for 54% of the total variance. The CFA analysis
49 based on the full sample confirmed this one-factor structure. Subgroup analyses by gender achieved
50 good model fit for configural and partial metric invariance, but not scalar invariance. A possible two-
51 construct model solution as outlined by previous researchers: dynamic-activities (personal care,
52 lifting, walking, sex and social) and static-activities (pain, sleep, standing, travelling and sitting) was
53 not preferred.

54

55 **Conclusions**

56 The ODI demonstrated a one-factor structure in a large LBP sample. A potential two-factor model was
57 considered, but not found appropriate for constructs of dynamic and static activity. The use of the
58 single summary score for the ODI is psychometrically supported. However,. practicality limitations

59 were reported for use in the clinical and research settings. Researchers are encouraged to consider a
60 shift towards newer, more sensitive and robustly developed instruments.

61

62 **Keywords**

63 Oswestry Disability Index, Confirmatory factor analysis, Patient-reported outcome instrument,
64 Validation, Spine Tango, Registry

65

66 INTRODUCTION

67 Measuring and monitoring the individual status and functional change in sufferers of low back pain
68 (LBP) is critical for its overall management ^{1,2}. However, this measurement is not standardized and
69 subsequently cannot systematically reflect the effectiveness of evidence-based interventions. There
70 are over 200 PROs available for LBP measurement with the Oswestry Disability Index (ODI) ^{3,4} one of
71 the most commonly used and advocated in clinical guidelines ^{2,4}. First published in 1980 ³, the ODI was
72 developed to guide treatment programmes and ensure critical LBP aspects were recorded and
73 progress monitored through measured changes in functional status. However, its development
74 followed a qualitative item-selection process rather than a scientific clinimetric methodology ^{3,4,6,7}.
75 Consequently the ODI presents a scale with ‘ordinal’ or ‘preference-based responses’ rather than
76 ‘interval’ or ‘precise measurement points’, which can affect its validity and capacity for standard
77 statistical analysis ⁸. Despite its 40 years of wide use, it has still not been conclusively proven whether
78 the ten ODI items can be summated into a single score ². The result is a lack of consensus regarding its
79 factor structure ⁹⁻¹¹, an important issue that needs resolution.

80 Factor structure is critical and demonstrates the underlying themes or factors present that must be
81 recognized to indicate a parsimonious structure ¹². Factor structure can be singular, enabling a single-
82 summated score; or two- or multi-factor, which requires separately reported scores ^{12,13}. The ODI has
83 always been reported as singular ^{10,14,15}; however, some researchers suggest a two-factor model of:
84 dynamic-activities (personal care, lifting, walking, sex and social) and static-activities (pain, sleep,
85 standing, travelling and sitting) ^{16,17}. With Rasch analysis, which considers the evenness or interval of
86 the scores, a suboptimal one-factor structure was found along with psychometric concerns of poor
87 coverage, plus a large floor and small ceiling effect ^{18,19}. If a PRO is to use a single-summated score, a
88 one-factor solution is required to ensure each question reports upon the same underlying construct
89 ¹¹⁻¹³ according to COSMIN standards ⁷. The gold standard to achieve this is confirmatory factor analysis
90 (CFA) which requires a large dataset for definitive analysis ^{12,20}. A CFA is validating a preceding
91 exploratory factor analysis (EFA), which expose the underlying traits, and requires 50–100 responses-

92 per-item and consequently a minimum sample of $n = 500$ – 1000 for the ODI ¹². There is a gap in the
93 literature as the published studies to date have performed only EFA and only in small samples. In
94 particular, cross-cultural adaptation studies are commonly carried out on samples below $n = 100$ ^{10,17}.
95 This is inadequate for EFA as the estimates become unstable ^{7,12}.
96 Consequently, to address the existing knowledge gaps a single robust study with a large sample size
97 greater than 10,000, or 1000 per item, would be appropriate to resolve the issue conclusively; whether
98 a one- or a multi-factor model has a better fit. The aims of this study were to analyze the ODI factor
99 structure in a LBP population using CFA in an adequately large sample that allows robust testing of
100 competing models, and to determine which model is consistent across genders.

101

102 **METHODS**

103 Ethical approval was not required for this post hoc analysis of anonymous data.

104

105 **Participants**

106 This study was carried out using the Spine Tango data pool. Spine Tango, the international spine
107 registry of EUROSPINE ²¹, the Spine Society of Europe is hosted at the University of Bern's Institute for
108 Social and Preventive Medicine. Completed baseline ODI-PROs ($n = 35,263$, 55.2% female, age = 15–
109 99, median 59-years) were obtained from symptomatic LBP patients included in the registry. The study
110 sample comprised patients with degenerative disease (76.1%), non-generative spondylolisthesis
111 (7.8%), pathological fracture (4.2%), repeat surgery (3.8%), deformity and traumatic fracture (2.7%
112 each), tumour and infection (1% each), and patients with other condition (<0.8%).

113

114 **Assessment tools**

115 The ODI contains ten pain-related, six responses options questions scored from zero (no pain) to five
116 (most severe pain). Scores are expressed as a percentage of total points, with $\leq 20\%$ indicating minimal

117 disability, 21–40% moderate disability, 41–60% severe disability, 61–80% crippled, and 81–100%
118 completely bed-bound ⁴.

119

120 **Factor analysis**

121 The EFA considers several statistics including: Eigenvalues, a special set of characteristic values
122 associated with a linear system of equations (generally >1.0 = statistically relevant); percentage of
123 variance explained by a particular factor ([10% = relevant); factor loading, a measure of how well any
124 item is represented by a factor (>0.30 = minimum); and ‘Scree Plot’, a visual representation chart of
125 Eigenvalues versus items (qualitatively assessed). For PRO’s to provide a one-factor solution and single
126 total score ^{13,15}, each criteria must be fulfilled and a single-factor solution needs to be obtained ¹².

127 When a two-factor solution is argued, the second eigenvalue must be >1 and at least 3–4 items load
128 appropriately on the second factor and also be interpretable. An EFA statistically checks an
129 instrument’s dimensionality where the factor structure must be theoretically meaningful ¹².

130 Subsequent CFA clarifies and validates the suggested EFA model/s using significantly larger samples
131 ¹².

132 Hence this study investigated the ODI factor structure through EFA from a randomly selected 10% sub-
133 group ($n = 3526$) using SPSS 22. Then CFA was conducted on the remaining 90% ($n = 31,736$, 90%)
134 using Mplus 7.11 ²⁰.

135 In CFA, model parameters were estimated using the maximum likelihood method which is robust to
136 non-normality ²⁰. The model fit was assessed using the Root Mean Square Error of Approximation
137 (RMSEA) and the Comparative Fit Index (CFI). A RMSEA value of 0.05 or lower suggests excellent fit,
138 and values 0.08 indicate acceptable fit ²². For the CFI, 0.90 is considered acceptable and 0.95 or above
139 reflects excellent model fit ²³. Additionally, modification indices (MI) were analysed to determine if
140 allowing error terms to co-vary would significantly improve the model fit, and during the CFA, errors
141 with MI exceeding 4.00 were allowed to correlate ²⁰.

142

143 **ODI references values (ODI_RV)**

144 To fully describe the level of severity of participants' disability, an ODI-RV was created.

145

146 **Sub-group analyses**

147 Multi-group analyses were conducted to examine whether the identified model through EFA and CFA
148 fits the data equally well for male and female participants. Namely, the degree to which a confirmatory
149 factor model measuring LBP with ten items per six-point response scale exhibited measurement and
150 structural invariance between male and female participants was assessed using Mplus 7.11²⁰.

151 The original CFA model was first analyzed using the remaining 90% sample. Then the initial configural
152 invariance model was compared with a series of models with increasing invariance constraints.
153 Specifically: (1) the first configural invariance model constrained the pattern of fixed and free
154 parameters to be equivalent across groups; (2) the second metric invariance model constrained factor
155 loadings to be equal across groups; (3) the scalar invariance model constrained all factor loadings and
156 intercepts to be equal across groups; (4) the residual variance invariance model constrained error
157 variance to be equal across groups; (5) the residual covariance invariance model constrained error
158 covariance to be equal across groups; (6) the factor variance invariance model constrained factor
159 variance to be equal across groups; and (7) the factor mean invariance model constrained factor mean
160 to be equal across groups.

161 Invariance between groups on a particular parameter is achieved when non-significant statistical
162 difference is found between a model without a parameter constrained to be equal across groups and
163 the model with the parameter constrained. Then the more parsimonious model is retained and
164 compared to the subsequent model with additional constraints.

165

166 **Assessing competing models**

167 The most common method to assess model equivalence is a Chi-square based Likelihood ratio test,
168 which compares the overall goodness of fit Chi-square values between the two models. However,

169 given Chi-square tests are highly sensitive to trivial differences in large samples ²⁴, other measures,
170 including the Akaike Information Criterion (AIC) and Δ CFI, were also used ²⁵. The Δ CFI was obtained by
171 subtracting the CFI of compared models, where 0.01 indicates a lack of invariance ²⁵. The AIC measures
172 the parsimony of two competing models, where lower values suggest better model fit ²⁶.
173 If a significant, meaningful difference between two compared models exists, then fewer constraints
174 are selected. This indicates a lack of invariance of the parameters in question across groups. The
175 measurement variance across male and female sub-groups was evaluated through multigroup
176 analyses.

177

178 **RESULTS**

179 **Odi_rv**

180 The ODI_RV was calculated from standardized scores classified into five categories: 'minimal',
181 'moderate', 'severe', 'crippling' and 'bed-bound/exaggerated' (Table 1).

182

183 **Explanatory factor analysis**

184 The initial EFA showed a one-factor structure which explained 54% of the total variance. The first
185 eigenvalue was 5.49 and all others were <1.0. Factor loading ranged from 0.58–0.81.

186

187 **Confirmatory factor analysis**

188 The CFA confirmed a one-factor structure. Factor loadings ranged 0.53–0.81. The CFI = 0.945 and
189 RMSEA = 0.075, suggesting adequate model fit. However, further examination of modification indices
190 indicated that allowing some error terms to co-vary would significantly improve model fit (Fig. 1).
191 Hence, the model was re-run to depict the second model with correlated errors (Fig. 1; Table 2). The
192 AIC and RMSEA values of the second model decreased, Δ CFI increased (\sim 0.04) and the difference in
193 Chi-square values between the two models was significant (Table 2). Consequently the second model,
194 with correlated errors, fit the data significantly better than the first model.

195 **Sub-group analyses**

196 Multi-group analyses comparing males ($n = 14,173$) and females ($n = 17,507$) demonstrated configural
197 invariance and partial metric invariance. The configural invariance model had good fit (CFI = 0.983,
198 RMSEA = 0.046), and partial metric invariance was achieved ($\Delta\text{Chi-square}_{\text{configural vs. partial metric}}(2) = 14.022$,
199 $p > 0.05$; $\Delta\text{CFI} < 0.001$; Table 3). Table 4 shows the unstandardized and standardized factor loadings that
200 are statistically similar between male and female (see ODI 2, 4, and 8). However, scalar invariance was
201 not achieved ($\Delta\text{Chi-square}_{\text{partial metric vs. scalar}}(2) = 101.005$, $p < 0.001$), although the DCFI was < 0.01 .

202

203 **DISCUSSION**

204 The findings from both the EFA and CFA confirmed that the ODI's one-factor structure was preferable
205 from both the statistical perspective and parsimony. This is critical as it ensures a valid, single-
206 summated score can be used. No appropriate two-factor model was found that is preferred to the
207 one-factor model, but ambiguity is present. Specifically, the two-factor solution, proposed recently of
208 dynamic and static-activities, was not preferred in the total population or either gender sub-group.
209 This study's findings support previous research for EFA in several samples^{10,15,16}. It also supports the
210 Rasch analysis that found a one-factor structure, but it was suboptimal¹⁸. In our study, while the Chi-
211 square test of the model fit was significant ($p < 0.001$), it is heavily impacted by large sample size and
212 further investigations may be optimal. The gender sub-group analysis indicated both configural
213 invariance and partial metric invariance were obtained between men and women specifying the
214 relationships of some items to the latent factor of disability were equivalent in both groups. However,
215 the scalar invariance was not observed. It suggests women tend to have a slightly higher item response
216 than men at the same absolute trait level of disability. The concerns with the ODI's practicality and
217 consequential clinimetric performance aspects affect both the limitations and implications from
218 clinical and research perspectives^{2,7}. The influence of pain on response options is overwhelming with
219 the iteration of similar optional answers in different sections limiting the patients' ability to express

220 their perceptions of their condition ^{9,11}. This is reflected in the large minimum detectable change
221 (MDC) and minimum clinically relevant difference (MCID), which determine responsiveness and
222 error ^{7,11}. These have been demonstrated in previous studies to be around 20–25% of baseline level
223 ^{1,9,11}. This is insufficient in comparison to several other regional PROs for which the MDC is in the order
224 of 10% or lower, and numerical rating scales have errors of around 15% in the same sample and
225 require only a single question¹⁴.

226 Consequently, the ODI as a modern viable PRO is less practical than simpler PROs that are easier to
227 use and have smaller error scores that reduce the ‘number needed to treat’ (NTT). This, consequently,
228 determines a smaller sample size and shorter time to provide meaningful results that verifies if true
229 change has occurred and ensures statistically significant outcomes for both the individual and
230 investigative research. The ODI is also unable to include objective parameters which limit post-
231 operative evaluation ^{1,11}. By comparison, recent computer based PROs have such values represented
232 or transferred into response options and algorithms that calculate a final single outcome score ⁸. The
233 practicality aspect of ‘patient demand’ to complete a PRO, expound the potential for completion
234 errors and inconsistency ^{10,11}. These include excessive completion time and scoring inaccuracies, a
235 consequence of a large number of response options and increased cognitive demand, that leads to
236 respondent uncertainty and reduced precision ^{1,9,11}. Solutions to overcome these issues include
237 shortening the PRO, modifications to improve practicality, modern scientific development
238 methodology ^{10,11} and a shift toward digital software systems such as computerized adaptive testing
239 (CAT) or computerized decision support systems (CDSS) ²⁷ in future randomized controlled trials that
240 incorporate objective and individual response options ^{1,11}.

241

242 **LIMITATIONS AND STRENGTHS**

243 This study’s limitations are several. As a secondary analysis, diagnostic sub-groups (e.g., spinal
244 stenosis, radiculopathy or disc degeneration) could not be considered due to limited diagnostic codes
245 within the data set. The implications of potential constructs of ‘dynamic’ and ‘static’ function, as

246 suggested by some researchers ¹⁷, could potentially have been present within the participants'
247 occupational, social, sporting or daily routine. However, this could not be ascertained from the
248 available data set. It is highly unlikely, from the statistical findings, that such considerations potentially
249 influenced the analysis. If so then this would affect the overall validity of the ODI in terms of the
250 capability of providing a single-summed score.

251 The dominant strength of this study is the very large sample size. The 10% EFA sample alone was over
252 tenfold larger than all previous factor analysis studies. This is certainly one of the important benefits
253 of registries besides implant tracking, detection of rare adverse events, early warning, benchmarking,
254 real-life perspective and so forth ²¹. Furthermore, a statistician independent of the data collectors is
255 responsible for the data analysis.

256

257 **CONCLUSION**

258 The findings are conclusive that the one-factor solution is preferable from the perspectives of both
259 the statistical analysis and parsimony. Consequently, the ongoing use of the ODI summary score is
260 psychometrically supported. However, the ODI, as an outcome instrument, continues to have
261 prominent limitations that include practicality and measurement error. Clinicians must be aware of
262 the completion burden for patients, and that a minimum detectable change is around 20–25% of the
263 baseline level. This may have consequences on the research. Researchers are encouraged to consider
264 a shift towards newer, more sensitive and robustly developed instruments.

265

266 **ACKNOWLEDGMENTS**

267 The participants of the Spine Tango Register are acknowledged for their continuous contribution that
268 makes possible such studies reflecting the daily practise of spine surgeons. The data of the following
269 centres were used (in alphabetic order of country, city, hospital and department): Dept. of Spinal
270 Surgery in Royal Adelaide Hospital (Australia); Dept. of Spinal Surgery in St. Andrew's Hospital in
271 Adelaide (Australia); Dept. of Orthopaedic Surgery in Landeskrankenhaus Krems (Austria); Dept. of

272 Orthopaedic Surgery in Orthopaedic Hospital Speising in Vienna (Austria); Dept. of Orthopaedic
273 Surgery in University Hospital St. Luc in Brussels (Belgium); Dept. of Orthopaedic Surgery in Grand
274 Hôpital de Charleroi (Belgium); Dept. of Neurosurgery in University Hospital Cologne (Germany); Dept.
275 of Orthopaedic Surgery in University Hospital of Cologne (Germany); Dept. of Spine Surgery in Hospital
276 Dortmund (Germany); Dept. of Orthopaedic Surgery in University Hospital of Greifswald (Germany);
277 Dept. of Spine Surgery and Neurotraumatology in St. Nobifacius Hospital Lingen (Germany); Dept. of
278 Neurosurgery in Orthopädisches Klinikum Markgröningen (Germany); Dept. of Neurosurgery in
279 Klinikum Offenbach (Germany); Dept. of Orthopaedic Surgery in Asklepios Klinikum Uckermark in
280 Schwedt (Germany); Dept. of Special Spine Surgery in Leopoldina Hospital of Schweinfurt (Germany);
281 Dept. of Spine Surgery in Krankenhaus der Barmherzigen Brüder in Trier (Germany); Dept. of Spine
282 Surgery in Clinica Cellini (Italy); Dept. of Spine Surgery in IRCCS Galeazzi in Milan (Italy); Dept. of
283 Neurosurgery in Sapienza University of Rome (Italy); Dept. of Spine Surgery in Centro Medico Puerta
284 de Hierro (Mexico); Dept. of Spine Surgery in SCTO in Chisinau (Moldova); Dept. of Neurosurgery in
285 Wojewódzki Szpital Specjalistyczny nr 2 in Jastrzębie-Zdrój (Poland); Dept. of Neurosurgery in
286 Specialized Medical Center S.A. Polanica (Poland); Dept. of Neurosurgery in Medical University Silesia
287 (Poland); Dept. of Neurosurgery in General Hospital Torun (Poland); Dept. of Orthopaedic Surgery and
288 Traumatology in Kliniczny in Wroclaw (Poland); Dept. of Orthopaedic Surgery in Tan Tock Seng
289 Hospital (Singapore); Dept. of Orthopaedic Surgery in University Hospital of Ljubljana (Slovenia); Dept.
290 of Neurosurgery in Bethesda Hospital of Basel (Switzerland); Dept. of Spine Surgery in Bethesda
291 Hospital of Basel (Switzerland); Dept. of Orthopaedic Surgery in Salem Hospital of Bern (Switzerland);
292 Dept. of Spine Surgery in University Hospital in Lausanne (Switzerland); Dept. of Spine Surgery in
293 Hirslandenklinik Birshof in Münchenstein (Switzerland); Dept. of Spine Surgery in The Spine Center
294 Thun (Switzerland); Dept. of Orthopaedic Surgery in Hospital Schwyz (Switzerland); Dept. of Spine
295 Surgery in Nottingham University Hospitals NHS Trust (UK); Spine Unit of Nuffield Oxford Centre (UK);
296 Division of Orthopaedic Surgery in SUNY Downstate Medical Center in New York (USA); Division of
297 Spine Surgery in NYU Hospital of New York (USA).

298 **Compliance with ethical statement.**

299

300 **Conflict of interest**

301 None of the authors has any potential conflict of interest.

302

303 **REFERENCES**

- 304 1. Cleland JA, Gillani R, Bienen EJ, Sadosky A (2011) Assessing dimensionality and responsiveness of
305 outcomes measures for patients with low back pain. *Pain Pract* 11(1):57–69
- 306 2. Chiarotto A, Maxwell LJ, Terwee CB, Wells G, Tugwell P, Ostelo R (2016) Roland-Morris Disability
307 Questionnaire and Oswestry Disability Index: which has better measurement properties for
308 measuring physical functioning in nonspecific low back pain? Systematic review and meta-
309 analysis. *Phys Ther* 96(10):1620–1637
- 310 3. Fairbank JCT, Couper J, Davies JB, O'Brien JP (1980) The oswestry low back pain disability
311 questionnaire. *Physiotherapy* 66(8):271–273
- 312 4. Fairbank JCT, Pynsent PB (2000) The Oswestry Disability Index. *Spine* 25(22):2940–2952
- 313 5. Guzman JZ, Cutler HS, Connolly J et al (2016) Patient-reported outcome instruments in spine
314 surgery. *Spine (Phila Pa 1976)* 41(5):429–437
- 315 6. Ostelo RW, Deyo R, Stratford P et al (2008) Interpreting change scores for pain and functional
316 status in low back pain: towards international consensus regarding minimal important change.
317 *Spine* 33(1):90–94
- 318 7. Terwee CB, Mokkink LB, Knol DL, Ostelo RW, Bouter LM, de Vet HC (2012) Rating the
319 methodological quality in systematic reviews of studies on measurement properties: a scoring
320 system for the COSMIN checklist. *Qual Life Res* 21(4):651–657
- 321 8. Gabel CP, Cuesta-Vargas AI, Osborne JO, Burkett B, Melloh M (2014) Confirmatory factory
322 analysis of the neck disability index indicates a one-factor model. *Spine J* 14(8):1410–1416
- 323 9. Mehra A, Baker D, Disney S, Pynsent PB (2008) Oswestry Disability Index scoring made easy. *Ann*
324 *R Coll Surg Engl* 90(6):497–499
- 325 10. Eranki V, Koul K, Fagan A (2013) Rationalization of outcome scores for low back pain: the Oswestry
326 disability index and the low back outcome score. *ANZ J Surg* 83(11):871–877
- 327 11. Gabel CP, Melloh M, Burkett B, Michener LA (2013) The Spine Functional Index: development and
328 clinimetric validation of a new whole-spine functional outcome measure. *Spine J*.
329 doi:10.1016/j.spinee.2013.09.055
- 330 12. Costello AB, Osborne J (2005) Best practices in exploratory factor analysis: four recommendations
331 for getting the most from your analysis. *Pract Assess, Res Eval* 10(7):1–9
- 332 13. Doward LC, McKenna SP (2004) Defining patient-reported outcomes. *Value Health* 7(S1):S4–S8
- 333 14. Hägg O, Fritzell P, Nordwall A, SLSS Group (2003) The clinical importance of changes in outcome
334 scores after treatment for chronic low back pain. *Eur Spine J* 12(1):12–20
- 335 15. van Hooff ML, Spruit M, Fairbank JC, van Limbeek J, Jacobs WC (2015) The Oswestry Disability
336 Index (version 2.1a): validation of a Dutch language version. *Spine (Phila Pa 1976)* 40(2):E83–E90

- 337 16. Guermazi M, Mezghani M, Ghroubi S et al (2005) The Oswestry index for low back pain translated
338 into Arabic and validated in a Arab population. [Article in French]. *Ann Readapt Med Phys*
339 48(1):1–10
- 340 17. Tan K, Zheng M, Yang BX et al (2009) Validating the Oswestry Disability Index in patients with low
341 back pain in Sichuan. [Article in Chinese]. *Sichuan Da Xue Xue Bao Yi Xue Ban* 40(3):559–561
- 342 18. Brodke DS, Goz V, Lawrence BD, Spiker WR, Neese A, Hung M (2016) Oswestry Disability Index: a
343 psychometric analysis with 1,610 patients. *Spine J* 17(3):321–327
- 344 19. Terwee CB, Bot SD, de Boer MR et al (2007) Quality criteria were proposed for measurement
345 properties of health status questionnaires. *J Clin Epidemiol* 60(1):34–42
- 346 20. Muthe´n LK, Muthe´n BO (1998–2015) *Mplus user’s guide*, 7th edn. Muthe´n & Muthe´n, Los
347 Angeles
- 348 21. Staub LP, Ryser C, Ro¨der C et al (2016) Total disc arthroplasty versus anterior cervical interbody
349 fusion: use of the Spine Tango registry to supplement the evidence from randomized control
350 trials. *Spine J* 16(2):136–145
- 351 22. Schumacher RE, Lomax RGA (1996) *A beginner’s guide to structural equation modeling*. Lawrence
352 Erlbaum, Mahwah
- 353 23. McDonald RP, Marsh HW (1990) Choosing a multivariate model: noncentrality and goodness of
354 fit. *Psychol Bull* 107:247–255
- 355 24. Cheung GW, Rensvold RB (2002) Evaluating goodness-of-fit indexes for testing measurement
356 invariance. *Struct Equ Model* 9(2):233–255
- 357 25. Byrne BM (2010) *Structural equation modeling with AMOS: basic concepts, applications, and*
358 *programming*, 2nd edn. Routledge, New York
- 359 26. Akaike H (1974) A new look at the statistical model identification. *IEEE Trans Autom Control*
360 19(6):716–723
- 361 27. Moja L, Passardi A, Capobussi M et al (2016) Implementing an evidence-based computerized
362 decision support system linked to electronic health records to improve care for cancer patients:
363 the ONCO-CODES study protocol for a randomized controlled trial. *Implement Sci* 11(1):153
364

365 **TABLES**366 **Table 1**

367 Percentiles of Oswestry Disability Index references values (ODI-RV) classified into five categories

	ODI-RV			Disability categories
	Total (<i>n</i> = 35,249)	Male (<i>n</i> = 15,801)	Female (<i>n</i> = 19,448)	
Percentile				
<20th	-0.902	-1.004	-0.799	Minimal
<40th	-0.284	-0.387	-0.284	Moderate
<60th	0.230	0.126	0.229	Severe
<80th	0.847	0.847	0.949	Crippling
<99th	2.390	2.390	2.287	Bed-bound/exaggerated

368

369 **Table 2**

370 Summary of the one-factor solution with or without error covariance using CFA.

371 χ^2 value indicates the difference between observed variance–covariance matrix and the model-implied variance–covariance matrix; p value indicates
 372 probability of the difference; and df stands for the degrees of freedom. *RMSEA* the root mean square error of approximation, is a measure of model fit, with
 373 a value of 0.05 or lower suggesting excellent fit, and values <0.08 indicating reasonable fit ²⁴; CFI stands for the Comparative Fit Index, with 0.90 being
 374 considered acceptable, and 0.95 or above reflecting excellent model fit ²⁴. AIC, the Akaike Information Criterion, is a comparative measure of fit, with lower
 375 values indicating a better model fit²⁵. $\Delta\chi^2$ is the difference in Chi-square values between the first model and the second model with correlated errors.
 376 Correlated errors in the second model represent that the unique variances of the associated indicators such as pain intensity and sleeping overlap (see Fig. 1
 377 for details)

Model	χ^2	df	p	CFI	RMSEA	AIC	Significance of $\Delta\chi^2$
First model	6942.724	35	<0.001	0.945	0.075	1019349.168	
Second model with correlated errors	2083.422	29	<0.001	0.983	0.045	1013579.776	$P < 0.001$

378

379 **Table 3**

380 Sub-group comparisons of CFA outputs—male vs. female participants

381 χ^2 value indicates the difference between observed variance–covariance matrix and the model implied variance–covariance matrix; p value indicates
 382 probability of the difference; and df stands for the degrees of freedom. RMSEA the root mean square error of approximation, is a measure of model fit, with
 383 a value of 0.05 or lower suggesting excellent fit, and values <0.08 indicating reasonable fit 24. CFI stands for the comparative fit index, with 0.90 being
 384 considered acceptable, and 0.95 or above reflecting excellent model fit 24. $\Delta\chi^2$ (14.022) is the difference in Chi-square values between the configural model
 385 and partial metric model, and $\Delta\chi^2$ (101.005) is the difference in Chi-square values between the partial metric model and scalar model. Partial metric invariance
 386 was achieved ($p>0.05$), whereas scalar invariance was not achieved ($p>0.001$)

Model	χ^2	df	p	$\Delta\chi^2$	CFI	p	RMSEA
Configural model	2186.526	58	<0.001		0.983		0.046
Partial metric model (item 2, 4 and 8)	2200.548	60	<0.001	14.022	0.983	>0.05	0.045
Scalar model	2301.553	62	<0.001	101.005	0.982	<0.001	0.045

387

388 **Table 4**

389 Factor loadings from sub-group analyses. * Factor loadings held equal across groups

Item	Unstandardized factor loading		Standardized factor loading	
	Males	Females	Males	Females
ODI 1	0.687	0.661	0.620	0.587
ODI 2*	0.796	0.796	0.720	0.694
ODI 3	0.936	1.005	0.694	0.702
ODI 4*	0.972	0.972	0.690	0.700
ODI 5	0.743	0.676	0.598	0.552
ODI 6	0.906	0.963	0.652	0.674
ODI 7	0.642	0.595	0.565	0.505
ODI 8*	1.054	1.054	0.788	0.780
ODI 9	1.195	1.283	0.813	0.813
ODI 10	1.310	1.412	0.764	0.762

390

391 **FIGURES**

392 **Figure 1**

393 The second model with correlated errors. Disability represents the ODI, and m1–10 stand for pain
 394 intensity, personal care, walking, lifting, sitting, standing, sleeping, social life, travelling, and sex life

