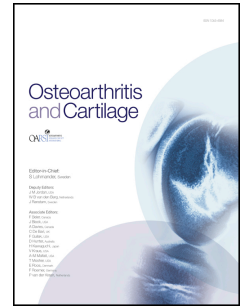


Journal Pre-proof

Estimating incidence and prevalence of hip osteoarthritis using electronic health records: a population-based cohort study

Ilgin G. Arslan, Jurgen Damen, Marcel de Wilde, Jacqueline J. van den Driest, Patrick J.E. Bindels, Johan van der Lei, Sita M.A. Bierma-Zeinstra, Dieuwke Schiphof



PII: S1063-4584(22)00674-4

DOI: <https://doi.org/10.1016/j.joca.2022.03.001>

Reference: YJOCA 5022

To appear in: *Osteoarthritis and Cartilage*

Received Date: 30 September 2021

Revised Date: 1 March 2022

Accepted Date: 7 March 2022

Please cite this article as: Arslan IG, Damen J, de Wilde M, van den Driest JJ, Bindels PJE, van der Lei J, Bierma-Zeinstra SMA, Schiphof D, Estimating incidence and prevalence of hip osteoarthritis using electronic health records: a population-based cohort study, *Osteoarthritis and Cartilage*, <https://doi.org/10.1016/j.joca.2022.03.001>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2022 The Author(s). Published by Elsevier Ltd on behalf of Osteoarthritis Research Society International.

Estimating incidence and prevalence of hip osteoarthritis using electronic health records: a population-based cohort study

Ilgin G. Arslan¹, Jurgen Damen¹, Marcel de Wilde², Jacoline J. van den Driest¹, Patrick J.E.

Bindels¹, Johan van der Lei², Sita M.A. Bierma-Zeinstra^{1,3}, Dieuwke Schiphof¹

¹*Department of General Practice, Erasmus MC University Medical Center, Rotterdam, The Netherlands*

²*Department of Medical Informatics, Erasmus University, Rotterdam, The Netherlands*

³*Department of Orthopaedics, Erasmus MC, University Medical Center, Rotterdam, The Netherlands*

Address for correspondence:

Ilgin G. Arslan

Postal address: P.O. Box 2040, 3000 CA Rotterdam, The Netherlands

E-mail: i.arslan@erasmusmc.nl

Phone number: +31 (0)10-7037741

ORCID ID: <https://orcid.org/0000-0003-2046-6177>

No specific funding was received from any bodies in the public, commercial or not-for-profit sectors to carry out the work described in this article. The authors have declared no conflicts of interest.

ABSTRACT

Objective: To determine the incidence and prevalence of hip osteoarthritis (OA) in electronic health records (EHRs) of Dutch general practices by using narrative and codified data.

Method: A retrospective cohort study was conducted using the Integrated Primary Care Information database. An algorithm was developed to identify patients with narratively diagnosed hip OA in addition to patients with codified hip OA. Incidence and prevalence estimates among people aged ≥ 30 were assessed from 2008 to 2019. The association of comorbidities with codified hip OA diagnosis was analysed using multivariable logistic regression.

Results: Using the hip OA narrative data algorithm (positive predicted value=72%) in addition to codified hip OA showed a prevalence of 1.76 to 1.95 times higher and increased from 4.03% in 2008 to 7.34% in 2019. The incidence was 1.83 to 2.41 times higher and increased from 6.83 to 7.78 per 1000 person-years from 2008 to 2019. Among codified hip OA patients, 39.4% had a previous record of narratively diagnosed hip OA, on average approximately 1.93 years earlier. Hip OA patients with a previous record of spinal OA, knee OA, hypertension, and hyperlipidaemia were more likely to be recorded with a hip OA code.

Conclusion: This study using Dutch EHRs showed that epidemiological estimates of hip OA are likely to be an underestimation. Using our algorithm, narrative data can be added to codified data for more realistic epidemiological estimates based on routine healthcare data. However, developing a valid algorithm remains a challenge, possibly due to the diagnostic complexity of hip pain in general practice.

Keywords: hip osteoarthritis, incidence, prevalence, epidemiology, electronic health records

Journal Pre-proof

1 INTRODUCTION

2 Osteoarthritis (OA) is one of the most prevalent joint diseases and has been ranked as the 10th
3 leading contributor to global disability.¹⁻³ The hip joint is often affected by OA, as it is one of the
4 most weight-bearing joints of the human body.⁴ In 2017, the global prevalence of hip OA was
5 estimated at 40 million people and the global incidence at 2 million people.⁵ There is no cure for
6 hip OA and current treatment focuses on reducing symptoms and improving function.⁶ The only
7 effective treatment is a joint replacement as an end-stage, which accounts for the majority of the
8 healthcare costs associated with hip OA.⁷ In 2017, 18.3% of the total healthcare costs for
9 musculoskeletal diseases in the Netherlands was due to OA.⁸ This is expected to increase due to
10 the ageing of the population and increasing obesity rate.⁵

11 Current incidence and prevalence of OA are estimated using primary care electronic
12 health records (EHRs) from routine healthcare data, largely focused on codified data containing
13 specific codes for specific diseases.⁹⁻¹³ However, EHRs also contain narrative data that include
14 free text notes from healthcare providers. In a previous study¹⁴ using primary care EHRs from the
15 Netherlands, we found that a substantial proportion of knee OA patients did not have a record
16 of codified knee OA, but had a record of a knee OA diagnosis in the free text of their EHR. Adding
17 these narratively diagnosed knee OA patients to codified knee OA patients yielded approximately
18 twofold higher prevalence and incidence estimates. Problems with under-recording of OA were
19 also found in UK primary care EHRs.¹⁵ Several reasons may contribute to this problem, such as
20 GPs giving lower priority to record diseases or symptoms¹⁵⁻¹⁷, which is likely in patients with OA
21 as multimorbidity is common¹⁸.

22 While misclassifications and under-recordings may have major impact on the accuracy of
23 epidemiological estimates, healthcare policy of hip OA is still based on epidemiological estimates
24 obtained from routine healthcare data using codified data alone. More accurate information on
25 epidemiological estimates is urgently needed to adequately respond to the large increase of the
26 burden of hip OA.⁵

27 Therefore, this study aimed to determine the incidence and prevalence of hip OA using
28 the complete EHR consisting of both codified and narrative data from a routine primary care
29 database in the Netherlands.

30

31

32 **METHODS**

33 **Design and setting**

34 This retrospective cohort study was conducted using the Integrated Primary Care Information
35 (IPCI) database which contains EHRs from Dutch general practices of approximately 2.5 million
36 patients. Details of this database have been published elsewhere.^{19, 20} In summary, EHRs from
37 the IPCI database comprise all medical journal entries written in free text by GPs, diagnoses using
38 the International Classification of Primary Care (ICPC) codes, laboratory findings, drug
39 prescriptions, referrals, and correspondence with other healthcare providers from primary and
40 secondary care (e.g. physiotherapist and orthopaedic surgeon). EHRs from the IPCI database
41 contain the majority of patients' medical information, as all citizens in the Netherlands are
42 obliged to register with a GP which acts as the first point of contact and the gatekeeper to
43 secondary care.^{21, 22}

44

45 **Study cohort**

46 We used a similar research method for the development of an algorithm based on narrative data
47 to identify under-recorded hip OA patients as we did in an earlier study¹⁴ in which we examined
48 the under-recording of knee OA. Patients were included during each study year from 1 January
49 2008 until 31 December 2019 if they were aged ≥ 30 with at least 12 months of valid database
50 history prior to the study entry. Patients with a codified diagnosis of hip OA were selected. The
51 codified diagnosis of hip OA was based on the ICPC code L89.

52 In addition, an algorithm was developed by our research group, including GPs, to identify
53 patients with keywords referring to hip OA in narrative data (i.e. the free text in their EHR)
54 without any record of codified hip OA (ICPC code L89). An overview of our workflow is illustrated
55 in Figure 1. In the first phase, the algorithm included patients with an ICPC code L13 (i.e. hip
56 complaints) plus keywords related to OA or keywords related to hip plus OA without ICPC code
57 L13, for example 'hip' plus 'osteoarthritis'. Keywords combined with terms indicating negation
58 (e.g. 'not' or 'no') were excluded, as were combinations with relatives (e.g. 'father has', 'mother
59 has'), patient's anxiety about a possible diagnosis of OA, and expressions of uncertainties
60 regarding the OA diagnosis by the GP or other healthcare providers in primary care or secondary
61 care (e.g. 'probably', 'differential diagnoses'). A random sample of 100 patients identified by the
62 algorithm was assessed by one author (IGA) to check for terminology variations and misspellings
63 of keywords. Textual alternations were made after discussion with all authors to improve the
64 algorithm.

65 In the second phase, we randomly selected 50 patients of these potential narratively
66 diagnosed hip OA patients without a record of codified hip OA. These cases were assessed on
67 true and false positive for having hip OA through a blinded medical record review by two authors,
68 IGA (physiotherapist and researcher) and JD (academic GP). True positive cases were defined by:
69 “Patients where the GP, healthcare provider from primary care (e.g. physiotherapist) or
70 secondary care (e.g. orthopaedist or radiologist) reported a hip OA diagnosis in the free text in
71 their EHR, with or without X-ray imaging”; a commonly used and generally accepted reference
72 standard.¹⁶ When the hip OA diagnosis was documented in a radiology report only,
73 documentation of hip pain in the EHR at the time of X-ray or MRI request was required to classify
74 as a true positive hip OA case. Hip OA as an incidental finding on X-ray or MRI after a traumatic
75 event was not considered as a true positive case, given the poor correlation between the severity
76 of structural damage of the joint and the severity of symptoms^{23, 24}. Consensus was reached
77 through discussion with the last author (DS, senior researcher experienced with IPCI database).
78 Results were then discussed with the research group and modifications to the algorithm were
79 made to reduce the number of false positive cases.

80 In the last phase, the positive predicted value (PPV) of the modified narrative data
81 algorithm was re-assessed using the same methods as in the second phase. To compare the
82 validity of the algorithm with that of codified hip OA, one author (IGA) assessed the PPV of a
83 random selection of 50 patients identified with codified hip OA (i.e. ICPC code L89) with the same
84 methods as for the PPV assessment of narratively diagnosed hip OA and with scrutiny by the co-
85 authors (JD or DS) if necessary. Different random samples of patients were used for all three
86 phases in the algorithm development process.

87

88 **Outcomes**

89 PPVs were calculated as the proportion of patients who were confirmed as having hip OA, based
90 on the information reported in the EHR. The annual lifetime prevalence was calculated as the
91 total number of people ever diagnosed as at 1 July each calendar year, divided by the total
92 number of patients in the population on that date, and multiplied by 100. The entire
93 retrospective record available for patients was used to estimate the prevalence. The annual
94 incidence rate was calculated by the number of new cases between 1 January and 31 December
95 (i.e. no previous diagnosis of hip OA) in each calendar year, divided by the number of person
96 years at risk between 1 January and 31 December each calendar year. This at risk period is the
97 period that a patient participated in the IPCI database without a recorded hip OA diagnosis until
98 the moment of death, changing practice, hip OA diagnosis, or end of participation in the IPCI
99 database. The entire retrospective record available for patients was used to exclude prior hip OA
100 when estimating the incidence rates. Thus, patients with a hip OA diagnosis in their medical
101 history (i.e. medical history before enrolment in the IPCI database or before 1 January 2008) were
102 defined as prevalent cases. See Supplementary File S1 for more information regarding the
103 medical history available for the study cohort. Prevalence and incidence estimates were
104 calculated separately for: 1) patients with codified hip OA diagnosis defined as at least one ICPC
105 code hip OA (i.e. L89), and 2) patients with narratively diagnosed hip OA according to the free-
106 text algorithm without any record of codified hip OA in their EHR. Incidence and prevalence
107 estimates were calculated stratified by sex. Further details of the study design are illustrated in
108 Figure 2.

109 To determine the effect of including narrative data in addition to codified data, annual
110 rate ratios between prevalence and incidence estimates of codified hip OA and codified plus
111 narratively diagnosed hip OA were calculated.

112 Furthermore, some of the patients identified with codified hip OA may have been
113 identified with hip OA at an earlier date based on narrative data. We explored the proportion of
114 patients with a narrative hip OA diagnosis prior to a codified hip OA diagnosis. The number of
115 days between the first narrative hip OA diagnosis and the first codified hip OA diagnosis was
116 calculated.

117 We explored differences in demographics and comorbidities between patients with
118 codified hip OA and patients with narratively diagnosed hip OA. In addition, based on previous
119 research¹⁵⁻¹⁷, we hypothesized that GPs may give patients with comorbidities lower priority to
120 also record OA with a code. Therefore, we analysed the association between concurrent
121 comorbidities (i.e. occurring before the first hip OA diagnosis) and codified hip OA among all
122 prevalent hip OA patients. Prevalent hip OA patients are either codified or narratively diagnosed
123 between 1 January 2008 and 31 December 2019. Narratively diagnosed hip OA patients are the
124 reference category of the outcome in this analyses. We selected the following common
125 comorbidities in patients with OA from an earlier systematic review¹⁸: 1) hypertension,
126 hyperlipidaemia, overweight, diabetes mellitus (i.e. disorders related to metabolic syndrome); 2)
127 heart/vascular diseases and events (i.e. stroke/TIA, peripheral arterial disease, and myocardial
128 infarction/angina pectoris), 3) asthma, 4) Chronic Obstructive Pulmonary Disease (COPD), 5) a
129 small selection of OA related to joints other than the hip (i.e. spinal OA and knee OA), 8) low back

130 pain. For the comorbidities we used the codified diagnosis based on ICPC-codes (see
131 Supplementary Table S2 for the full list of ICPC-codes). This analysis was adjusted for age and sex.

132

133 **Statistics**

134 Binomial 95% confidence intervals (CIs) were calculated for the PPVs. Prevalence and incidence
135 estimates were standardized for age and sex using the annual distribution for the whole Dutch
136 population as given by the StatLine database of Statistics Netherlands from 2008 up to 2019²⁵.

137 The Poisson distribution was used to provide 95% CIs for prevalence and incidence estimates.

138 Descriptive characteristics were reported as means and standard deviations (SDs), medians and
139 interquartile ranges (IQRs), and counts (n) and percentages (%), as appropriate. Multivariable

140 logistic regression was performed to determine the association of comorbidities with the codified
141 diagnosis among patients with hip OA (either narratively diagnosed or codified diagnosed); the

142 results were expressed as odds ratios (ORs) including 95% CIs. The significance level throughout

143 was set at two-tailed $P < .05$. Statistical analyses were performed using R Studio Software V.4.0.2.

144

145

146 **RESULTS**

147 **Validity assessment**

148 *Narrative data algorithm*

149 An overview of our workflow for the development of the narrative data algorithm is illustrated
150 in Figure 1 and full details in Supplementary Data S3. The first version of the algorithm yielded a
151 PPV of 60% (95%CI = 46.4% to 73.6%) (Phase 2). False positive cases were found frequently due
152 to codified hip complaints (i.e. ICPC code L13) plus keywords for OA in the lower back or sacroiliac
153 joint, and were therefore excluded in the second revised algorithm. We also excluded the
154 keyword 'prosthesis', as this was often found after a hip fracture and not due to hip OA.
155 Subsequently, the PPV of this final narrative data algorithm resulted into 72% (95% CI = 59.6% to
156 84.4%) (Phase 3). In the final algorithm, false positive cases were still frequently found due to
157 keywords for OA in the lower back or sacroiliac joint in combination with a keyword related to
158 the hip joint or codified hip complaints, but also due to unclear diagnosis of hip OA and hip OA
159 as an incidental finding on X-ray to rule out a hip fracture after traumatic event. For 80.6% (29
160 out of 36) of the true-positive narratively diagnosed hip OA patients, an X-ray was used to confirm
161 the diagnosis, either requested by the GP or documented in the correspondence from an
162 orthopaedic surgeon in secondary care to the GP.

163 *Codified hip OA diagnosis*

164 The PPV of codified diagnosed hip OA was 98% (95% CI= 94.1% to 100%). The reason for the false
165 positive case was a coding error where the GP recorded the ICPC code L89 (hip OA) instead of
166 L90 (knee OA). For 87.8% (43 out of 49) of the true-positive codified diagnosed hip OA patients,
167 an X-ray was used to confirm the diagnosis, either requested by the GP or documented in the
168 correspondence from an orthopaedic surgeon, rheumatologist, internist, or urologist in
169 secondary care to the GP.

170

171 **Study cohort**

172 The study cohort consisted of 117,758 patients with hip OA. A total of 63,470 patients had a
173 record of codified hip OA with a mean age of 68.2 (SD=11.7) and 34.3% were men. The remaining
174 54,288 patients did not have any record of codified hip OA, but were identified with narratively
175 diagnosed hip OA alone. These patients were younger (mean age=65.4 (SD=12.8)) and comprised
176 a slightly greater percentage of men (36.0%) compared to codified hip OA patients.

177

178 *Narrative diagnosis prior to codified diagnosis*

179 Of the patients identified with codified hip OA, 39.4% (n=25030) was at an earlier time point
180 diagnosed narratively with hip OA; on average 1.93 years earlier (median number of days = 706;
181 IQR = 48 to 2378).

182

183 **Prevalence**

184 The standardized prevalence of codified hip OA in 2008 was 2.07% (95%CI 2.06-2.08) and
185 increased to 4.01% (95%CI 4.00-4.02) in 2019 (Figure 3A). The standardized prevalence of
186 narratively diagnosed hip OA alone (i.e. without any record of codified hip OA) was estimated to
187 be 1.96% (95%CI 1.96-1.97) in 2008 and increased to 3.33% (95%CI 3.32-3.34) in 2019 (Figure
188 3B). The annual crude and standardized prevalence proportions are presented in Supplementary
189 Table S4, as well as the accurate number of included people each year in analysis.

190 Adding narrative data to codified data showed prevalence proportions with a rate ratio
191 between 1.76 and 1.95 during the study period (Table 1) and increased from 4.03% (95%CI 4.02-
192 4.04) in 2008 to 7.34% (95%CI 7.32-7.35) in 2019 (Figure 4A).

193

194 **Incidence**

195 The standardized incidence of codified hip OA declined from 3.74 per 1000 person-years (95%CI
196 3.70-3.78) in 2008 to 3.22 per 1000 person-years (95%CI 3.19-3.25) in 2019 (Figure 5A) and
197 peaked in 2013 with 4.19 per 1000 person-years (95%CI 4.15-4.23). In contrast, the standardized
198 incidence of narratively diagnosed hip OA alone increased consistently year by year with 2.72 per
199 1000 person-years (95%CI 2.68-2.75) in 2008 to 3.86 per 1000 person-years (95%CI 3.82-3.89) in
200 2019 (Figure 5B). The annual crude and standardized incidence rates are presented in
201 Supplementary Table S4.

202 Adding narrative data to codified data showed incidence rates with a rate ratio between
203 1.83 and 2.41 during the study period (Table 1). The incidence increased from 6.83 per 1000
204 person-years (95%CI 6.78-6.88) in 2008 to 7.78 per 1000 person-years (95%CI 7.78-7.83) in 2019
205 and was highest in 2011 with 7.89 per 1000 person-years (95%CI 7.84-7.94) (Figure 4B).

206 Prevalence and incidence estimates for all case definitions were at any given time point
207 higher for women than for men. Sex stratified estimates are presented in Supplementary Table
208 S5.

209

210 **Factors associated with a record of codified hip OA**

211 In general, multivariable analysis showed small to no statistically significant associations of
212 demographic variables and concurrent comorbidities with codified hip OA (Figure 6). Among the
213 concurrent comorbidities, spinal OA (OR 1.13 [95%CI 1.07-1.19]), knee OA (OR 1.10 [95%CI 1.05-
214 1.14]), hyperlipidaemia (OR 1.11 [95%CI 1.07-1.15]), and hypertension (OR 1.10 [95%CI 1.07-
215 1.13]) were associated with a record of codified hip OA. Concurrent stroke/TIA, diabetes, and low
216 back pain reduced the likelihood of being recorded with codified hip OA, but with small
217 associations. The remaining comorbidities showed no statistically significant associations. Full
218 details are provided in Supplementary Table S6.

219

220

221 **DISCUSSION**

222 This study developed an algorithm to determine the incidence and prevalence of hip OA in EHRs
223 of Dutch general practices by using a combination of narrative and codified data. Adding narrative
224 data based on this algorithm to codified data showed prevalence and incidence estimates of
225 almost twice as many on average from 2008-2019. Our algorithm had a positive predicted value
226 of 72%. False positive cases mainly occurred due to keywords for OA in the lower back or
227 sacroiliac joint combined with keyword related to the hip joint or codified hip complaints, unclear
228 diagnosis of hip OA, and hip OA as an incidental finding on X-ray to rule out a hip fracture after
229 traumatic event. Contrary to current guidelines^{24,26-29}, an X-ray was used to confirm the diagnosis
230 in most of the hip OA patients.

231 A previous record of spinal OA and knee OA showed a positive association with codified
232 hip OA. It may be that GPs are more prone to record hip OA with a code when the patient is
233 already known to have OA in joints other than the hip. Furthermore, a previous record of
234 hyperlipidaemia and hypertension increased the likelihood of hip OA patients being recorded
235 with a hip OA code. The Dutch healthcare system includes reimbursement schemes for
236 cardiovascular risk management. Patients included in this program are routinely invited to visit
237 their GP to monitor their health status, including screening on hypertension and hyperlipidaemia.
238 It may be that patients who are routinely monitored are more likely to have a record of codified
239 hip OA. Previous research¹⁵⁻¹⁷ hypothesized that GPs may under-record codified OA because they
240 give it lower priority than other diseases. Although we found that a record of concurrent
241 stroke/TIA, diabetes, and low back pain reduced the likelihood of hip OA patients being recorded
242 with codified hip OA, these associations were too small to support this hypothesis.

243 The current study found that hip OA was increasingly under-recorded over time, since the
244 incidence of codified hip OA diagnosis decreased over time, while that of narratively diagnosed
245 hip OA alone increased. However, it should be noted that these patients with narratively
246 diagnosed hip OA alone may be recorded with codified hip OA in the future, since almost 40% of
247 codified hip OA patients had a previous record of narratively diagnosed hip OA. In contrast, Swain
248 et al.¹¹ found an increase of codified hip OA and a decrease of codified 'unspecified' OA over time
249 in EHRs from the UK. The authors suggested that this may be due to better recording of codified
250 hip OA, since hip OA patients are increasingly being recorded with codified hip OA rather than
251 unspecified OA.

252 Similar to our previous study¹⁴ on knee OA, the current study showed that adding
253 narrative data to codified data yielded almost twice as many hip OA patients than the standard
254 approach of using codified data alone. However, the development of the algorithm to identify
255 narratively diagnosed hip OA patients in the current study was more complex than for narratively
256 diagnosed knee OA patients in our previous study¹⁴. The algorithm for hip OA included false-
257 positive cases resulting from keywords for spinal OA combined with hip complaints, which was
258 not present in the knee OA algorithm. This can be explained by a strong association of low back
259 pain with hip OA compared to knee OA.³⁰ Also, false-positive cases in the hip OA algorithm
260 occurred due to keywords for hip prosthesis after a hip fracture rather than for hip OA. These
261 false-positive cases were not present in the knee OA algorithm, as arthroplasty is far more
262 commonly used in patients with acute femur fracture than in knee fractures.^{31, 32} Although
263 exclusion of these combinations increased the PPV from 60% to 72%, the validity of the narrative
264 data algorithm for hip OA remained lower than for knee OA (i.e. PPV=94%). This reflects the
265 greater clinical diagnostic challenge of hip OA compared to knee OA. The differential diagnosis of
266 hip pain presented to a GP is much broader than in knee pain, e.g. hip pain is sometimes difficult
267 to distinguish from trunk pain and is often associated with a variety of hip conditions, such as OA,
268 gluteal tendinopathy, and femoral acetabular impingement syndrome.³³⁻³⁵ While current
269 guidelines do not recommend imaging to diagnose OA in clinical practice, but recommend using
270 history taking and physical examination instead^{24, 26-29}, we found in the current study that an X-
271 ray was used for most hip OA patients to confirm the diagnosis. This overuse of X-rays for
272 diagnosing hip OA in the general practice may reflect the clinical diagnostic complexity of hip OA.

273 It may also indicate the demand of patients, asking their GP to confirm a likely chronic diagnosis
274 with potential major implications for the patient.

275 Furthermore, similar to the findings in our previous study¹⁴ on knee OA, around 40% of
276 the codified hip OA patients in the current study had a previous record of a narrative diagnosis.
277 Capturing hip OA patients earlier may help policymakers to plan and prioritize resources more
278 adequately to keep healthcare affordable. Remarkably, the time between the narrative diagnosis
279 and codified diagnosis was shorter for hip OA than for knee OA (1.9 years vs 3 years,
280 respectively).¹⁴ This difference may relate to findings from a previous research in which the
281 symptom duration at the time of initial presentation was found to be shorter for hip OA than for
282 knee OA (2.7 years and 3.9 years, respectively).³⁶ However, to date, the reason for this difference
283 in clinical presentation is unclear.

284 A previous study¹⁵ found an under-recording of codified OA in UK primary care EHRs in a
285 quarter of severe OA patients aged 40 with total hip and knee replacements. However, these
286 results do not apply to the less severe OA patients (i.e. without joint replacement) where under-
287 recording may be even more present since patients with less severe OA are less likely to have a
288 codified OA diagnosis³⁷. To the best of our knowledge, the current study is the first that presented
289 the under-recording of hip OA across the entire spectrum of severity.

290 The Dutch National Institute for Public Health and the Environment (RIVM) published
291 prevalence and incidence estimates of codified hip OA based codified data alone retrieved from
292 Nivel Primary Care Registrations.¹³ Comparing their estimates with our results is difficult because
293 of the differences in age restriction. We therefore reproduced our analyses without restriction

294 on age as estimates published by RIVM, which showed similar estimates; i.e. crude prevalence in
295 2019, 1.97% for men and 3.44% for women in the current study versus 1.96% for men and 3.34%
296 for women published by RIVM. Nevertheless, estimates published by RIVM are probably
297 underestimated, since they only include codified hip OA patients.

298 A strength of this study is the use of a representative sample of the Dutch population
299 from IPCI database.^{19, 20} Limitations of this study include that, although we captured a substantial
300 part of under-recorded hip OA patients by adding narrative data to codified data, our prevalence
301 and incidence estimates might still be an underestimation due to the restrictiveness of the
302 algorithm. On the other hand, the PPV of 72% of the narrative data algorithm might imply an
303 overestimation of 28% of the hip OA patients identified with narrative data, as they possibly do
304 not have hip OA. In addition, we were able to calculate the PPV of the diagnoses, but not other
305 features of the algorithm, such as negative predicted value or sensitivity, and future research on
306 this is required. Also, an important aspect to consider when interpreting our results is that under-
307 recording of hip OA could be related to several factors, such as the type of general practice and
308 the type of information systems, as Dutch GPs are free to choose among competing information
309 systems that significantly differ in user interfaces and features²⁰. Future research into this is
310 warranted to better understand factors contributing to under-recording of diseases in routine
311 healthcare data.

312 Current healthcare policy on prevention and management is based on routine primary
313 care data using codified data alone from EHRs. Findings from the current study and previous
314 studies^{14, 15} demonstrating the under-recording of OA indicate a serious underestimation of

315 epidemiological estimates and other estimates obtained from EHR-based studies (i.e. association
316 studies, descriptive management policy studies). This leads to inaccurate outcomes and
317 eventually inaccurate healthcare policy making. Narrative data can be added to codified data in
318 EHR-based OA research. In that way, policy makers will have a more realistic picture of the
319 current and future burden of OA and can better respond to its predicted large increase.⁵
320 However, it should be noted that the use of narrative data may not always be feasible, since
321 coding systems and the use of narrative data fields built into EHRs may differ between countries
322 and systems. Data protection may even limit access to narrative data fields, making other
323 alternatives to identify under-recorded hip OA patients in EHR data more suitable, for example
324 using process, referral, and intervention codes. In addition, developing an algorithm based on
325 patient characteristics (i.e. age and occupation) in combination with symptomatic codes (i.e. hip
326 complaints ICPC code L13 in the Netherlands) may potentially help to identify patients with OA
327 in joints without an OA code.

328

329 **CONCLUSIONS**

330 This study developed an algorithm to determine the incidence and prevalence of hip OA in EHRs
331 of Dutch general practices by using a combination of narrative and codified data. The positive
332 predicted value of narratively diagnosed hip OA patients alone was 72%. Adding narrative data
333 to codified data yielded prevalence and incidence estimates of almost twice as many on average
334 from 2008-2019. A previous record of spinal OA, knee OA, hypertension, and hyperlipidaemia
335 increased the likelihood of hip OA patients being recorded with a hip OA code. This study showed

336 the importance of using narrative data in addition to codified data in EHR-based OA research to
337 produce realistic epidemiologic estimates. However, developing a valid algorithm to identify hip
338 OA patients based on narrative data remains a challenge, possibly due to the diagnostic
339 complexity of hip pain in general practice.

Journal Pre-proof

ACKNOWLEDGEMENTS

None.

AUTHOR CONTRIBUTIONS

IGA, JD, MdW, JJvdD, PJEB, DS, and SMAB-Z participated in the design of the study. IGA, JD, and DS reviewed electronic health records for validity assessment. IGA conducted statistical analysis. IGA, JD, MdW, JJvdD, PJEB, JvdL, DS, and SMAB-Z gave their comment on the first version of the manuscript and approval of the final manuscript.

ROLE OF FUNDING SOURCE

No specific funding was received from any bodies in the public, commercial or not-for-profit sectors to carry out the work described in this article.

CONFLICT OF INTEREST

The authors have declared no conflicts of interest.

ETHICAL APPROVAL INFORMATION

This study was approved by the Board of Directors of the IPCI database.

DATA SHARING STATEMENT

The aggregated data are available on request from the corresponding author.

REFERENCES

1. World Health Organization. Musculoskeletal conditions. WHO 2020.
2. Hunter DJ, Bierma-Zeinstra S. Osteoarthritis. *Lancet* 2019; 393: 1745-1759.
3. Osteoarthritis Research Society International. Osteoarthritis: A Serious Disease. 2016: p. 103.
4. Zhang Y, Jordan JM. Epidemiology of osteoarthritis. *Clin Geriatr Med* 2010; 26: 355-369.
5. Global Burden Disease, Injury I, Prevalence C. Global, regional, and national incidence, prevalence, and years lived with disability for 354 diseases and injuries for 195 countries and territories, 1990-2017: a systematic analysis for the Global Burden of Disease Study 2017. *Lancet* 2018; 392: 1789-1858.
6. Hunter DJ, Bierma-Zeinstra S. Osteoarthritis. *The Lancet* 2019; 393: 1745-1759.
7. Pivec R, Johnson AJ, Mears SC, Mont MA. Hip arthroplasty. *Lancet* 2012; 380: 1768-1777.
8. Dutch National Institute for Public Health and the Environment (RIVM). Health care expenses osteoarthritis. vol. 20212021.
9. Cross M, Smith E, Hoy D, Nolte S, Ackerman I, Fransen M, et al. The global burden of hip and knee osteoarthritis: estimates from the global burden of disease 2010 study. *Ann Rheum Dis* 2014; 73: 1323-1330.
10. Spitaels D, Mamouris P, Vaes B, Smeets M, Luyten F, Hermens R, et al. Epidemiology of knee osteoarthritis in general practice: a registry-based study. *BMJ Open* 2020; 10: e031734.
11. Swain S, Sarmanova A, Mallen C, Kuo CF, Coupland C, Doherty M, et al. Trends in incidence and prevalence of osteoarthritis in the United Kingdom: findings from the Clinical Practice Research Datalink (CPRD). *Osteoarthritis Cartilage* 2020; 28: 792-801.
12. Turkiewicz A, Petersson IF, Bjork J, Hawker G, Dahlberg LE, Lohmander LS, et al. Current and future impact of osteoarthritis on health care: a population-based study with projections to year 2032. *Osteoarthritis Cartilage* 2014; 22: 1826-1832.
13. National Institute for Public Health and the Environment N. Public Health Foresight Study 2018 (VTV-2018): diseases. 2018.
14. Arslan IG, Damen J, de Wilde M, van den Driest JJ, Bindels PJE, van der Lei J, et al. Incidence and prevalence of knee osteoarthritis using codified and narrative data from electronic health records: a population-based study. *Arthritis Care Res (Hoboken)* 2022.
15. Yu D, Jordan KP, Peat G. Underrecording of osteoarthritis in United Kingdom primary care electronic health record data. *Clin Epidemiol* 2018; 10: 1195-1201.
16. Shrestha S, Dave AJ, Losina E, Katz JN. Diagnostic accuracy of administrative data algorithms in the diagnosis of osteoarthritis: a systematic review. *BMC Med Inform Decis Mak* 2016; 16: 82.
17. Jencks SF, Williams DK, Kay TL. Assessing Hospital-Associated Deaths From Discharge Data: The Role of Length of Stay and Comorbidities. *Jama* 1988; 260: 2240-2246.
18. Swain S, Sarmanova A, Coupland C, Doherty M, Zhang W. Comorbidities in Osteoarthritis: A Systematic Review and Meta-Analysis of Observational Studies. *Arthritis Care Res (Hoboken)* 2020; 72: 991-1000.
19. Vlug AE, van der Lei J, Mosseveld BM, van Wijk MA, van der Linden PD, Sturkenboom MC, et al. Postmarketing surveillance based on electronic patient records: the IPCI project. *Methods Inf Med* 1999; 38: 339-344.
20. van der Lei J, Duisterhout JS, Westerhof HP, van der Does E, Cromme PV, Boon WM, et al. The introduction of computer-based patient records in The Netherlands. *Ann Intern Med* 1993; 119: 1036-1041.
21. Kroneman M, Boerma, W., Van den Berg, M., Groenewegen, P., De Jong, J., Van Ginneken, E. The Netherlands: health system review., vol. 18. *Health Systems in Transition* 2016:1-239.

22. Kringos D, Boerma W, Bourgueil Y, Cartier T, Dedeu T, Hasvold T, et al. The strength of primary care in Europe: an international comparative study. *Br J Gen Pract* 2013; 63: e742-750.
23. Lawrence JS, Bremner JM, Bier F. Osteo-Arthrosis: Prevalence in the Population and Relationship between Symptoms and X-ray Changes. *Annals of the Rheumatic Diseases* 1966; 25: 1-24.
24. Sakellariou G, Conaghan PG, Zhang W, Bijlsma JWW, Boyesen P, D'Agostino MA, et al. EULAR recommendations for the use of imaging in the clinical management of peripheral joint osteoarthritis. *Annals of the Rheumatic Diseases* 2017; 76: 1484-1494.
25. CBS Open Data Statline. Population dynamics: month and year.
26. Kolasinski SL, Neogi T, Hochberg MC, Oatis C, Guyatt G, Block J, et al. 2019 American College of Rheumatology/Arthritis Foundation Guideline for the Management of Osteoarthritis of the Hand, Hip, and Knee. *Arthritis Care Res (Hoboken)* 2020; 72: 149-162.
27. Bannuru RR, Osani MC, Vaysbrot EE, Arden NK, Bennell K, Bierma-Zeinstra SMA, et al. OARSI guidelines for the non-surgical management of knee, hip, and polyarticular osteoarthritis. *Osteoarthritis and Cartilage* 2019; 27: 1578-1589.
28. National Institute for Health & Clinical Excellence. Osteoarthritis: the care and management of osteoarthritis in adults. London 2014.
29. Nederlands Huisartsen Genootschap. Niet-traumatische knieklachten Nederlands Huisartsen Genootschap 2016.
30. Stupar M, Côté P, French MR, Hawker GA. The association between low back pain and osteoarthritis of the hip and knee: a population-based cohort study. *J Manipulative Physiol Ther* 2010; 33: 349-354.
31. Bohm ER, Tufescu TV, Marsh JP. The operative management of osteoporotic fractures of the knee: To fix or replace? *Journal of Bone and Joint Surgery - Series B* 2012; 94 B: 1160-1169.
32. Ries MD. Primary arthroplasty for management of osteoporotic fractures about the knee. *Current Osteoporosis Reports* 2012; 10: 322-327.
33. Ferguson RJ, Prieto-Alhambra D, Walker C, Yu D, Valderas JM, Judge A, et al. Validation of hip osteoarthritis diagnosis recording in the UK Clinical Practice Research Datalink. *Pharmacoepidemiol Drug Saf* 2019; 28: 187-193.
34. Birrell F, Lunt M, Macfarlane GJ, Silman AJ. Defining hip pain for population studies. *Annals of the rheumatic diseases* 2005; 64: 95-98.
35. Hall M, van der Esch M, Hinman RS, Peat G, de Zwart A, Quicke JG, et al. How does hip osteoarthritis differ from knee osteoarthritis? *Osteoarthritis and Cartilage* 2022; 30: 32-41.
36. Dabare C, Le Marshall K, Leung A, Page CJ, Choong PF, Lim KK. Differences in presentation, progression and rates of arthroplasty between hip and knee osteoarthritis: Observations from an osteoarthritis cohort study-a clear role for conservative management. *International Journal of Rheumatic Diseases* 2017; 20: 1350-1360.
37. Jordan KP, Tan V, Edwards JJ, Chen Y, Englund M, Hubertsson J, et al. Influences on the decision to use an osteoarthritis diagnosis in primary care: a cohort study with linked survey and electronic health record data. *Osteoarthritis Cartilage* 2016; 24: 786-793.

Under-recording of hip osteoarthritis

FIGURES LEGENDS

Figure 1. Workflow diagram for the development of the narrative data algorithm

Figure 2. Details of the study design

Figure 3. Standardized prevalence of hip OA based on codified data (A) and narrative data alone (B)

Figure 4. Standardized (A) prevalence and (B) incidence of hip OA based narrative data alone in addition to codified data

Figure 5. Standardized incidence of hip OA based on codified data (A) and narrative data alone (B)

Figure 6. Characteristics associated with codified hip OA diagnosis among all hip OA patients (either codified diagnosed or narratively diagnosed without a hip OA code)

TABLES

Table 1. Prevalence and incidence of hip OA based on codified data versus a combination of codified and narrative data

Standardized prevalence [95% CI]				Standardized incidence [95% CI]			
Year	Codified data	Codified + narrative data	Rate ratio	Year	Codified data	Codified + narrative data	Rate ratio
2008	2.07 [2.06-2.08]	4.03 [4.02-4.04]	1.95	2008	3.74 [3.70-3.78]	6.83 [6.78-6.88]	1.83
2009	2.23 [2.22-2.24]	4.19 [4.18-4.21]	1.88	2009	3.82 [3.79-3.86]	7.08 [7.03-7.14]	1.85
2010	2.43 [2.42-2.44]	4.52 [4.51-4.54]	1.87	2010	3.90 [3.86-3.93]	7.48 [7.43-7.53]	1.92
2011	2.67 [2.66-2.68]	4.97 [4.96-4.99]	1.86	2011	4.08 [4.04-4.11]	7.89 [7.84-7.94]	1.94
2012	2.86 [2.85-2.87]	5.28 [5.27-5.30]	1.85	2012	4.04 [4.01-4.08]	7.51 [7.46-7.56]	1.86
2013	3.07 [3.06-3.08]	5.50 [5.48-5.51]	1.79	2013	4.19 [4.15-4.23]	7.76 [7.70-7.81]	1.85
2014	3.30 [3.29-3.31]	5.83 [5.82-5.85]	1.77	2014	3.87 [3.83-3.90]	7.42 [7.37-7.47]	1.92
2015	3.47 [3.46-3.48]	6.11 [6.09-6.12]	1.76	2015	3.70 [3.66-3.74]	7.50 [7.45-7.55]	2.03
2016	3.62 [3.61-3.63]	6.41 [6.40-6.43]	1.77	2016	3.49 [3.46-3.53]	7.22 [7.17-7.27]	2.07
2017	3.76 [3.75-3.77]	6.71 [6.69-6.72]	1.78	2017	3.56 [3.52-3.59]	7.68 [7.63-7.73]	2.16
2018	3.92 [3.90-3.93]	7.06 [7.04-7.07]	1.80	2018	3.39 [3.36-3.43]	7.46 [7.41-7.51]	2.20
2019	4.01 [4.00-4.02]	7.34 [7.32-7.35]	1.83	2019	3.22 [3.19-3.25]	7.78 [7.72-7.83]	2.41

Note: Standardized prevalence and incidence estimates are standardized for age and sex distribution of the total population from the Netherlands.

FIGURES

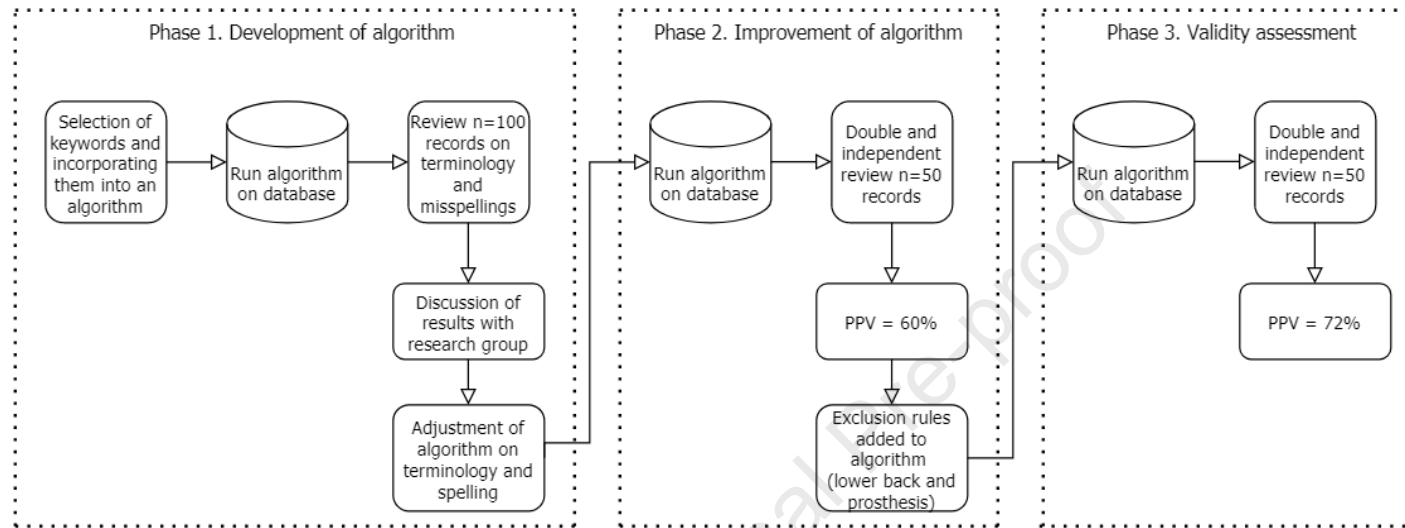


Figure 1. Workflow diagram for the development of the narrative data algorithm

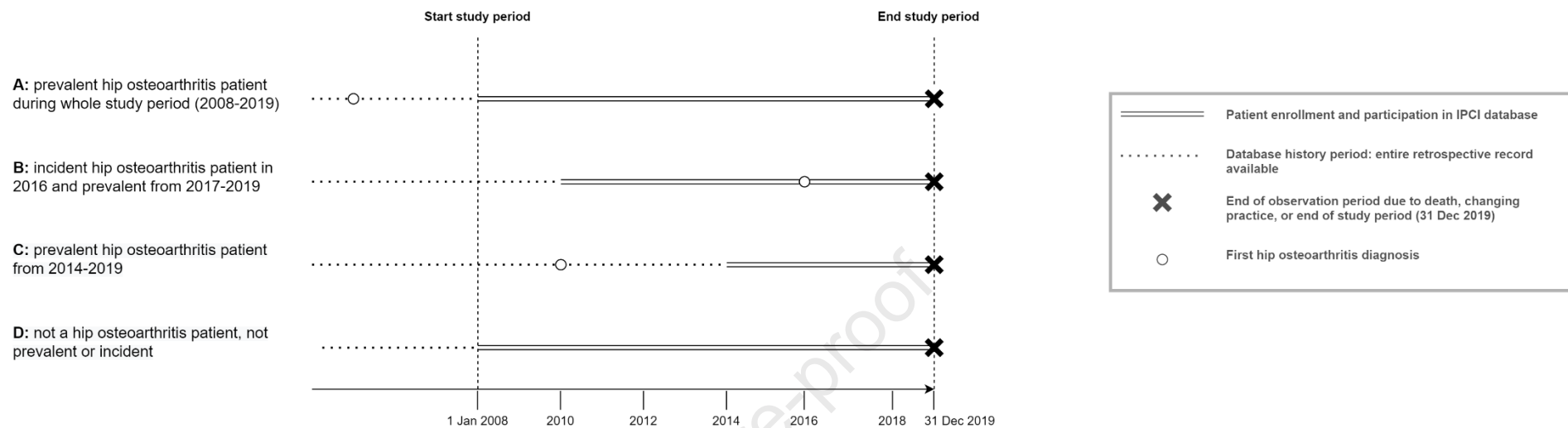
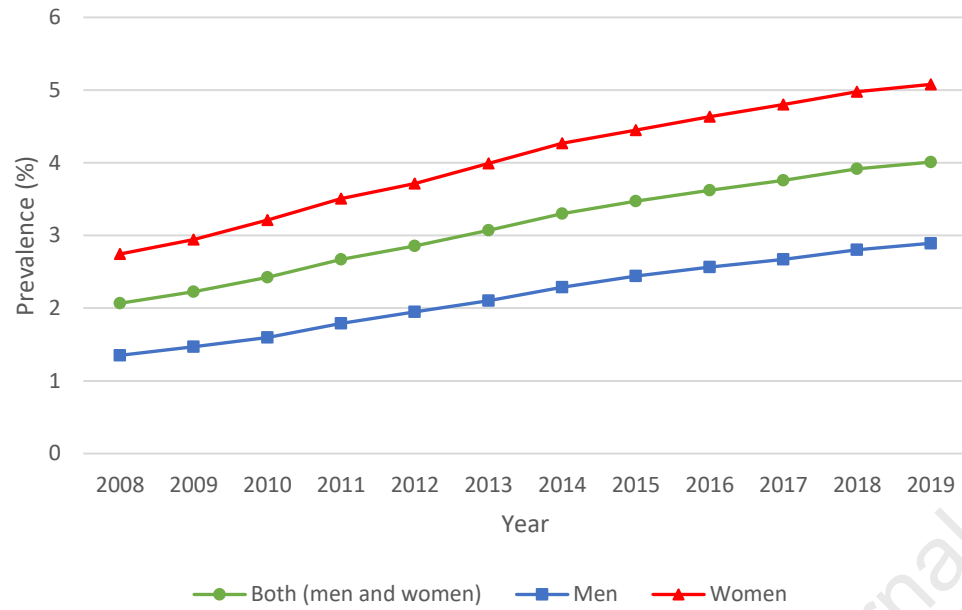


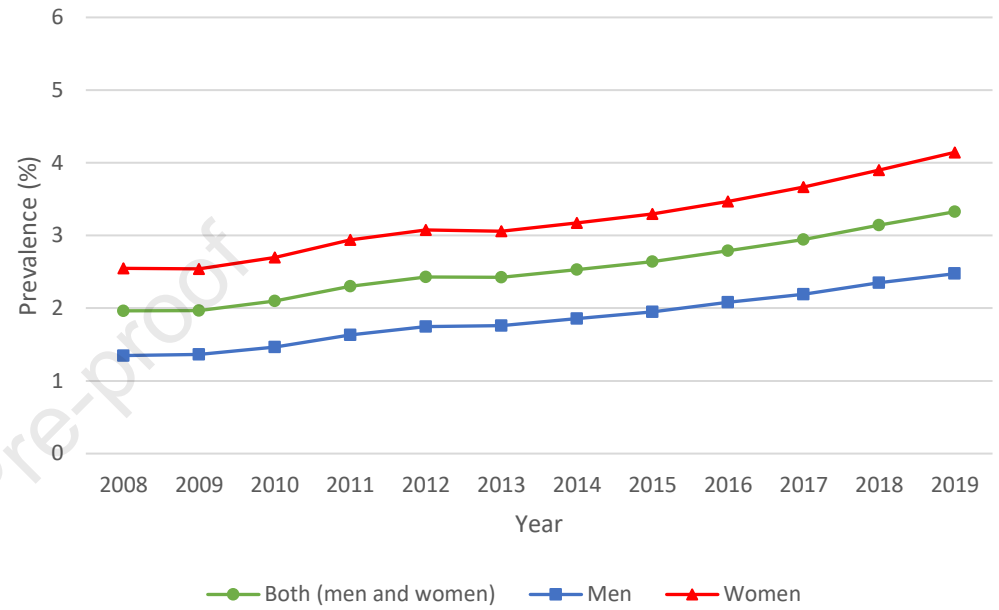
Figure 2. Details of the study design

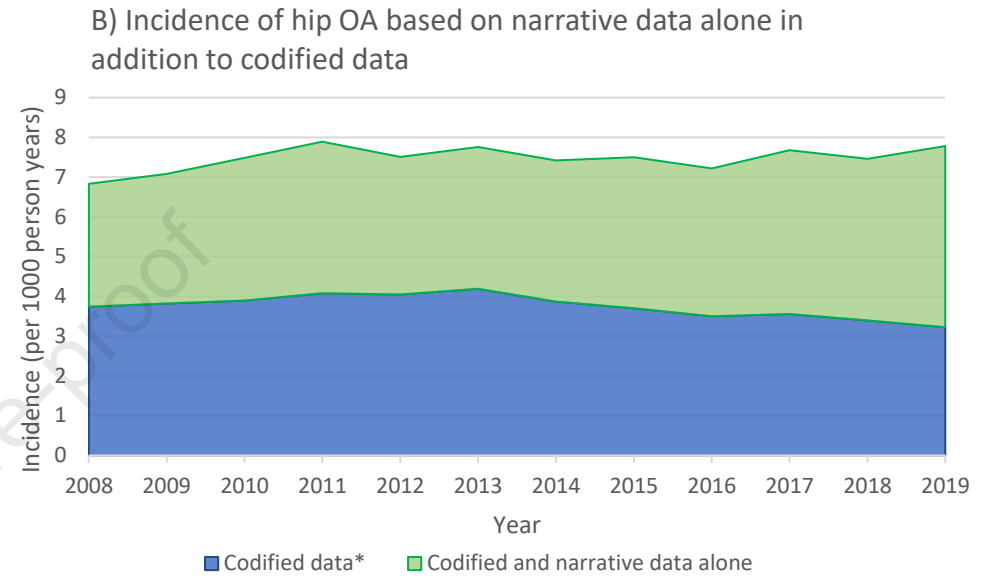
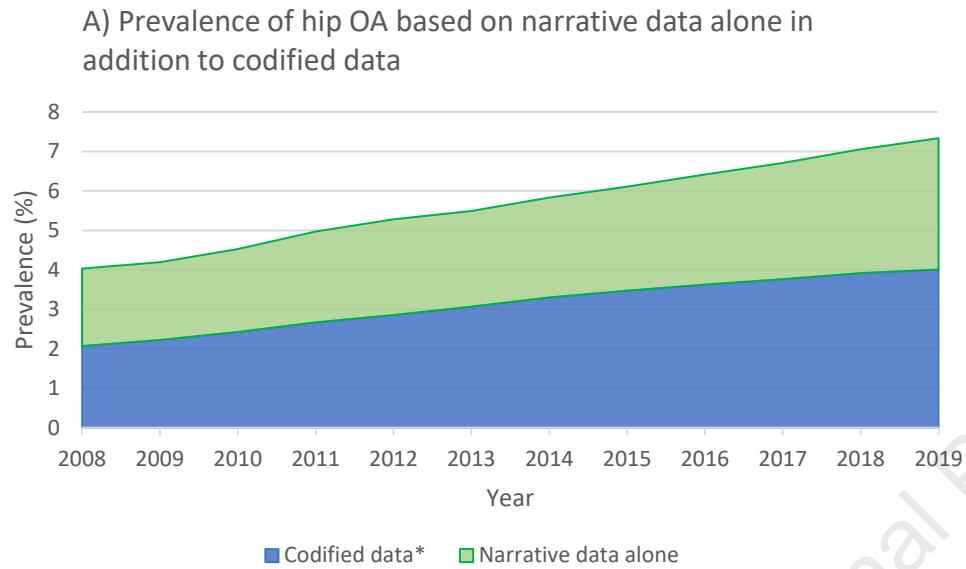
Notes. Figure 2 shows four examples of patients in the study cohort (A-D). The study period started on 1 January 2008 until 31 December 2019. The IPCI database is an open cohort, meaning that patients can also enter the database after the start of study period and stop before the end of study period due to death or changing practice. Patients were followed from the start of study period (patient A and patient D) or from the moment they entered the IPCI database if this moment was after 1 January 2008 (patient B and patient C). Patients were followed until the end of the study period (patient A, B, C and D) or until the moment of death or changing practice when this moment was before 31 December 2019. A first hip OA diagnosis was defined as incident when the first diagnosis was given within the study period and participation in IPCI database (patient B). The incidence rate was calculated annually by the number of new cases in each calendar year, divided by the number of person years at risk between in each calendar year. For example, when calculating the incidence rate of the year 2016, patient B is included in the numerator and patient B and D are included in the denominator. The prevalence was calculated annually as the total number of people ever diagnosed as at 1 July each calendar year, divided by the total number of patients in the population on that date, and multiplied by 100. For example, when calculating the prevalence of the year 2014, patient A and C are included in the numerator and patient A-D are included in the denominator.

A) Prevalence of hip OA based on codified data



B) Prevalence of hip OA based on narrative data alone

**Figure 3.** Standardized prevalence of hip OA based on codified data (A) and narrative data alone (B)



* Among patients identified with codified hip OA, 39.4% were previously diagnosed narratively with hip OA, which was approximately 1.9 years prior to the first codified hip OA diagnosis. These patients are not counted in the annual lifetime prevalence proportions of narrative data alone.

Figure 4. Standardized (A) prevalence and (B) incidence of hip OA based narrative data alone in addition to codified data

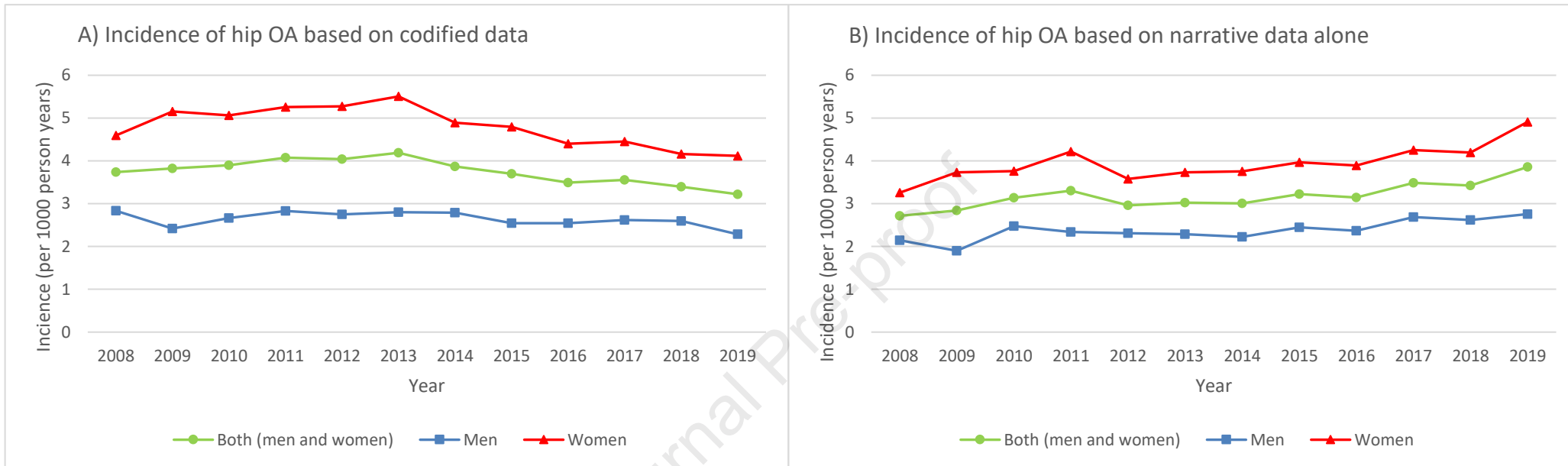


Figure 5. Standardized incidence of hip OA based on codified data (A) and narrative data alone (B)

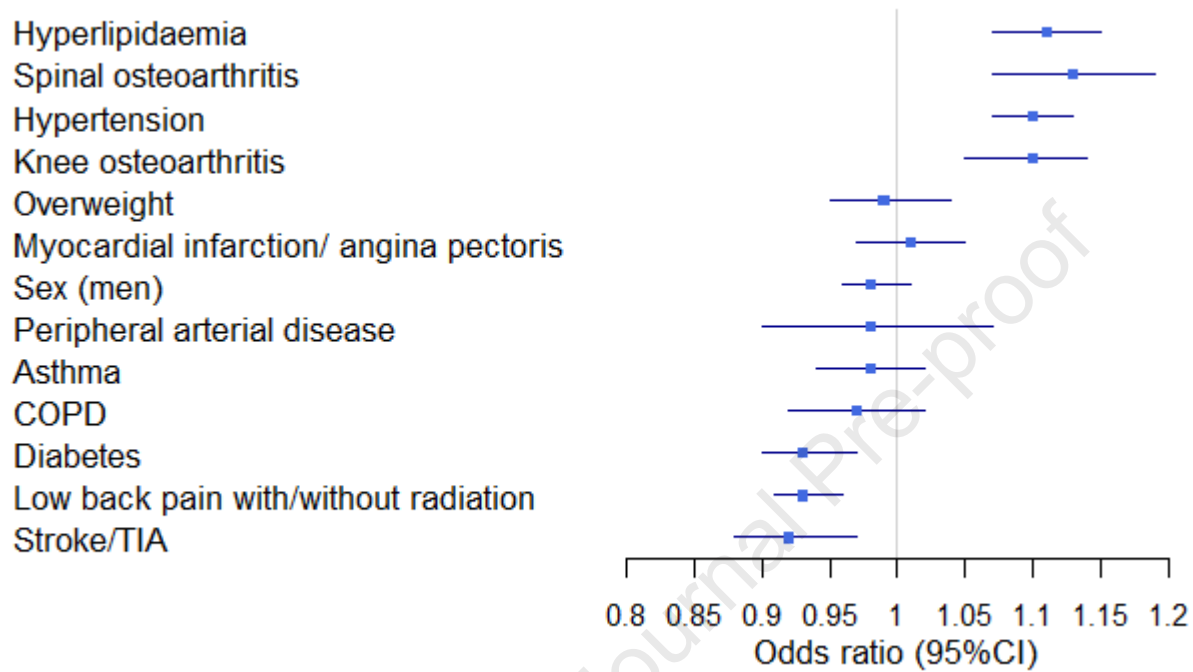


Figure 6. Characteristics associated with codified hip OA diagnosis among all hip OA patients (either codified diagnosed or narratively diagnosed without a hip OA code)

Note. Full details are provided in Supplementary Table S6.