# Identifying Parkinson's disease and parkinsonism cases using routinely-collected healthcare data: a systematic review

Authors: Zoe Harding[1¶†], Tim Wilkinson[2,3¶*], Anna Stevenson[4,5], Sophie Horrocks[1], Amanda Ly[3], Christian Schnier[3], David P Breen[6], Kristiina Rannikmäe[2,3], Cathie LM Sudlow[2,3], on behalf of Dementias Platform UK

1. College of Medicine & Veterinary Medicine, University of Edinburgh, Edinburgh, UK
2. Centre for Clinical Brain Sciences, University of Edinburgh, Edinburgh, UK
3. Centre for Medical Informatics, Usher Institute of Population Health Sciences and Informatics, University of Edinburgh, Edinburgh, UK
4. Institute of Genetics and Molecular Medicine, University of Edinburgh, Edinburgh, UK
5. Centre for Cognitive Ageing and Cognitive Epidemiology, Edinburgh, UK
6. Edmond J. Safra Program in Parkinson's Disease and the Morton and Gloria Shulman Movement Disorders Clinic, Toronto Western Hospital, Toronto, Canada

*Corresponding author.

Email: tim.wilkinson@ed.ac.uk

¶ZH and TW are joint first authors

# Abstract

**Background:** Population-based, prospective studies can provide important insights into Parkinson's disease (PD) and other parkinsonian disorders. Participant follow-up in such studies is often achieved through linkage to routinely-collected healthcare datasets. We systematically reviewed the published literature on the accuracy of these datasets for this purpose.

**Methods:** We searched four electronic databases for published studies that compared PD and parkinsonism cases identified using routinely-collected data to a reference standard. We extracted study characteristics and two accuracy measures: positive predictive value (PPV) and/or sensitivity.

**Results:** We identified 18 articles, resulting in 27 measures of PPV and 14 of sensitivity. For PD, PPVs ranged from 56-90% in hospital datasets, 53-87% in prescription datasets, 81-90% in primary care datasets and was 67% in mortality datasets. Combining diagnostic and medication codes increased PPV. For parkinsonism, PPVs ranged from 36-88% in hospital datasets, 40-74% in prescription datasets, and was 94% in mortality datasets. Sensitivities ranged from 15-73% in single datasets for PD and 43-63% in single datasets for parkinsonism.

**Conclusions:** In many settings, routinely-collected datasets generate good PPVs and reasonable sensitivities for identifying PD and parkinsonism cases. Further research is warranted to investigate primary care and medication datasets, and to develop algorithms that balance a high PPV with acceptable sensitivity.

2

# Introduction

47

48    Despite well-established pathological features, the aetiologies of Parkinson's Disease (PD)

49    and other parkinsonian conditions remain poorly understood and disease-modifying

50    treatments have proved elusive. Large, prospective, population-based cohort studies with

51    biosample collections (e.g., UK Biobank, German National Cohort, US Precision Medicine

52    Initiative) provide a robust methodological framework with statistical power to investigate

53    the complex interplay between genetic, environmental and lifestyle factors in the aetiology

54    and natural history of neurological disorders such as PD and other parkinsonian disorders[1–

55    3].

56         Linkage to routinely-collected healthcare data – which are administrative datasets

57    collected primarily for healthcare purposes rather than to address specific research

58    questions[4] –provides an efficient means of long term follow-up in order to identify large

59    numbers of incident cases in such studies[1]. Furthermore, participant linkage to such

60    datasets can be used in randomised controlled trials as a cost-effective and comprehensive

61    method of follow-up for disease outcomes[5]. These data are coded using systems such as

62    the International Classification of Diseases (ICD)[6], the Systematized Nomenclature of

63    Medicine – Clinical Terms (SNOMED-CT) system[7], and the UK primary care Read system[8].

64         Before such datasets can be used to identify PD and parkinsonism cases in

65    prospective studies, their accuracy must be determined. Important measures are the

66    positive predictive value (PPV, the proportion of those coded positive that are true disease

67    cases) and sensitivity (the proportion of true disease cases that are coded positive).

68    Specificity and negative predictive value (NPV) are less relevant as specificity will be high

69    when precise diagnostic codes are used and NPV, which is related to disease prevalence, will

3

70    be high in population-based studies where most individuals do not develop the disease of

71    interest.

72         We systematically reviewed published studies evaluating the accuracy of routinely-

73    collected healthcare data for identifying PD and parkinsonism cases.

74

75

## 76    Methods

### 77    Study Protocol

78    We prospectively published the protocol for this systematic review

79    (www.crd.york.ac.uk/PROSPERO, number: 2016:CRD42016033715)[9].

80

### 81    Search Strategy and Eligibility Criteria

82    We searched the electronic databases MEDLINE (Ovid), EMBASE (Ovid), CENTRAL (Cochrane

83    Library) and Web of Science (Thomson Reuters) for articles published in any language

84    between 01.01.1990 and 23.06.2017 that compared codes for PD or parkinsonism from

85    routinely-collected healthcare data to a clinical expert-derived reference standard (see

86    Supplementary File S1 for search strategy). Studies had to provide either a PPV and/or a

87    sensitivity estimate, or sufficient raw data to calculate these. Where articles assessed more

88    than one dataset or evaluated both PPV and sensitivity, we included these as separate

89    studies. Hereafter we will refer to published papers as 'articles' and these separate analyses

90    as 'studies'. We chose the date limits based on our judgement that accuracy estimates from

91    studies published prior to 1990 would have limited current applicability. We also screened

92    bibliographies of included studies and relevant review papers to identify additional

93    publications. Studies had to have ≥10 coded cases, due to the limited precision of studies

94    below this size. Studies reporting sensitivity values had to be population-based (i.e.

95    community-based as opposed to hospital-based) with comprehensive attempts to detect all

96    disease cases. Where multiple studies investigated overlapping populations, we included

97    the study with the larger population size.

98

## Study Selection

100    Two authors (AS and SH) independently screened all titles and abstracts generated by the

101    search, and reviewed full text articles of all potentially eligible studies to determine if the

102    inclusion criteria were met. In the case of disagreement or uncertainty, we reached a

103    consensus through discussion and, where necessary, involvement of a senior third author

104    (CLMS).

105

## Data Extraction

107    Using a standardized form, two authors (TW and ZH) independently extracted the following

108    data from each study: first author, year of publication, time period during which coded data

109    were collected, country of study, study population, study size (defined as the total number

110    of code positive cases for PPV [true positives plus false positives] and the total number of

111    true positives for sensitivity [true positives and false negatives]), type of routine data used

112    (e.g., hospital admissions, mortality or primary care), coding system and version used,

113    specific codes used to identify cases, diagnostic coding position (e.g. primary or secondary

114    position), parkinsonian subtypes investigated, and the method used to make the reference

115    standard diagnosis.

116        We recorded the reported PPV and/or sensitivity estimates, as well as any

117        corresponding raw data. After discussion, any remaining queries were resolved with a senior

118        third author (CLMS). When necessary, we contacted study authors to request additional

119        information.

120

## Quality Assessment

122        We adapted the Quality Assessment of Diagnostic Accuracy Studies 2 (QUADAS-2)[10] tool

123        to evaluate the risk of bias in the estimates of accuracy and any concerns about the

124        applicability of each article to our specific research question (Supplementary Table S2). Two

125        authors (TW and ZH) independently assigned quality ratings, with any discrepancies

126        resolved through discussion. We performed this evaluation in the context of our specific

127        review question and not as an indication of the overall quality of the articles.

128

## Statistical Analysis/Data Synthesis

130        We tabulated the extracted data, and calculated 95% confidence intervals for the accuracy

131        measures from the raw data using the Clopper-Pearson (exact) method. Due to substantial

132        heterogeneity in study settings and methodologies, we did not perform a meta-analysis, as

133        we considered any summary estimate to be potentially misleading. Instead, we assessed the

134        full range of results in the context of study methodologies, populations and specific data

135        sources. We also reported any within-study comparisons in which a single variable was

136        changed to examine its effect on PPV or sensitivity. We performed analyses using the

137        statistical software StatsDirect3.

138

139

# Results

## Study Characteristics

142 18 published articles fulfilled our inclusion criteria[11–28]. A flow diagram of the study

143 selection process is shown in Fig 1. We obtained key additional information from the

144 authors of two studies[20,24]. Of the 18 included articles, 13 reported PPV[11,13–24], four

145 reported sensitivity[25–28] and one reported both[12]. Four articles contained more than

146 one study[11–13,17]. One of these consisted of multiple sub-studies, using different

147 methods to evaluate datasets across several countries, so we included these as six separate

148 studies[13]. In total, there were 27 measures of PPV and 14 of sensitivity. Study

149 characteristics are summarised in Tables 1 and 2 respectively.

150

151 **Fig 1: PRISMA Flow Diagram**

152

153

154

155

156

157

158

159

160

161

**Table 1. Characteristics of studies reporting positive predictive value, stratified by dataset type**

| First author & year of publication | Year of study | Country | Study population composition | Study size (n) | Routine dataset used | Coding system | Codes used to identify cases | Diagnostic coding position | Reference standard |
|---|---|---|---|---|---|---|---|---|---|
| Hospital-derived datasets: | | | | | | | | | |
| Butt 2014[10] | 1991-2011 | Canada | Population of Ontario ≥20yrs | Inpatient: 79 Outpatient: 435 | Hospital: inpatient Hospital: outpatient | ICD-9 (pre-2002) ICD-10 (post-2002) | Parkinsonism: ICD-9: 332.0, 332; ICD-10: G20, G21.0–0.4, G21.8–9, G22, F02.3 | Not specified | Medical record review |
| Feldman 2012[11] | 1964-2004 | Sweden | Twins across Sweden >50yrs | PD: 72 Parkinsonism: 75 | Hospital: inpatient | ICD-7 (1961-67) ICD-8 (1968-86) ICD-9 (1987-96) ICD-10 (1997-2009) | PD: ICD-7: 350; ICD-8: 342.00; ICD-9: 332.0; ICD-10: G20 Parkinsonism: ICD-8: 342.08, 342.09; ICD-9: 333.0; ICD-10: G21.4, G21.8, G21.9, G23.1, G23.2, G23.9, G25.9 | Any | Screening interview, medical record review and examination by physician |
| Gallo [a] 2015[12] | 1994-2010 | Sweden | Hospital patients, EPIC study participants | 62 | Hospital: unclear | ICD-9 (pre-1996) ICD-10 (post-1996) | PD: ICD-9: 332; ICD-10: G20, G21 | Not specified | Medical record review |
| Gallo [b] 2015[12] | 1991-2010 | Sweden | Hospital patients, EPIC study participants | 299 | Hospital: inpatient and outpatient | ICD-9 ICD-10 | PD: ICD-9: 332; ICD-10: G20 | Not specified | Medical record review |
| Kestenbaum 2015[13] | 2009-2014 | USA | Tertiary referral centre patients | 100 | Hospital: unclear | ICD-9 | PD: 332.0 | Not specified | Medical record review |
| Swarztrauber 2005[14] | 1998-2002 | USA | Veterans hospital patients | 175 | Hospital: inpatient and outpatient | ICD-9-CM | Parkinsonism: 332.0, 332.1, 333.0 | Not specified | Medical records review |
| Szumski 2009[15] | 2001-2004 | USA | Veterans hospital patients | 577 | Hospital: outpatient | ICD-9-CM | PD: 332.0 | Not specified | Medical record review |
| Wei 2016[16] | Unclear | USA | Hospital patients | 100 | Hospital: inpatient and outpatient | ICD-9 | PD: 332.0 | Not specified | Medical records review |
| Wermuth 2015[17] | 1996-2009 | Denmark | Neurological hospital patients | 2625 | Hospital: inpatient and outpatient | ICD-8 ICD-10 | PD: ICD-8: 342, ICD-10: G20 | Primary | Medical record review |
| White 2007[18] | 1998-2000 | USA | Veterans hospital patients | 782 | Hospital: inpatient and outpatient | ICD-9-CM | Parkinsonism: 332.0, 332.1 | Any | Medical record review |
| Primary care-derived datasets | | | | | | | | | |
| Hernán 2004[19] | 1995-2001 | UK | GP patients | 106+ | Primary care | Read code | Not specified (investigated PD) | Not applicable | Medical record review |
| Prescription-derived datasets | | | | | | | | | |
| Butt 2014[10] | 1991-2011 | Canada | Population of Ontario ≥65 | 395 | Prescriptions | Not specified | Parkinsonism: Levodopa; MAO-B inhibitors; dopamine agonists; COMT inhibitors | Not applicable | Medical record review |

| Meara 1999[20] | Not stated | UK | GP patients | PD: 402 Parkinsonism: 402 | Prescriptions (from primary care) | Not specified | PD: Not specified Parkinsonism: Not specified | Not applicable | History and examination by physician and medical record review |
|---|---|---|---|---|---|---|---|---|---|
| Wei 2016[16] | Unclear | USA | Hospital patients | 100 | Prescriptions | Not specified | PD: Rotigotine; Entacapone; Selegiline hydrochloride; Pergolide; Rasagiline | Not applicable | Medical record review |
| **Mortality datasets** | | | | | | | | | |
| Feldman. A 2012[11] | 1998-2007 | Sweden | Twins across Sweden >50yrs | PD: 18 Parkinsonism: 18 | Mortality | ICD-10 | PD: G20 Parkinsonism: G21.4, G21.8, G21.9, G23.1, G23.2, G23.9, G25.9 | Any | Screening interview, medical record review and examination by physician |
| **Combined datasets (accuracy measures for constituent datasets unable to be separated)** | | | | | | | | | |
| Bower 1999[21] | 1976-1990 | USA | Population of Olmsted county | 2472 | Synthesised medical information | H-ICDA | Parkinsonism: H-ICDA 53 diagnostic codes | Not specified | Medical record review |
| Gallo [c] 2015[12] | 1998-2010 | Spain | EPIC study participants | 39 | Prescriptions; Primary care; Mortality; Hospital: inpatient | ATC/DDD index ICD-9 | PD: ICD-9: 332, 332.0, 332.1; ATC/DDD index N04, N04A, N04B | Not specified | Medical record review |
| Gallo [d] 2015[12] | Unclear - 2010 | Spain | EPIC study participants | 41 | Primary care; Prescriptions; Mortality | ICPC ATC/DDD index ICD-9 | PD: ICPC N87; ATC/DDD index N04, N04A, N04B; ICD-9: 332.x | Not specified | Medical record review |
| Gallo [e] 2015[12] | 1998-2010 | Spain | EPIC study participants | 99 | Hospital: inpatient; Primary care; Prescriptions; Mortality | ICD-9 ICPC2 ATC/DDD index ICD-10 | PD: ICD-9: 332; ICPC2-WICC N87; ATC/DDD index N04x; ICD-10: G20 | Not specified | Medical record review |
| Gallo [f] 2015[12] | 1992-2008 | Italy | EPIC study participants | 81 | Hospital: inpatient; Mortality; Prescriptions | ICD-9 ICD-10 ATC/DDD index | PD: ICD-9 332; ATC/DDD index: N04, N04A, N04B; ICD-10 G20 | Not specified | Medical record review |
| Savica 2013[22] | 1991-2005 | USA | Population of Olmsted county | 4957 | Synthesised medical information | H-ICDA ICD-9 | Parkinsonism: H-ICDA 38 diagnostic codes, ICD-9: 331.9, 332.0, 332.1, 333.0, 333.1, 781.0, 781.3 | Not specified | Medical record review |
| Thacker 2016[23] | 2005-2015 | USA | Patients from a single medical institution | 129 | Hospital: inpatient and outpatient Primary care | ICD-9 | PD: 332, 332.0 | Primary | Medical records review |

163 Year of study: the time period during which coded data was collected. Study size: the total number of code positive cases (true positives plus
164 false positives). Where both PD and parkinsonism were investigated in one article, study sizes for both are displayed. Study population
165 composition: population cohort from which cases were identified.
166 ICD codes for Parkinson's disease - ICD-7 350; ICD-8 342.00; ICD-9(-CM) 332.0; ICD-10 G20.
167 ICD codes for other Parkinsonism - ICD-8: 342.08 (other defined Parkinsonism), 342.09 (unspecified Parkinsonism); ICD-9(-CM): ICD-9-CM:
168 332.1 (secondary Parkinson's disease), 333.0 (other degenerative diseases of the basal ganglia); ICD-10: G21.4 (vascular Parkinsonism), G21.8
169 (other defined secondary Parkinsonism), G21.9 (unspecified secondary Parkinsonism), G23.1 (progressive supranuclear ophthalmoplegia),
170 G23.2 (striatonigral degeneration), G23.9 (unspecified degenerative disease of basal ganglia), G25.9 (unspecified extrapyramidal and

171 movement disorder). Additional ICD codes – ICD-9: 331.9 (cerebral degeneration), 333.1 (essential and other specified forms of tremor), 781.0
172 (abnormal involuntary movements), 781.3 (lack of coordination).
173 [+] Exact study size unknown, reported as 7% of 1521 (could range from 99-115) – authors contacted, but data unavailable.
174 Abbreviations: PD - Parkinson's Disease; EPIC - European Prospective Investigation into Cancer and Nutrition study; ICD- International
175 Classification of Diseases; H-ICDA - Hospital Adaptation of ICDA; ATC/DDD index - Anatomical Therapeutic Chemical Classification System with
176 Defined Daily Doses; ICPC - International Classification of Primary Care.
177

178

179

180

181

182

183

184

185

186

187

188
189

**190    Table 2. Characteristics of studies reporting sensitivity, stratified by dataset type**

| First author, year of publication | Year of study | Country | Study population composition | Study size (n) | Routine dataset used | Coding system | Codes used to identify cases | Diagnostic coding position | Reference standard |
|---|---|---|---|---|---|---|---|---|---|
| Mortality certificate-derived datasets: | | | | | | | | | |
| Benito-León 2014[24] | 1994-2007 | Spain | Three communities near Madrid | 82 | Mortality | ICD-9 (pre 1999) ICD-10 (post 1999) | Not specified (investigated PD) | Primary | Screening (in-person, telephone and mail questionnaire) and neurological examination |
| Beyer 2001[25] | 1993-1996 | Norway | County (Rogaland) | 84 | Mortality | ICD-9 or ICPC | Not specified (investigated PD) | Primary + Any | Semi-structured interview and a clinical examination |
| Fall 2003[26] | 1989-1998 | Sweden | Central district of Östergötland | 121 | Mortality | ICD-9 | Not specified (investigated PD) | Primary + Any | Examination and medical record review |
| Feldman 2012[11] | 1998-2008 | Sweden | Twins across Sweden >50yrs | Parkinsonism: 127 PD: 77 | Mortality | ICD-10 | PD: G20 Parkinsonism: G21.4, G21.8, G21.9, G23.1, G23.2, G23.9, G25.9 | Any | Screening interview, medical record review and examination |
| Williams-Gray 2013[27] | 2000-2012 | UK | County (Cambridgeshire) | 63 | Mortality | Not specified | Not specified (investigated PD) | Primary + Any | History and neurological examination |
| Hospital-derived datasets: | | | | | | | | | |
| Feldman 2012[11] | 1964-2009 | Sweden | Twins across Sweden >50yrs | Parkinsonism: 194 PD: 132 | Hospital: inpatient | ICD-7 (1961-67) ICD-8 (1968-86) ICD-9 (1987-96) ICD-10 (1997-2009) | PD: ICD-7: 350; ICD-8: 342.00; ICD-9: 332.0; ICD-10: G20 Parkinsonism: ICD-8: 342.08, 342.09; ICD-9: 333.0; ICD-10: G21.4, G21.8, G21.9, G23.1, G23.2, G23.9, G25.9 | Any | Screening interview, medical record review and examination |

191    Year of study: the time period during which coded data was collected. Study size: the total number of true positive according to the reference
192    standard (true positives and false negatives). Where both PD and parkinsonism were investigated in one article, study sizes for both are
193    displayed. Study population composition: population cohort from which cases were identified.
194    ICD codes for Parkinson's disease - ICD-7 350; ICD-8 342.00; ICD-9 332.0; ICD-10 G20.
195    ICD codes for other Parkinsonism - ICD-8: 342.08 (other defined Parkinsonism), 342.09 (unspecified Parkinsonism); ICD-9: 333.0 (other
196    degenerative diseases of the basal ganglia); ICD-10: G21.4 (vascular Parkinsonism), G21.8 (other defined secondary Parkinsonism), G21.9
197    (unspecified secondary Parkinsonism), G23.1 (progressive supranuclear ophthalmoplegia), G23.2 (striatonigral degeneration), G23.9
198    (unspecified degenerative disease of basal ganglia), G25.9 (unspecified extrapyramidal and movement disorder)
199

200         Study size varied considerably, ranging from 39-4957. All 18 articles were based in

201   high-income countries. Three were from the UK[20,21,28], six from mainland

202   Europe[12,13,18,25–27], eight from the USA[14–17,19,22–24], and one from Canada[11].

203   There were 12 PPV estimates and two sensitivity estimates from hospital data[11–19], two

204   PPV and 10 sensitivity estimates from mortality data[12,25–28], two PPV estimates from

205   primary care data[20], four PPV estimates from prescription data[11,17,21] and seven PPV

206   estimates and two sensitivity estimates from combining datasets from different

207   sources[12,13,22–24]. There were no sensitivity estimates from primary care or prescription

208   data.

209         PD was evaluated in 13 articles, with eight estimating PPV[13,14,16–18,20,21,24],

210   four estimating sensitivity[25–28] and one estimating both[12]. Parkinsonism was evaluated

211   by seven articles, of which six estimated PPV[11,15,19,21–23] and one assessed both PPV

212   and sensitivity[12]. All of the parkinsonism articles combined PD with other causes of

213   parkinsonism.

214         The methods of reference standard used could be broadly divided into two

215   categories: patient history and examination (majority of studies reporting sensitivity) and

216   medical record review (majority of studies reporting PPV). In addition, where entire

217   populations were under study, some studies incorporated a screening method (e.g.,

218   telephone interview) to identify potential cases[12,25].

219         Where reported, codes used to identify PD cases were consistent and appropriate to

220   the ICD version used. However, the range of codes used to identify other parkinsonian

221   conditions varied considerably, reflecting the broad range of pathologies that can lead to

222   parkinsonism. Seven studies did not specify the exact codes used[17,20,21,25–28]. ICD

223   versions used reflected the time period over which the studies were conducted. 19 studies

12

224    used ICD-9 (or ICD-9-CM, a clinically modified version used in the USA, and identical to ICD-9

225    with respect to parkinsonian diagnoses)[11–17,19,23–27], 11 used ICD-10[11–13,18,25],

226    three used ICD-8[12,18], and two used ICD-7[12]. One of the primary care studies used

227    Read-coded data[20]. Four studies, including the three that evaluated prescription data, did

228    not specify the coding system used[11,17,21,28].

229        The diagnostic coding position assessed also varied. Three studies assessed primary

230    diagnoses alone[18,24,25], eight used any diagnostic position[12,19,26–28], while 13 did

231    not specify the coding position[11,13–17,22,23]. Diagnostic position was not applicable in

232    the studies of primary care and prescription data due to the nature of these

233    datasets[11,17,20,21].

234

## Quality Assessment

236    Only two articles were judged to be of low risk of bias or applicability concerns in the

237    QUADAS-2 assessment[11,12] (Supplementary Table S3). The commonest concerns were:

238    selection bias, lack of reporting of the codes used to identify disease cases, insufficiently

239    rigorous reference standards, inappropriate inclusions and exclusions, or patients being lost

240    to follow-up.

241

## Positive predictive value

243    For PD, there were 17 PPV estimates in total (Fig 2)[12–14,16–18,20,21,24]. These

244    comprised seven PPV estimates of hospital data alone[12–14,16–18], one of mortality data

245    alone[12], two for prescription data alone[17,21], one of primary care data alone[20], one

246    of prescription data and primary care data in combination[20], and five of datasets used in

247    combination[13,24]. PPVs ranged from 36-90% across all studies. Nine of the 17 estimates

13

248    were >75%. The single study of Read coding in primary care data alone reported a PPV of

249    81%, increasing to 90% with the presence of a relevant medication code in addition to a

250    diagnostic code[20]. The two studies of medication data alone reported PPVs of 53% and

251    87%[17,21]. The single, small study of mortality data had a PPV of 67%[12].

252

253

254    **Fig 2: Positive predictive values (PPVs) of coded diagnoses**

255    Study size: total number of code-positive cases (true positives + false positives). *Exact

256    sample size unknown, most conservative estimate used. Box sizes reflect Mantel-Haenszel

257    weight of study (inverse variance, fixed effects).

258

259

260         Several within-study comparisons were available from three studies identifying PD

261    (Table 3)[12,16,17]. Two of these investigated the change in PPV for hospital data to identify

262    PD when algorithms containing additional criteria were used[12,16]. Both showed a

263    moderate increase in PPV if a relevant diagnosis code was recorded more than once, or if a

264    specialist department assigned such a code. One study reported an increase in PPV when

265    only primary position diagnoses were assessed[12]. Another showed that incorporating

266    selected medication codes with diagnosis codes increased the PPV from 76% to 86%,

267    although this was at the expense of reduced case ascertainment[16]. Finally, one study

268    showed that the combination of a diagnostic code in hospital data with a relevant

269    medication code increased the PPV when compared to using either dataset alone (94%

270    versus 87% and 89% respectively)[17].

14

271        For parkinsonism there were 10 PPV estimates in total (Fig 2)[11,12,15,19,21–23].

272    These comprised five estimates from hospital data alone[11,12,15,19], two from

273    prescription data alone[11,21], one from mortality data alone[12], and two from using

274    datasets in combination[22,23]. PPVs ranged from 40-94% in the single datasets and from

275    22-28% in the combination datasets. The two studies of parkinsonism in prescription data

276    produced very different PPV estimates of 40% and 74%[11,21]. One of these studies

277    reported that the PPV of medication data to identify any parkinsonian disorder was

278    considerably higher than that for PD (74% and 53% respectively)[21].

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

| Criteria applied: | PPV (95% CI) | Number of cases identified |
|---|---|---|
| **Parkinson's Disease** | | |
| a) **Feldman 2012** (hospital inpatient data) | | |
| Parkinson's disease ICD code only | 71 (59-81) | 72 |
| Exclusion of patients with other (non-Parkinson's disease) parkinsonian codes | 70 (58-81) | 67 |
| Code frequency ≥2 hospital admissions | 76 (61-88) | 42 |
| Code in primary diagnostic position | 83 (70-92) | 53 |
| Code assigned in specialist department (neurological/neurosurgical/geriatric) | 83(63-95) | 24 |
| b) **Szumski 2009** (hospital outpatient data) | | |
| Parkinson's disease ICD codes only | 76 (72-79) | 579 |
| Code frequency ≥2 at any clinic | 79(76-83) | 409 |
| Code assigned in any neurology clinic | 79 (75-83) | 352 |
| Code assigned in movement disorder speciality clinic | 87 (81-92) | 177 |
| Code + prescribed antiparkinsonian medication | 86 (82-89) | 408 |
| c) **Wei 2016** | | |
| Parkinson's disease ICD codes only | 89 (81-94) | 100 |
| Prescription only | 87 (78-93) | 100 |
| ICD code and prescription | 94* | Unknown* |
| **Parkinsonism** | | |
| d) **Butt 2014** [†] | | |
| Hospital inpatient ICD code ever | 87 (79-96) | 63 |
| Hospital outpatient ICD code ever | 55 (49-60) | 297 |
| Prescription ever | 40 (35-44) | 395 |
| Outpatient code frequency ≥2 in one year | 83 (77-89) | 169 |
| Outpatient code frequency ≥2 in one year by a specialist | 87 (81–92) | 134 |
| Outpatient code AND Prescription | 85 (79-90) | 174 |
| Prescription AND outpatient code within +/- 6 months | 87 (82-92) | 166 |

295     **Table 3: Within-study analyses: algorithm development**

296     The effect of additional criteria to identify PD cases on PPV and the number of cases

297     identified. * Sample size and confidence intervals unknown for this accuracy measure.

298

## Sensitivity

300  For PD, there were 11 sensitivity estimates in total (Fig 3)[12,25–28]. Of these, nine were

301  sensitivity estimates for mortality data alone, consistently showing that codes in the primary

302  position only gave low sensitivities of 11-23%, rising to 53-60% when codes from any

303  position were included[12,25–28]. A single study reported the sensitivity of hospital data to

304  be 73%, increasing to 83% when hospital and mortality data were combined. There were no

305  sensitivity estimates for primary care or prescription data.

306       For parkinsonism, there were three sensitivity estimates, all from one study[12].

307  Hospital admissions and mortality data combined gave higher sensitivity (71%) compared

308  with either mortality or hospital data alone (43% and 63% respectively).

309

310

311  **Fig 3: Sensitivity estimates of coded diagnoses**

312  Study size: total number of true positives according to reference standard (true positives +

313  false negatives). *Unknown sample size and confidence intervals. Box sizes reflect Mantel-

314  Haenszel weight of study (inverse variance, fixed effects).

315

316

317

318

319

320

17

## Discussion

We have demonstrated that existing validation studies show a wide variation in the accuracy of routinely-collected healthcare data for the identification of PD and parkinsonism cases. Despite this, in the right setting, achieving high PPVs is possible. Sensitivity is generally lower than PPV, but is increased by combining data sources.

False positives (participants who receive a disease code but do not have the disorder) may arise in routinely-collected coded datasets for several reasons. Firstly, the clinician may incorrectly diagnose the condition. Given that PD and other parkinsonian disorders are largely clinical diagnoses made without a definitive diagnostic test, there is the potential for diagnostic inaccuracies. Clinicopathological studies have shown discrepancies between clinical diagnoses in life and neuropathological confirmation[29] and there is evidence that accuracy increases when diagnoses are made by movement disorder specialists[30–32]. Secondly, diagnoses may be incorrectly recorded in medical records, or errors may arise during the coding process. Similarly, false negatives (patients who have the condition but do not receive a code) may arise due to under-diagnosis, omission of the diagnosis from the medical records (e.g., because the condition is not the primary reason for hospital admission), or errors during the coding process.

The pharmacological treatment of PD is largely focussed on improving motor function and patients are treated with a limited number of drugs. This has allowed antiparkinsonian drugs to be used as 'tracers' in epidemiological studies[33,34]. There are potential problems with using prescription data as a proxy for PD diagnosis. This approach may disproportionately under-identify patients with early stage disease who do not yet require treatment. Also, a response to a trial of dopaminergic drugs may be used as part of

344    the diagnostic assessment in potential PD cases, meaning some patients prescribed

345    antiparkinsonian medications will not be subsequently diagnosed with PD. Furthermore,

346    antiparkinsonian can be prescribed for indications other than PD (such as dopamine

347    agonists for restless legs syndrome, endocrine disorders and other forms of parkinsonism).

348    The specific drugs licensed for use in parkinsonian conditions varies between countries and

349    may change over time. Therefore, an algorithm incorporating prescription data would need

350    to be continually revised to match prescribing patterns. Results from our review suggest

351    that prescription data alone has a low PPV for PD case ascertainment[21]; however, when

352    drug codes are combined with diagnostic codes, PPV increases but with reduced case

353    ascertainment[16,20]. Furthermore, prescription datasets appear to have a higher PPV

354    when identifying any parkinsonian disorder rather than specifically PD[21].

355

356    This study has several strengths and limitations. Our review benefits from prospective

357    protocol publication, comprehensive search criteria, and independent duplication of each

358    stage by two authors. Despite this, relevant studies may still have been missed, especially if

359    a validation study was a subsection of a paper with a wider aim. As all eligible studies were

360    included, the results may have been influenced by studies of lower quality. Only two articles

361    were found to be at low risk of bias or applicability concerns[11,12], and it is likely that

362    biases in study design would have affected the results. For example, one study with the

363    lowest PPV[23] used very broad ICD-9 codes such as 781.0 (abnormal involuntary

364    movements) and 781.3 (lack of coordination).

365        Since there is no method of diagnosing PD with certainty in life, there is likely to be

366    some misclassification of the reference standards used in the studies. The application of

367    stringent diagnostic criteria to reference standard diagnoses, although often necessary for

368    research purposes, may lead to some patients being misclassified as 'false positives' when

369    they do in fact have the condition. This may lead to underestimation of the PPV in some of

370    the studies. When considering the ideal reference standard for validation studies, there is a

371    trade-off between the robustness of the reference standard and validating sufficient cases

372    to produce precise accuracy estimates. For example, in-person neurological examination

373    may have greater diagnostic certainty than medical record review but this becomes difficult

374    as the cohort size increases.

375        Many of the studies reported cases with insufficient information to meet the

376    reference standard and the handling of these varied. Some studies excluded such cases,

377    others classified them as false positives, while some did not specify how they handled such

378    missing data. Excluding such cases may introduce selection bias, whereas counting them as

379    false positives may underestimate PPV.

380        The effect of possible publication bias on the results is difficult to estimate, but

381    disproportionate publication of studies which report more favourable accuracy measures

382    may lead to over-estimation of the performance of the codes. In addition, estimates of PPV

383    are dependent upon the prevalence of the condition in the study population but it was not

384    possible to assess the prevalence of PD within each study population.

385

386    Our review highlights several areas requiring further research. Given that the management

387    of PD is largely delivered in outpatients or the community, primary care data may be an

388    effective method of identifying cases. Whilst studies have suggested that PD diagnoses

389    made in primary care are less accurate than those made in a specialist setting[35,36],

390    primary care records combine notes made by primary care clinicians with prescription

391    records and correspondence from secondary care. Codes from primary care should

392    therefore include diagnoses made by specialists, thus increasing their accuracy. We found

393    only one small study of primary care data, reporting a promising PPV of 81%, improving to

394    90% with the inclusion of medication codes[20]. No studies investigated the sensitivity of

395    primary care data. Further research into the accuracy of primary care data is needed.

396         Two studies investigated using algorithmic combinations of codes from different

397    sources to improve PPV[12,16]. These investigated the additional benefit of the inclusion of

398    factors such as only including codes that appeared more than once, selecting codes in the

399    primary position only, combining diagnostic codes with prescription data, and only including

400    diagnoses made in specialist clinics. These methods increased PPV but at a cost to the

401    number of cases identified.  The development of algorithms that maximize PPV whilst

402    maintaining a reasonable sensitivity (e.g., by combining multiple complimentary datasets)

403    merits further evaluation.

404         To our knowledge, no studies have evaluated the accuracy of routinely-collected

405    healthcare data for solely identifying atypical parkinsonian syndromes such as PSP and MSA.

406    Further work is needed to understand whether these datasets provide a valuable resource

407    for studying these less common diseases.

408

409    In conclusion, our review summarises existing knowledge of the accuracy of routinely-

410    collected healthcare data for identifying PD and parkinsonism, and highlights approaches to

411    increase accuracy and areas where further research is required. Given the wide range of

412    results observed, prospective cohorts may wish to perform their own validation studies

413    based on their specific setting and research question.

414

415

## Acknowledgements

## References

1. Sudlow C, Gallacher J, Allen N, Beral V, Burton P, Danesh J, et al. UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. PLoS Med. 2015;12: e1001779. doi:10.1371/journal.pmed.1001779

2. German National Cohort (GNC) Consortium. The German National Cohort: aims, study design and organization. Eur J Epidemiol. 2014;29: 371–382. doi:10.1007/s10654-014-9890-7

3. Jaffe S. Planning for US Precision Medicine Initiative underway. The Lancet. 2015;385: 2448–2449. doi:10.1016/S0140-6736(15)61124-2

4. Benchimol EI, Smeeth L, Guttmann A, Harron K, Moher D, Petersen I, et al. The REporting of studies Conducted using Observational Routinely-collected health Data (RECORD) Statement. PLOS Med. 2015;12: e1001885. doi:10.1371/journal.pmed.1001885

5. Mc Cord KA, Al-Shahi Salman R, Treweek S, Gardner H, Strech D, Whiteley W, et al. Routinely collected data for randomized trials: promises, barriers, and implications. Trials. 2018;19: 29. doi:10.1186/s13063-017-2394-5

6. World Health Organization. The ICD-10 classification of mental and behavioural disorders: clinical descriptions and diagnostic guidelines. Geneva: World Health Organization; 1992.

7. SNOMED International. SNOMED CT [Internet]. [cited 18 May 2017]. Available: http://www.snomed.org/snomed-ct

8. NHS Digital. Read Codes [Internet]. [cited 17 May 2017]. Available: https://digital.nhs.uk/article/1104/Read-Codes

9. Stevenson A, Wilkinson T, Sudlow CLM, Ly A. The accuracy of electronic health datasets in identifying Parkinson's disease cases: a systematic review [Internet]. 28 Jan 2016

447     [cited 17 May 2017]. Available:
448     http://www.crd.york.ac.uk/PROSPERO/display_record.asp?ID=CRD42016033715

449   10. Whiting PF, Rutjes AWS, Westwood ME, Mallett S, Deeks JJ, Reitsma JB, et al. QUADAS-
450        2: a revised tool for the quality assessment of diagnostic accuracy studies. Ann Intern
451        Med. 2011;155: 529–536. doi:10.7326/0003-4819-155-8-201110180-00009

452   11. Butt DA, Tu K, Young J, Green D, Wang M, Ivers N, et al. A validation study of
453        administrative data algorithms to identify patients with Parkinsonism with prevalence
454        and incidence trends. Neuroepidemiology. 2014;43: 28–37. doi:10.1159/000365590

455   12. Feldman AL, Johansson ALV, Gatz M, Flensburg M, Petzinger GM, Widner H, et al.
456        Accuracy and sensitivity of Parkinsonian disorder diagnoses in two Swedish national
457        health registers. Neuroepidemiology. 2012;38: 186–193. doi:10.1159/000336356

458   13. Gallo V, Brayne C, Forsgren L, Barker RA, Petersson J, Hansson O, et al. Parkinson's
459        Disease Case Ascertainment in the EPIC Cohort: The NeuroEPIC4PD Study.
460        Neurodegener Dis. 2015;15: 331–338. doi:10.1159/000381857

461   14. Kestenbaum M, Ford B, Louis ED. Estimating the Proportion of Essential Tremor and
462        Parkinson's Disease Patients Undergoing Deep Brain Stimulation Surgery: Five-Year Data
463        From Columbia University Medical Center (2009–2014). Mov Disord Clin Pract. 2015;2:
464        384–387. doi:10.1002/mdc3.12185

465   15. Swarztrauber K, Anau J, Peters D. Identifying and distinguishing cases of parkinsonism
466        and Parkinson's disease using ICD-9 CM codes and pharmacy data. Mov Disord. 2005;20:
467        964–970. doi:10.1002/mds.20479

468   16. Szumski NR, Cheng EM. Optimizing algorithms to identify Parkinson's disease cases
469        within an administrative database. Mov Disord. 2009;24: 51–56.
470        doi:10.1002/mds.22283

471   17. Wei W-Q, Teixeira PL, Mo H, Cronin RM, Warner JL, Denny JC. Combining billing codes,
472        clinical notes, and medications from electronic health records provides superior
473        phenotyping performance. J Am Med Inform Assoc. 2016;23: e20-27.
474        doi:10.1093/jamia/ocv130

475   18. Wermuth L, Cui X, Greene N, Schernhammer E, Ritz B. Medical Record Review to
476        Differentiate between Idiopathic Parkinson's Disease and Parkinsonism: A Danish
477        Record Linkage Study with 10 Years of Follow-Up. Parkinsons Dis. 2015;2015: 781479.
478        doi:10.1155/2015/781479

479   19. White D, Moore S, Waring S, Cook K, Lai E. Identifying incident cases of parkinsonism
480        among veterans using a tertiary medical center. Mov Disord. 2007;22: 915–923.
481        doi:10.1002/mds.21353

482   20. Hernán MA, Logroscino G, Rodríguez LAG. A prospective study of alcoholism and the risk
483        of Parkinson's disease. J Neurol. 2004;251 Suppl 7: vII14-17. doi:10.1007/s00415-004-
484        1705-4

21. Meara J, Bhowmick BK, Hobson P. Accuracy of diagnosis in patients with presumed Parkinson's disease. Age Ageing. 1999;28: 99–102.

22. Bower JH, Maraganore DM, McDonnell SK, Rocca WA. Incidence and distribution of parkinsonism in Olmsted County, Minnesota, 1976-1990. Neurology. 1999;52: 1214–1220.

23. Savica R, Grossardt BR, Bower JH, Ahlskog JE, Rocca WA. Incidence and pathology of synucleinopathies and tauopathies related to parkinsonism. JAMA Neurol. 2013;70: 859–866. doi:10.1001/jamaneurol.2013.114

24. Thacker T, Wegele AR, Pirio Richardson S. Utility of electronic medical record for recruitment in clinical research: from rare to common disease. Mov Disord Clin Pract. 2016;3: 507–509. doi:10.1002/mdc3.12318

25. Benito-León J, Louis ED, Villarejo-Galende A, Romero JP, Bermejo-Pareja F. Under-reporting of Parkinson's disease on death certificates: a population-based study (NEDICES). J Neurol Sci. 2014;347: 188–192. doi:10.1016/j.jns.2014.08.048

26. Beyer MK, Herlofson K, Arsland D, Larsen JP. Causes of death in a community-based study of Parkinson's disease. Acta Neurol Scand. 2001;103: 7–11.

27. Fall P-A, Saleh A, Fredrickson M, Olsson J-E, Granérus A-K. Survival time, mortality, and cause of death in elderly patients with Parkinson's disease: a 9-year follow-up. Mov Disord. 2003;18: 1312–1316. doi:10.1002/mds.10537

28. Williams-Gray CH, Mason SL, Evans JR, Foltynie T, Brayne C, Robbins TW, et al. The CamPaIGN study of Parkinson's disease: 10-year outlook in an incident population-based cohort. J Neurol Neurosurg Psychiatr. 2013;84: 1258–1264. doi:10.1136/jnnp-2013-305277

29. Adler CH, Beach TG, Hentz JG, Shill HA, Caviness JN, Driver-Dunckley E, et al. Low clinical diagnostic accuracy of early vs advanced Parkinson disease: clinicopathologic study. Neurology. 2014;83: 406–412. doi:10.1212/WNL.0000000000000641

30. Hughes AJ, Daniel SE, Kilford L, Lees AJ. Accuracy of clinical diagnosis of idiopathic Parkinson's disease: a clinico-pathological study of 100 cases. J Neurol Neurosurg Psychiatr. 1992;55: 181–184.

31. Hughes AJ, Daniel SE, Lees AJ. Improved accuracy of clinical diagnosis of Lewy body Parkinson's disease. Neurology. 2001;57: 1497–1499.

32. Hughes AJ, Daniel SE, Ben-Shlomo Y, Lees AJ. The accuracy of diagnosis of parkinsonian syndromes in a specialist movement disorder service. Brain. 2002;125: 861–870.

33. Brandt-Christensen M, Kvist K, Nilsson FM, Andersen PK, Kessing LV. Use of antiparkinsonian drugs in Denmark: results from a nationwide pharmacoepidemiological study. Mov Disord. 2006;21: 1221–1225. doi:10.1002/mds.20907
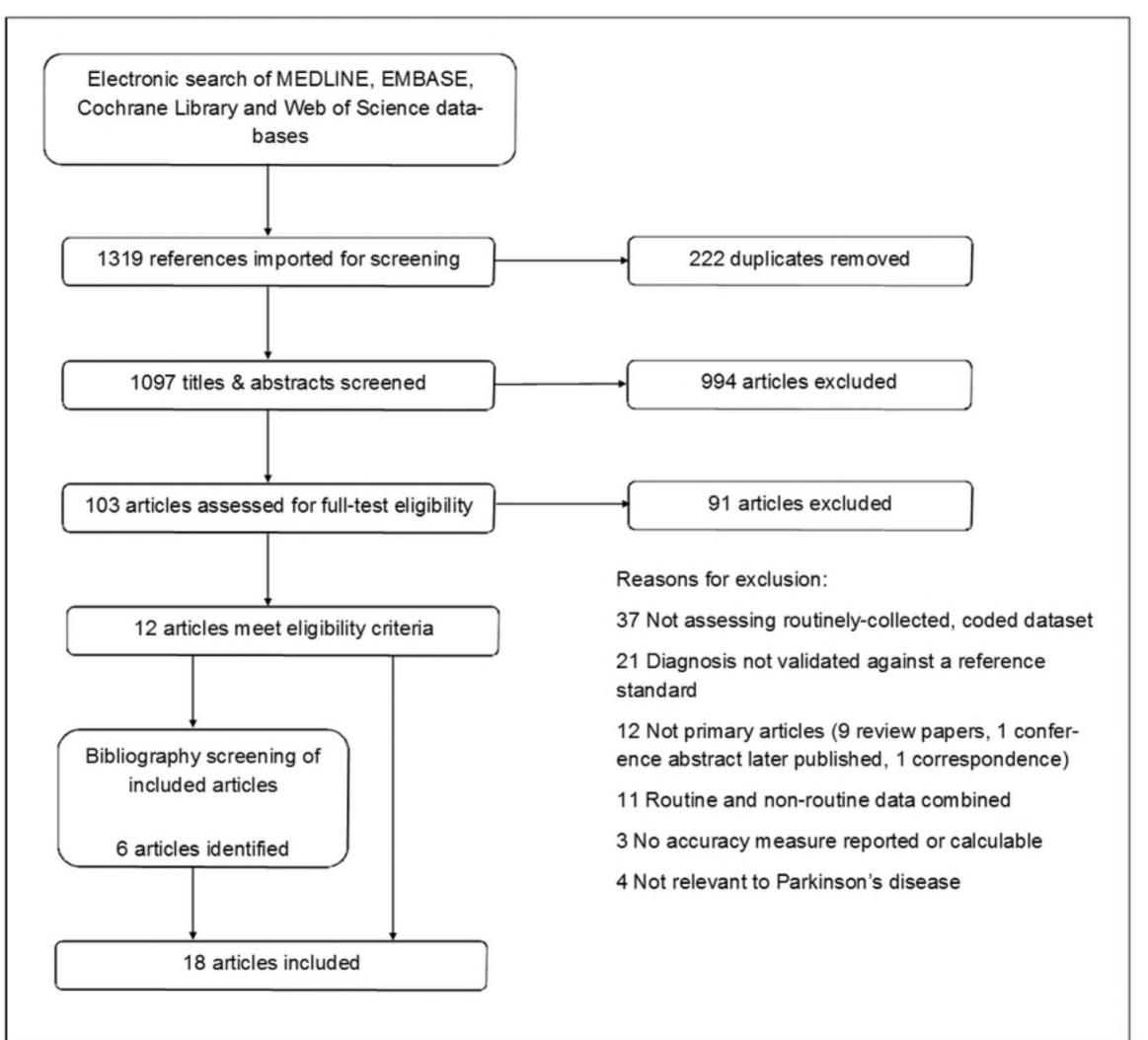
521   34. Chiò A, Magnani C, Schiffer D. Prevalence of Parkinson's disease in Northwestern Italy:
522        comparison of tracer methodology and clinical ascertainment of cases. Mov Disord.
523        1998;13: 400–405. doi:10.1002/mds.870130305

524   35. Newman EJ, Breen K, Patterson J, Hadley DM, Grosset KA, Grosset DG. Accuracy of
525        Parkinson's disease diagnosis in 610 general practice patients in the West of Scotland.
526        Mov Disord. 2009;24: 2379–2385. doi:10.1002/mds.22829

527   36. Schrag A, Ben-Shlomo Y, Quinn N. How valid is the clinical diagnosis of Parkinson's
528        disease in the community? J Neurol Neurosurg Psychiatr. 2002;73: 529–534.

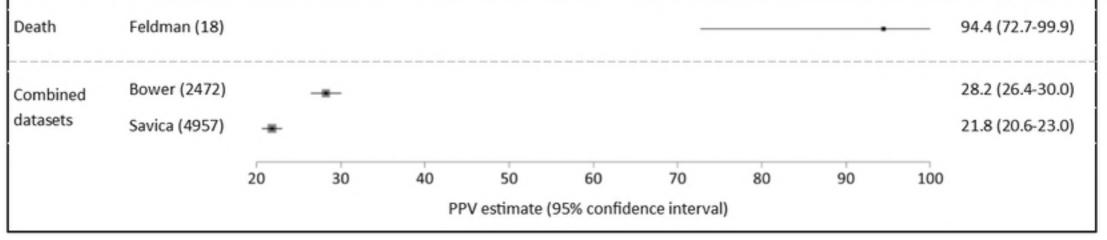529

## 530  **Supporting information**

531  **S1 File. Search strategy.**

532

533  **S2 File. QUADAS-2 assessment.**

534

535  **S3 Table. QUADAS-2 summary results.**

536

537  **S4 Checklist. PRISMA checklist.**

Electronic search of MEDLINE, EMBASE, Cochrane Library and Web of Science databases

↓

1319 references imported for screening → 222 duplicates removed

↓

1097 titles & abstracts screened → 994 articles excluded

↓

103 articles assessed for full-test eligibility → 91 articles excluded

↓

12 articles meet eligibility criteria

Reasons for exclusion:

37 Not assessing routinely-collected, coded dataset

21 Diagnosis not validated against a reference standard

12 Not primary articles (9 review papers, 1 conference abstract later published, 1 correspondence)

11 Routine and non-routine data combined

3 No accuracy measure reported or calculable

4 Not relevant to Parkinson's disease

↓

Bibliography screening of included articles

6 articles identified

↓

18 articles included

| Dataset type | Author (study size) | | PPV (95%CI) |
|---|---|---|---|

**Parkinson's disease**

| | | | |
|---|---|---|---|
| Hospital | Feldman (72) | | 70.8 (58.9-81.0) |
| | Gallo [a] (62) | | 90.3 (80.1-96.4) |
| | Gallo [b] (299) | | 55.5 (49.7-61.2) |
| | Kestenbaum (100) | | 88.0 (80.0-93.6) |
| | Szumski (579) | | 75.6 (71.6-78.8) |
| | Wei (100) | | 89.0 (81.2-94.4) |
| | Wermuth (2625) | | 78.8 (77.2-80.3) |
| Primary care | Hernan [diagnoses] (99*) | | 81.0 (71.7-88.0) |
| | Hernan [diagnoses] | | 90.0 (82.2-95.0) |
| Prescription | Meara (402) | | 53.0 (48.0-58.0) |
| | Wei (100) | | 87.0 (78.8-92.9) |
| Death | Feldman (18) | | 66.7 (41.0-86.7) |
| Combined datasets | Gallo [c] (39) | | 82.1 (66.5-92.5) |
| | Gallo [d] (41) | | 51.2 (35.1-67.1) |
| | Gallo [e] (99) | | 53.5 (43.2-63.6) |
| | Gallo [f] (81) | | 35.8 (25.4-47.2) |
| | Thacker (129) | | 55.0 (46.0-63.8) |

**Any Parkinsonian disorder**

| | | | |
|---|---|---|---|
| Hospital | Butt [inpatient] (79) | | 75.9 (65.0-84.9) |
| | Butt [outpatient] (435) | | 43.9 (39.2-48.7) |
| | Feldman (75) | | 88.0 (78.4-94.4) |
| | Swarztrauber (75) | | 88.0 (82.2-92.4) |
| | White (782) | | 76.0 (72.8-78.9) |
| Prescription | Butt (395) | | 39.5 (34.6-44.5) |
| | Meara (402) | | 74.0 (69.8-78.6) |
| Death | Feldman (18) | | 94.4 (72.7-99.9) |
| Combined datasets | Bower (2472) | | 28.2 (26.4-30.0) |
| | Savica (4957) | | 21.8 (20.6-23.0) |

PPV estimate (95% confidence interval)

| Data type (author (study size)) | Sensitivity (95%CI) |
|---|---|
| **Parkinson's disease** | |
| Death [primary position] | |
| Benito-Leon (82) | 14.6 (7.8-24.2) |
| Beyer (84) | 22.6 (14.2-33.0) |
| Fall (121) | 10.7 (5.9-17.7) |
| Feldman (77) | 19.5 (11.3-30.1) |
| Williams-Gray (63) | 20.0 (10.3-30.9) |
| Death [any position] | |
| Beyer (84) | 56.0 (44.7-66.8) |
| Fall (121) | 52.9 (43.6-62.0) |
| Feldman (77) | 57.1 (45.4-68.4) |
| Williams-Gray (63) | 60.0 (47.0-72.4) |
| Hospital Feldman (132) | 72.7 (64.3-80.1) |
| Combined hospital and death Feldman* | 83.1 |
| **Any Parkinsonian disorder** | |
| Death [any position] Feldman (127) | 43.3 (34.5-52.4) |
| Hospital Feldman (194) | 63.4 (56.2-70.2) |
| Combined hospital and death Feldman* | 70.9 |



Sensitivity estimate (95% confidence interval)