

# Intra-database Validation of Case-identifying Algorithms Using Reconstituted Electronic Health Records From Healthcare Claims Data

**Nicolas Henri Thurin** (✉ [nicolas.thurin@u-bordeaux.fr](mailto:nicolas.thurin@u-bordeaux.fr))

Univ. Bordeaux, INSERM CIC-P1401, Bordeaux PharmacoEpi, Bordeaux, France <https://orcid.org/0000-0003-3589-0819>

**Pauline Bosco-Levy**

Univ. Bordeaux, INSERM CIC-P1401, Bordeaux PharmacoEpi

**Patrick Blin**

Univ. Bordeaux, INSERM CIC-P1401, Bordeaux PharmacoEpi

**Magali Rouyer**

Univ. Bordeaux, INSERM CIC-P1401, Bordeaux PharmacoEpi

**Jérémy Jové**

Univ. Bordeaux, INSERM CIC-P1401, Bordeaux PharmacoEpi

**Stéphanie Lamarque**

Univ. Bordeaux, INSERM CIC-P1401, Bordeaux PharmacoEpi

**Lignot Séverine**

Univ. Bordeaux, INSERM CIC-P1401, Bordeaux PharmacoEpi

**Régis Lassalle**

Univ. Bordeaux, INSERM CIC-P1401, Bordeaux PharmacoEpi

**Abdelilah Abouelfath**

Univ. Bordeaux, INSERM CIC-P1401, Bordeaux PharmacoEpi

**Bignon Emmanuelle**

Univ. Bordeaux, INSERM CIC-P1401, Bordeaux PharmacoEpi

**Diez Pauline**

Univ. Bordeaux, INSERM CIC-P1401, Bordeaux PharmacoEpi

**Marine Gross-Goupil**

Département of medical Oncology, Hôpital saint-André, CHU de Bordeaux

**Michel Soulié**

Department of Urology, University Hospital of Rangueil, CHU de Toulouse

**Mathieu Roumigué**

Department of Urology, University Hospital of Rangueil, CHU de Toulouse

**Sylvestre Le Moulec**

Department of oncology, Clinique Marzet

**Marc Debouverie**

Department of neurology, CHRU de Nancy, Université de Lorraine, EA 4360 APEMAC

**Bruno Brochet**

CRC SEP, Neurology department, CHU de Bordeaux ; Univ. Bordeaux, INSERM U1215, Neurocentre magendie

**Francis Guillemin**

Université de Lorraine, EA 4360 APEMAC ; CHRU de Nancy, INSERM CIC 1433 Epidémiologie Clinique

**Céline Louapre**

Sorbonne Université, Institut du cerveau, ICM, Hôpital de la Pitié Salpêtrière, INSERM S 1127, CNRS UMR 7225 ; Neurology Department, Pitié Salpetriere Hospital, APHP

**Elisabeth Maillart**

Neurology Department, Pitie Salpetriere Hospital, APHP

**Olivier Heinzlef**

department of Neurology, Hôpital CHI de Poissy/Saint-germain-en-Laye

**Nicholas Moore**

Univ. Bordeaux, INSERM CIC-P1401, Bordeaux PharmacEpi

**Droz-Perroteau Cécile**

Univ. Bordeaux, INSERM CIC-P1401, Bordeaux PharmacEpi

---

**Research article**

**Keywords:** Validation study, case-identifying algorithm, claims database, reconstituted Electronic Health Record, Multiple Sclerosis, Prostate Cancer, Positive Predictive Value, Negative Predictive Value

**Posted Date:** January 15th, 2021

**DOI:** <https://doi.org/10.21203/rs.3.rs-145262/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published on May 1st, 2021. See the published version at <https://doi.org/10.1186/s12874-021-01285-y>.

1 **Intra-database validation of case-identifying algorithms using reconstituted electronic**  
2 **health records from healthcare claims data**

3

4 **Authors:** Nicolas H. Thurin<sup>1†</sup>, Pauline Bosco-Levy<sup>1†</sup>, Patrick Blin<sup>1</sup>, Magali Rouyer<sup>1</sup>, Jérémy  
5 Jové<sup>1</sup>, Stéphanie Lamarque<sup>1</sup>, Séverine Lignot<sup>1</sup>, Régis Lassalle<sup>1</sup>, Abdelilah Abouelfath<sup>1</sup>,  
6 Emmanuelle Bignon<sup>1</sup>, Pauline Diez<sup>1</sup>, Marine Gross-Goupil<sup>2</sup>, Michel Soulié<sup>3</sup>, Mathieu  
7 Roumigué<sup>3</sup>, Sylvestre Le Moulec<sup>4</sup>, Marc Debouverie<sup>5,6</sup>, Bruno Brochet<sup>7,8</sup>, Francis  
8 Guillemin<sup>6,9</sup>, Céline Louapre<sup>10,11</sup>, Elisabeth Maillart<sup>11</sup>, Olivier Heinzlef<sup>12</sup>, Nicholas Moore<sup>1</sup>,  
9 Cécile Droz-Perroteau<sup>1</sup>

10

11 † Contributed equally

12

13 **Affiliations:**

14 <sup>1</sup> Univ. Bordeaux, INSERM CIC-P1401, Bordeaux PharmacoEpi, Bordeaux, France

15 <sup>2</sup> Department of Medical Oncology, Hôpital Saint André, CHU de Bordeaux, Bordeaux, France

16 <sup>3</sup> Department of Urology, University Hospital of Rangueil, CHU de Toulouse, Toulouse,  
17 France

18 <sup>4</sup> Department of oncology, Clinique Marzet, Pau, France

19 <sup>5</sup> Department of neurology, CHRU de Nancy, Nancy, France

20 <sup>6</sup> Université de Lorraine, EA 4360 APEMAC, Nancy, France

21 <sup>7</sup> CRC SEP, Neurology department, CHU de Bordeaux, Bordeaux, France

22 <sup>8</sup> Univ. Bordeaux, INSERM U1215, Neurocentre Magendie, Bordeaux, France

23 <sup>9</sup> CHRU de Nancy, INSERM CIC 1433 Epidémiologie clinique, Nancy, France

24 <sup>10</sup> Sorbonne Université, Institut du cerveau, ICM, Hôpital de la Pitié Salpêtrière, INSERM

25 UMR S 1127, CNRS UMR 7225, Paris, France

26 <sup>11</sup> Neurology Department, Pitie Salpetriere Hospital, APHP, Paris, France,

27 <sup>12</sup> Department of Neurology, Hôpital CHI de Poissy/Saint-Germain-en-Laye, Paris, France

28

29 **Corresponding author :** Nicolas H. Thurin, [nicolas.thurin@u-bordeaux.fr](mailto:nicolas.thurin@u-bordeaux.fr)

30

31 **Email addresses:** Nicolas H. Thurin [nicolas.thurin@u-bordeaux.fr](mailto:nicolas.thurin@u-bordeaux.fr); Pauline Bosco-Levy

32 [pauline.bosco-levy@u-bordeaux.fr](mailto:pauline.bosco-levy@u-bordeaux.fr); Patrick Blin [patrick.blin@u-bordeaux.fr](mailto:patrick.blin@u-bordeaux.fr); Magali Rouyer

33 [magali.rouyer@u-bordeaux.fr](mailto:magali.rouyer@u-bordeaux.fr); Jérémy Jové [jeremy.jove@u-bordeaux.fr](mailto:jeremy.jove@u-bordeaux.fr); Stéphanie Lamarque

34 [stephanie.lamarque@u-bordeaux.fr](mailto:stephanie.lamarque@u-bordeaux.fr); Séverine Lignot [severine.lignot@u-bordeaux.fr](mailto:severine.lignot@u-bordeaux.fr); Régis

35 Lassalle [regis.lassalle@u-bordeaux.fr](mailto:regis.lassalle@u-bordeaux.fr); Abdelilah Abouelfath [36 \[bordeaux.fr\]\(mailto:bordeaux.fr\); Emmanuelle Bignon \[emmanuelle.bignon@u-bordeaux.fr\]\(mailto:emmanuelle.bignon@u-bordeaux.fr\); Pauline Diez](mailto:abdelilah.abouelfath@u-</a></p></div><div data-bbox=)

37 [pauline.diez@u-bordeaux.fr](mailto:pauline.diez@u-bordeaux.fr); Marine Gross-Goupil [marine.gross-goupil@chu-bordeaux.fr](mailto:marine.gross-goupil@chu-bordeaux.fr);

38 Michel Soulié [soulie.m@chu-toulouse.fr](mailto:soulie.m@chu-toulouse.fr); Mathieu Roumigué [roumiguie.m@chu-toulouse.fr](mailto:roumiguie.m@chu-toulouse.fr);

39 Sylvestre Le Moulec [sylvestre.lemoulec@gmail.com](mailto:sylvestre.lemoulec@gmail.com); Marc Debouverie [40 \[nancy.fr\]\(mailto:nancy.fr\); Bruno Brochet \[bruno.brochet@chu-bordeaux.fr\]\(mailto:bruno.brochet@chu-bordeaux.fr\); Francis Guillemain](mailto:m.debouverie@chru-</a></p></div><div data-bbox=)

41 [francis.guillemain@chru-nancy.fr](mailto:francis.guillemain@chru-nancy.fr); Céline Louapre [celine.louapre@aphp.fr](mailto:celine.louapre@aphp.fr); Elisabeth Maillart

42 [elisabeth.maillart@aphp.fr](mailto:elisabeth.maillart@aphp.fr); Olivier Heinzlef [Olivier.Heinzlef@ght-yvelinesnord.fr](mailto:Olivier.Heinzlef@ght-yvelinesnord.fr); Nicholas

43 Moore [nicholas.moore@u-bordeaux.fr](mailto:nicholas.moore@u-bordeaux.fr); Cécile Droz-Perroteau [cecile.droz@u-bordeaux.fr](mailto:cecile.droz@u-bordeaux.fr)

44

#### 45 **ORCID**

46 Nicolas H. Thurin 0000-0003-3589-0819; Pauline Bosco-Levy 0000-0001-9763-5974; Patrick

47 Blin 0000-0003-4005-7928; Magali Rouyer 0000-0002-2560-4412; Régis Lassalle 0000-0001-

48 6726-6215; Bruno Brochet 0000-0003-3824-2796; Céline Louapre 0000-0002-4987-1531;

49 Elisabeth Maillart 0000-0001-7699-0328; Nicholas Moore 0000-0003-1212-2817; Cécile

50 Droz-Perroteau 0000-0002-7697-1167

51

52

53 ABSTRACT

54

55 **Background:** Diagnosis performances of case-identifying algorithms developed in healthcare  
56 database are usually assessed by comparing identified cases with an external data source. When  
57 this is not feasible, intra-database validation can present an appropriate alternative.

58 **Objectives:** To illustrate through two practical examples how to perform intra-database  
59 validations of case-identifying algorithms using reconstituted Electronic Health Records  
60 (rEHRs).

61 **Methods:** Patients with 1) multiple sclerosis (MS) relapses and 2) metastatic castration-  
62 resistant prostate cancer (mCRPC) were identified in the French nationwide healthcare database  
63 (SNDS) using two case-identifying algorithms. A validation study was then conduct to estimate  
64 diagnostic performances of these algorithms through the calculation of their positive predictive  
65 value (PPV) and negative predictive value (NPV). To that end, anonymized rEHRs were  
66 generated based on the overall information captured in the SNDS over time (e.g. procedure,  
67 hospital stays, drug dispensing, medical visits) for a random selection of patients identified as  
68 cases or non-cases according to the predefined algorithms. For each diseases, an independent  
69 validation committee reviewed the rEHRs of 100 cases and 100 non-cases in order to adjudicate  
70 on the status of the selected patients (true case/ true non-case), blinded with respect to the result  
71 of the corresponding algorithm.

72 **Results:** Algorithm for relapses identification in MS showed a 95% PPV and 100% NPV and  
73 for mCRPC identification, a 97% PPV and 99% NPV.

74 **Conclusion:** The use of rEHRs to conduct an intra-database validation appears to be a valuable  
75 tool to estimate the performances of a case-identifying algorithm and assess its validity, in the  
76 absence of alternative.

77

78 **KEYWORDS**

79

80 Validation study; case-identifying algorithm; claims database; reconstituted Electronic Health  
81 Record; Multiple Sclerosis; Prostate Cancer; Positive Predictive Value; Negative Predictive  
82 Value

83

84 **DECLARATIONS**

85

86 **Ethics approval and consent to participate:** In both of the presented examples, rEHR  
87 adjudication was used to answer the main research question of the corresponding study for  
88 which the approval from the French data protection agency (*Commission Nationale*  
89 *Informatique & Libertés – CNIL*) was obtained. Moreover, data were fully anonymized before  
90 adjudication.

91

92 **Consent for publication:** Not applicable

93

94 **Availability of data and material:** As per law raw SNDS data cannot be shared. Access to  
95 SNDS data requires approval from the *Comité Ethique et Scientifique pour les Recherches, les*  
96 *Etudes et les Evaluations dans le domaine de la Santé (CESREES)* in charge of assessing  
97 scientific quality of the project, and authorization from the *Commission Nationale de*

98 *l'Informatique et des Libertés* (CNIL) which is the French data protection authority, and then  
99 an agreement with the SNDS data holder (CNAM).

100

101 **Competing interests:**

102 M. G.-G. declares personal fees and non-financial support from Janssen, Sanofi, Astellas, Ipsen,  
103 Amgen and Pfizer.

104 M. S. and M. R. declare personal fees and non-financial support from Janssen, Sanofi, Astellas,  
105 Ipsen, Amgen, Ferring, and Astra-Zeneca.

106 E. M. declares personal fees and non-financial support from Biogen, Novartis, Roche, Merck,  
107 Sanofi-Genzyme.

108 B. B. declares personal fees and non-financial support from Biogen, Genzyme, Bayer, Medday,  
109 Actelion, Roche, Celgene, Novartis, Merck.

110 F. G. and M.D. declare personal fees and non-financial support from Biogen.

111 C. L. declares consulting or travel fees from Biogen, Novartis, Roche, Sanofi, Teva and Merck  
112 Serono, and research grant from Biogen

113 O. H. declares personal fees and non-financial support from Biogen, Merck, Novartis, Roche,  
114 Genzyme

115 All remaining authors have declared no conflicts of interest.

116

117 **Funding:** Presented examples were drawn from two studies funded by Janssen-Cilag, France  
118 and Biogen, France.

119

120 **Authors' contributions:** All authors contributed to the study conception and design. Material  
121 preparation, data collection and analysis were performed by J.J., R.L. and A.A.. All authors  
122 discussed the results. The first draft of the manuscript was written by P. B.-L. and N. H. T.

123 and all authors commented on previous versions of the manuscript. All authors read and  
124 approved the final manuscript.

125

## 126 **Acknowledgements**

127 Both of the presented examples were drawn from studies carried out by the Bordeaux  
128 PharmacoEpi platform in collaboration with Janssen-Cilag, France and Biogen, France and  
129 supervised by independent Scientific Committees. The authors thank all the members of these  
130 scientific committees for their support and advices. The authors thank ADERA for legal, human  
131 resource and management support that made these studies possible.

132

133

## 134 INTRODUCTION

135 For the last two decades, the use of healthcare databases has considerably increased in health  
136 research field [1]. This trend is fueled by the growing recognition that randomized clinical trials,  
137 while essential, are not the unique and exhaustive answer to therapeutic efficacy and safety  
138 issues. The wealth of information that healthcare databases contain, made them robust tool for  
139 many epidemiology-related fields of research, especially in pharmacoepidemiology, where  
140 epidemiologic approaches are applied to well-defined and/or large population to assess the use  
141 and the effects of drugs in real-world practice [2, 3]. The very large amount of data collected  
142 prospectively and systematically in extended period of times, mainly for billing purposes,  
143 enables the assessment of infrequent or delayed adverse events as well as therapeutic long-term  
144 effectiveness, which is complex to assess in classical randomized trials, field cohort or registry  
145 [4]. However, the use of secondary data collected for other purposes than epidemiologic  
146 research, is not devoid of significant limitations [5, 6]. Data quality is a major issue that may  
147 impact case identification by inducing a selection or misclassification bias. In studies conducted

148 on healthcare databases, the population or the health outcome of interest are generally identified  
149 using in- and/or out-patient diagnosis codes. To enhance accuracy, algorithms including  
150 multiple elements specifically related to the studied medical condition (*e.g.* medical procedures,  
151 drug dispensing, laboratory test or radiological exam) in addition to the diagnosis code, may  
152 also be developed and implemented [7]. Whatever the approach used, the coding quality may  
153 be nuanced in terms of how the codes are applied or how physician records are interpreted by  
154 the medical reviewer entering the codes. The financial pressure induced by activity based  
155 payment may also lead to encourage the income-maximizing coding of diagnoses and  
156 procedures in hospitals at the expense of clinical accuracy [8], although more and more quality  
157 audits are carried out to improve coding reliability [9–11]. The validity of algorithm used to  
158 identify health outcome in administrative and claims data has always been a matter of concern  
159 for researchers, especially in a context of active surveillance and assessment of marketed  
160 medical product [12, 13]. Several different types of validation studies may be conducted to  
161 assess the fidelity of the codes or algorithms used for cases identification. In all of them, cases  
162 identified by the algorithm are compared with a presumably more reliable external diagnostic  
163 source or gold standard [14][7]. These gold standards are most of the time the information that  
164 have originated the records in the database (*e.g.* medical charts or registries) and which contain  
165 measure of the disease status based on clinical, biological and/or imaging criteria. The  
166 performance of a case-identifying code or algorithm is commonly reported in term of positive  
167 predictive value, sensitivity and specificity. Although necessary, these validation studies are  
168 time-consuming and require significant resources and expertise to review diagnoses of clinical  
169 data sources. Setting up such a process is also not always possible since the access to the original  
170 data source is often complicated or even impossible because of technical or legal issues.

171 Healthcare databases, are constantly updated with all patient healthcare encounters – medical  
172 visits and procedures, lab tests or medical imaging, drugs dispensing, hospital stays, *etc.* – over

173 a considered period of time, or sometimes a lifetime. They may, by their richness and their  
174 depth, contain information not available in medical charts. Hence, they may provide a holistic  
175 overview of the patient's journeys in real-life settings. These longitudinal patient records could  
176 be seen as reconstituted Electronic Health Records (rEHRs) and so constitutes a valuable  
177 alternative to medical chart in validation study of case-identifying algorithms.

178

179 The objective of this paper is to illustrate through two examples of validation studies conducted  
180 in the French nationwide healthcare database, the *Système National des Données de Santé*  
181 (SNDS) [15], how to perform intra-database validations of case-identifying algorithm using  
182 anonymized rEHRs.

183

## 184 METHODS

### 185 *Data source*

186 Two validation studies were conducted using data from the SNDS, which currently covers more  
187 than 99% of the French population from birth (or immigration) to death (or emigration), even  
188 if a subject moves, changes occupation or retires [15, 16]. Using a unique pseudonymized  
189 identifier, the SNDS merges all reimbursed outpatient claims from all French healthcare  
190 insurance schemes with hospital-discharge summaries from public and private hospitals, and  
191 the national death registry. As a consequence, the SNDS contains information on all reimbursed  
192 medical and paramedical encounters. For each expenditure, the prescriber and caregiver  
193 specialties as well as the corresponding date are provided. The exact quantity of drug dispensed  
194 and reimbursed can be identified at the product level with the exact form and dosage. Performed  
195 laboratory test and procedures are available but without results. Registration for Long Term  
196 Disease (LTD) – status that ensures a full coverage for all related medical expenses – hospital

197 discharge diagnosis and cause of death are defined using codes from the International  
198 Classification of Diseases, 10<sup>th</sup> revision (ICD-10).

199

#### 200 *General method*

201 In the frame of different projects approved by the French regulatory authorities (*Comité*  
202 *d'Expertise pour les Recherches, les Etudes et les Evaluations dans le domaine de la Santé,*  
203 *CEREES* and *Commission Nationale Informatique & Libertés, CNIL*), two algorithms were  
204 developed in the SNDS in collaboration with clinical experts of the field to identify: 1) multiple  
205 sclerosis (MS) relapses and 2) metastatic castration-resistant prostate cancer (mCRPC). The  
206 same methodology for intra-database validation was then applied to each of them in order to  
207 ascertain that patients identified as cases or non-cases by the algorithm were respectively true  
208 cases and true non-cases.

209 In a first step, anonymized longitudinal rEHRs were generated based on SNDS data for a  
210 random selection of 100 patients identified by the algorithm as case and 100 patients identified  
211 by the algorithm as non-case (Figure 1). To ensure that individual data contained in these rEHRs  
212 did not lead to patient re-identification, new patient identifiers were assigned, calendar dates  
213 were replaced by the delay elapsed since inclusion, location details were deleted and only age  
214 classes were displayed. In a second step, a validation committee consisting of medical experts  
215 of the field, proceeded to a double review of the rEHRs in order to adjudicate on the true case  
216 or true non-case status of the selected patients, blinded with respect to the algorithm result. In  
217 case of discrepancy, committee members discussed to reach a consensus. In a final step, experts  
218 conclusions were compared with algorithm results to estimate its diagnostic performance  
219 through the positive predictive value (PPV), the negative predictive value (NPV) and their  
220 corresponding 95% confidence intervals (95%CI). The formulae for PPV and NPV were:

221

$$PPV = \frac{TP}{TP + FP}$$

222 
$$NPV = \frac{TN}{TN + FN}$$

223 where TP and FP are respectively true and false positives, and TN and FN true and false  
224 negatives. Corresponding formula for 95% CIs were:

225 
$$[95\%CI]_{PPV} = PPV \pm z_{1-\frac{\alpha}{2}} * \sqrt{\frac{PPV(1 - PPV)}{n_{positive}}}$$

226 
$$[95\%CI]_{NPV} = NPV \pm z_{1-\frac{\alpha}{2}} * \sqrt{\frac{NPV(1 - NPV)}{n_{negative}}}$$

227 where  $n_{positive}$  and  $n_{negative}$  are respectively the number of algorithm-based positive and negative  
228 assessed cases, and  $z_{(1-\alpha/2)}$  the z-value for standard normal distribution with left-tail probability  
229  $(1-\alpha/2)$ . Here  $z_{(1-\alpha/2)} = 1.96$  for a type-1 error  $\alpha = 0.05$ .

230 The limitation of the number of rEHRs to be assessed per group to 100 allowed to estimate PPV  
231 and NPV with a margin of error <10% for values above 50%, and made the adjudication of  
232 cases possible by experts in less than 48h.

233 Following the validation study, experts' inputs were used to adjust algorithm settings and  
234 further improve its discriminatory ability. Overall estimated performances indicators were then  
235 updated.

236

### 237 *Case examples*

#### 238 *1) Relapses identification in Multiple Sclerosis (MS) patients*

239 The algorithm for identifying relapse in MS was initially developed in the EVIDEMS study  
240 whose objective was to assess the effectiveness of dimethyl fumarate versus other MS drugs  
241 (*i.e.* teriflunomide, fingolimod or immunomodulatory injectable drugs) on relapses after  
242 treatment initiation [17]. The study cohort included all patients identified in the SNDS by a first  
243 dispensing of MS drug (*i.e.* dimethyl fumarate, indicated for MS) between July 2015 and  
244 December 2017, with 4.5-year history and 1 to 3.5 years of follow-up. Relapses were identified

245 using a complex algorithm including dispensing of high dose of corticosteroids  
246 (methylprednisolone or betamethasone) for outpatient, and hospitalizations with MS relapse  
247 diagnosis whether or not combined with high dose of corticosteroids.

248

## 249 2) *Patients with metastatic Castration-Resistant Prostate Cancer (mCRPC)*

250 The algorithm for mCRPC patients identification was initially developed in the CAMERRA  
251 study whose objectives were to assess mCRPC burden and describe mCRPC-specific treatment  
252 lines [18, 19]. The study cohort included all patients with a prostate cancer identified in the  
253 SNDS by a specific hospital discharge or LTD diagnosis code or a specific treatment dispensing  
254 between January 2009 and December 2014, with 5-year history and 3 years of follow-up.  
255 Patients with mCRPC were identified using an algorithm integrating time indicators related to  
256 metastases management and castration resistance. Both indicators relied on the detection of  
257 specific procedures (e.g. imaging, surgery or radiotherapy), drug dispensing (e.g. androgen  
258 deprivation therapy, metastases-targeted treatment, chemotherapy) or specific hospitalizations.  
259 A complete description of the algorithm and its validation are available elsewhere [18]. In the  
260 CAMERRA validation study, so as to ensure the presence of all categories of non-mCRPC  
261 patients, three groups of non-mCRPC patients were identified: 34 with non-metastatic  
262 hormone-sensitive prostate cancer, 33 with metastatic hormone-sensitive prostate cancer, and  
263 33 with non-metastatic castration-resistant prostate cancer. A single NPV relying on the overall  
264 non-mCRPC population was then estimated for the algorithm by weighting false-negative cases  
265 according to the actual distribution of the 3 categories of non-mCRPC patient in the prostate  
266 cancer population. In a last stage, PPV, NPV and the observed prevalence of mCRPC among  
267 the study population (prostate cancer patients), were used to derived the sensitivity and  
268 specificity [20].

269

270 RESULTS

271

272 *Diagnostic performance of the MS relapse algorithm*

273 A sample of 200 patients was randomly selected from the initial study population; 100 of them  
274 had at least one relapse and 100 did not have any relapse according to the algorithm. The  
275 validation committee confirmed 95 patients with relapses (true cases) among the algorithm-  
276 identified cases and 96 without relapse among the algorithm-identified non-cases, resulting in  
277 a PPV of 95.0% (95%IC = [91; 99]) and a NPV of 96.0% (95%IC = [92; 100]) (Appendix 1-  
278 A). After the update of algorithm settings based on experts' conclusions, NPV reached 100.0%

279 (TABLES

280 Table 1).

281

282 *Diagnostic performance of the mCRPC algorithm*

283 A sample of 200 patients was randomly selected from the initial population with prostate  
284 cancer; 100 of them were identified as mCRPC and 100 as non-mCRPC according to the  
285 algorithm. Experts confirmed 92 of the 100 algorithm-identified mCRPC cases and 93 of the  
286 100 algorithm-identified non-mCRPC cases, resulting in a PPV and NPV of respectively 92.0%  
287 (95%CI= [87; 97]) and 99.0% (95%CI= [98; 100]), after weighting according to non-mCRPC  
288 cases distribution (Appendix 1-B). Following the algorithm adjustment based on expert  
289 feedback, PPV reached 97% (95%CI= [93; 100]) (Table 2). Based on an observed proportion  
290 mCRPC/prostate cancer of 3.4%, sensitivity and specificity of the final algorithm were  
291 respectively estimated at 80% and 100% [18].

292

293 DISCUSSION

294

295 Based on two practical examples relying on the French nationwide healthcare database, this  
296 paper illustrates an innovative method to assess case-identifying algorithms conducting an  
297 intra-database validation study. In both examples, this validation study showed that algorithms  
298 had high diagnostic performances, with excellent PPV and NPV. To our knowledge, there are  
299 no previous examples of the use of rEHRs to assess the performances of algorithms for case  
300 identification, therefore results are difficult to compare with existing data. Because by law,  
301 returning to individual medical records from SNDS data is forbidden, most of the currently  
302 published French validation studies were limited to the comparison of hospital discharge codes  
303 extracted from local hospital databases – before their de-identification and integration to the  
304 SNDS – with traditional sources of information such as medical charts or registries, and leading  
305 to a PPV varying from 80 to 90% but tending to decrease according to the granularity of the  
306 required information [21–27].

307 In the present case, intra-database validation provides the opportunity to assess algorithms that  
308 rely on multiple elements from SNDS, enabling to improve discriminatory abilities compared  
309 to single identification criterion or to overcome the absence of a direct-identifying diagnostic  
310 code [28]. Experts of the validation committee reported that rEHRs proceeding from SNDS  
311 data were on certain points more informative than the usual medical charts and contained a high  
312 level of details as well as an accurate chronology regarding patients' journeys, which generally  
313 made the adjudication of the cases non-ambiguous. Clinicians insights also allowed to refine  
314 the algorithm, adjusting its settings to further improve its performances. This suggests that  
315 SNDS data are comprehensive enough to develop a complex algorithm and to validate it.

316 Validation studies based on medical charts review stay the best way to evaluate claims database  
317 algorithms. However, it requires a lot of human time and reliable significant funding, which are  
318 often missing, to be able most often to estimate only the PPV. Wherever feasible, validation  
319 studies relying on linkage between administrative databases and medical registries or electronic

320 medical record databases are a good alternative, but they are rarely fully representative of the  
321 whole database population, and stay quite long and expensive. Conversely, rEHR review offer  
322 a time- and cost-efficient way to conduct validation studies, using the data source accessed by  
323 the algorithm. Files to review are standardized and structured, allowing the assessment of  
324 hundreds of cases in a limited time: 50 cases per expert-day in the two presented examples.  
325 Moreover, as both cases and non-cases are accessible, this approach enables the calculation of  
326 other indicators than PPV (e.g. NPV, sensitivity, specificity), with a full representativeness of  
327 the population covered by the database.

328 We acknowledge that validating an algorithm in the same database that was used to develop it  
329 may be questionable. The suitability of using an unique data source to generate and evaluate a  
330 hypothesis has been previously discussed in the scientific literature, even if the scope was  
331 slightly different [29–32]. Walker AM., and Wang SV. and colleagues argued that for such an  
332 approach to be considered valid “test data need to be independent of hypothesis-generating  
333 data” [29, 30]. Though it is consensual that re-using data to perform quality check, reevaluate  
334 findings and strengthen hypotheses (e.g. sensitivity analyses) in the frame of  
335 pharmacoepidemiology studies belongs to good research practice, the fact that they can also be  
336 used to validate hypotheses is more challenged, especially because of the potential lack of  
337 argument to establish causality [31]. In hypothesis-evaluating treatment effectiveness studies,  
338 the reuse of data sources is usually not recommended upon the main argument that it leads to  
339 replication rather than confirmation [30, 32]. In the present work the lack of argument to  
340 establish causality is not an issue, as we do not seek it; the unique objective is to prove that  
341 cases identified by the algorithm are true cases. Moreover, here, the independence is ensured  
342 by the unrelated approaches used in the identification and the confirmation of the cases: to  
343 classify a case, the algorithm picks up information in the database as previously defined in a  
344 statistical analysis plan. When experts do so, they choose the relevant information for

345 themselves in the rEHR. Relying on the same data, the elements considered can be the same (or  
346 not), but the approach and the selection process are different, resulting in independent bodies  
347 of proof.

348 Obviously, preference must be given to external data source to conduct validation study, but  
349 when it is not feasible, the re-use of the original data source appears as a valuable alternative.

350 Moreover, it should be borne in mind that the decision whether or not to proceed with intra-  
351 database validation for case-identifying algorithm will strongly depend both on the nature of  
352 the outcome of interest and on the characteristics of the considered database. Two conditions  
353 must be fulfilled to ensure an effective application of the method: 1) the health outcome of  
354 interest must be managed by a specific sequence of cares and encounters; 2) the considered  
355 healthcare database must capture in an exhaustive way a sufficient number of medical elements  
356 in line with the outcome of interest.

357 Outcome validation should not rely on a unique diagnostic or procedure code but on several  
358 tangible elements. As a consequence, intra-database validation should only be considered for  
359 health outcomes that are managed in usual clinical practice by a well-defined chronological  
360 sequence of cares (procedures, drug dispensings, hospital stays, medical visits, *etc.*) since  
361 diagnostic evidence – such as images or laboratory results – may be absent of the database. The  
362 succession of healthcare encounters, that individually may be unspecific of the outcome when  
363 taken together, gives rise to a specific healthcare pathway. This is particularly true for serious  
364 outcomes, mobilizing large healthcare resources. Thus, chronic conditions such as MS (see  
365 example 1) or cancer (see example 2), or serious acute outcomes for which the management  
366 follows consensual and structured guidelines (e.g. myocardial infarction) [33] seem to be better  
367 suited to intra-database validation, compared to non-serious outcomes involving few and  
368 unspecific healthcare resources (e.g. acute sore throat) [34], or serious but rare diseases with no  
369 clinical practice guidelines [35]. Particular attention must be paid to the clinical guidelines

370 which were ongoing at the time of the study, since they drive patients journeys and thus,  
371 experts' judgment.

372 Furthermore, in order to ensure that rEHRs provide sufficient and reliable information to enable  
373 case adjudication, the underlying healthcare database must capture a sufficient number of  
374 medical elements in an exhaustive way over a suitable period of time. Data collected must be,  
375 at least, in line with the type of care involved in the management of the outcome of interest (e.g.  
376 validation of a myocardial infarction identification requires outpatient data). Ideally, outpatient  
377 and inpatient healthcare encounters should be included, and the quality of the captured  
378 information regularly assessed. Data completeness, at least over the study period, is mandatory  
379 to ensure that the absence of record is synonym of an absence of encounter. The SNDS is  
380 particularly well suited to this situation since it fulfills all these requirements: it includes in-  
381 and out-patient information of all reimbursed healthcare encounters, most of the time lifelong,  
382 and the quality of coding is ensured by regular internal and external audits [9–11].

383

## 384 CONCLUSION

385

386 Homogeneous healthcare databases such as the SNDS captures healthcare journey of patients  
387 lifelong. Although these data cannot replace the anamnesis and the clinical information reported  
388 in patient medical charts, this succession of healthcare records appears to be comprehensive  
389 enough to generate consistent rEHRs assessable by experts, allowing to conduct validation  
390 studies without using external information. It should be made clear that intra-database  
391 validation based on rEHRs review does not pretend to replace traditional methods of validation  
392 relying on medical charts review. However, as illustrated here through the MS relapse example  
393 and the mCRPC example, in the absence of alternative, such method appears to be a valuable  
394 tool to estimate the performances of a case-identifying algorithms and assess their validity. The

395 development in the coming years of data linkages allowing to gather claims data, registries,  
396 electronic health records, *etc.* [36, 37], will further enrich data available for experts to review  
397 in rEHRs and may blur the line between intra-validation and external medical chart review.

398

399

#### 400 LIST OF ABBREVIATIONS

401 CEREEES: Comité d'Expertise pour les Recherches, les Etudes et les Evaluations dans le  
402 domaine de la Santé; CNIL: Commission Nationale Informatique & Libertés; FN: false  
403 negative; FP: false positive; ICD-10: International Classification of Diseases, 10th revision;  
404 LTD: Long Term Disease; mCRPC: metastatic castration-resistant prostate cancer; MS:  
405 Multiple Sclerosis; NPV: negative predictive value; PPV: positive predictive value; rEHR:  
406 reconstituted Electronic Health Record; SNDS: Système National des Données de Santé; TN:  
407 True negative; TP: True positive; 95%CI: 95% confidence interval

408

409

#### 410 REFERENCES

411

412 1. Ray WA (2011) Improving Automated Database Studies. *Epidemiology* 22:302.

413 <https://doi.org/10.1097/EDE.0b013e31820f31e1>

414 2. Gavriellov-Yusim N, Friger M (2014) Use of administrative medical databases in  
415 population-based research. *J Epidemiol Community Health* 68:283–287.

416 <https://doi.org/10.1136/jech-2013-202744>

417 3. Strom BL (2019) What Is Pharmacoepidemiology? In: *Pharmacoepidemiology*. John  
418 Wiley & Sons, Ltd, pp 1–26

419 4. Hennessy S (2006) Use of Health Care Databases in Pharmacoepidemiology. *Basic*

- 420 Clin Pharmacol Toxicol 98:311–313. [https://doi.org/10.1111/j.1742-7843.2006.pto\\_368.x](https://doi.org/10.1111/j.1742-7843.2006.pto_368.x)
- 421 5. Hashimoto RE, Brodt ED, Skelly AC, Dettori JR (2014) Administrative Database  
422 Studies: Goldmine or Goose Chase? *Evid-Based Spine-Care J* 05:74–76.  
423 <https://doi.org/10.1055/s-0034-1390027>
- 424 6. Grimes DA (2010) Epidemiologic research using administrative databases: garbage in,  
425 garbage out. *Obstet Gynecol* 116:1018–1019.  
426 <https://doi.org/10.1097/AOG.0b013e3181f98300>
- 427 7. Lanes S, Brown JS, Haynes K, et al (2015) Identifying health outcomes in healthcare  
428 databases. *Pharmacoepidemiol Drug Saf* 24:1009–1016. <https://doi.org/10.1002/pds.3856>
- 429 8. Georgescu I, Hartmann FGH (2013) Sources of financial pressure and up coding  
430 behavior in French public hospitals. *Health Policy* 110:156–163.  
431 <https://doi.org/10.1016/j.healthpol.2013.02.003>
- 432 9. Gilleron V, Gasnier-Duparc N, Hebbrecht G (2018) Certification des comptes: Une  
433 incitation à la traçabilité des processus de contrôle. *Revue Hospitaliere de France* 582:6
- 434 10. Marescaux C (2011) Entre soin et contrôle de gestion : place du DIM dans  
435 l'organisation hospitalière. *Inf Psychiatr Volume* 87:487–491
- 436 11. Caeyseele T, Bruandet A, Delaby F, Theis D (2016) Création d'un outil de gestion des  
437 contrôles qualités du codage au DIM du CHRU de Lille. *Rev DÉpidémiologie Santé Publique*  
438 64:S20. <https://doi.org/10.1016/j.respe.2016.01.066>
- 439 12. Carnahan RM (2012) Mini-Sentinel's systematic reviews of validated methods for  
440 identifying health outcomes using administrative data: summary of findings and suggestions  
441 for future research: HEALTH OUTCOME ALGORITHM SUMMARY. *Pharmacoepidemiol*  
442 *Drug Saf* 21:90–99. <https://doi.org/10.1002/pds.2318>
- 443 13. Carnahan RM, Moores KG (2012) Mini-Sentinel's systematic reviews of validated  
444 methods for identifying health outcomes using administrative and claims data: methods and

445 lessons learned: HEALTH OUTCOME ALGORITHM REVIEW METHODS.  
446 *Pharmacoepidemiol Drug Saf* 21:82–89. <https://doi.org/10.1002/pds.2321>

447 14. van Walraven C, Bennett C, Forster AJ (2011) Administrative database research  
448 infrequently used validated diagnostic or procedural codes. *J Clin Epidemiol* 64:1054–1059.  
449 <https://doi.org/10.1016/j.jclinepi.2011.01.001>

450 15. Bezin J, Duong M, Lassalle R, et al (2017) The national healthcare system claims  
451 databases in France, SNIIRAM and EGB: Powerful tools for pharmacoepidemiology.  
452 *Pharmacoepidemiol Drug Saf* 26:954–962. <https://doi.org/10.1002/pds.4233>

453 16. Tuppin P, Rudant J, Constantinou P, et al (2017) Value of a national administrative  
454 database to guide public decisions: From the système national d’information interrégimes de  
455 l’Assurance Maladie (SNIIRAM) to the système national des données de santé (SNDS) in  
456 France. *Rev Epidemiol Sante Publique* 65 Suppl 4:S149–S167.  
457 <https://doi.org/10.1016/j.respe.2017.05.004>

458 17. Bosco-Levy P, Debouverie M, Brochet B, et al Comparative effectiveness of dimethyl  
459 fumarate versus other disease-modifying therapies in Multiple Sclerosis: a large population-  
460 based cohort study. Under Submission

461 18. Thurin NH, Rouyer M, Gross-Goupil M, et al (2020) Epidemiology of metastatic  
462 castration-resistant prostate cancer: A first estimate of incidence and prevalence using the  
463 French nationwide healthcare database. *Cancer Epidemiol* 69:101833.  
464 <https://doi.org/10.1016/j.canep.2020.101833>

465 19. Gross-Goupil M, Thurin NH, Rouyer M, et al (2020) Survival outcome in patients  
466 with metastatic castration-resistant prostate cancer according to first-line treatment. *J Clin*  
467 *Oncol* 38:5570–5570. [https://doi.org/10.1200/JCO.2020.38.15\\_suppl.5570](https://doi.org/10.1200/JCO.2020.38.15_suppl.5570)

468 20. Bollaerts K, Rekkas A, Smedt TD, et al (2020) Disease misclassification in electronic  
469 healthcare database studies: Deriving validity indices—A contribution from the ADVANCE

- 470 project. PLOS ONE 15:e0231333. <https://doi.org/10.1371/journal.pone.0231333>
- 471 21. Bosco-Lévy P, Duret S, Picard F, et al (2019) Diagnostic accuracy of the International  
472 Classification of Diseases, Tenth Revision, codes of heart failure in an administrative  
473 database. *Pharmacoepidemiol Drug Saf* 28:194–200. <https://doi.org/10.1002/pds.4690>
- 474 22. Coureau G, Baldi I, Savès M, et al (2012) [Performance evaluation of hospital claims  
475 database for the identification of incident central nervous system tumors compared with a  
476 cancer registry in Gironde, France, 2004]. *Rev Epidemiol Sante Publique* 60:295–304.  
477 <https://doi.org/10.1016/j.respe.2012.02.003>
- 478 23. Giroud M, Hommel M, Benzenine E, et al (2015) Positive Predictive Value of French  
479 Hospitalization Discharge Codes for Stroke and Transient Ischemic Attack. *Eur Neurol*  
480 74:92–99. <https://doi.org/10.1159/000438859>
- 481 24. Goueslard K, Cottenet J, Benzenine E, et al (2020) Validation study: evaluation of the  
482 metrological quality of French hospital data for perinatal algorithms. *BMJ Open* 10:e035218.  
483 <https://doi.org/10.1136/bmjopen-2019-035218>
- 484 25. Mezaache S, Derumeaux H, Ferraro P, et al (2017) Validation of an algorithm  
485 identifying incident primary immune thrombocytopenia in the French national health  
486 insurance database. *Eur J Haematol* 99:344–349. <https://doi.org/10.1111/ejh.12926>
- 487 26. Palmaro A, Gauthier M, Conte C, et al (2017) Identifying multiple myeloma patients  
488 using data from the French health insurance databases. *Medicine (Baltimore)* 96:.  
489 <https://doi.org/10.1097/MD.00000000000006189>
- 490 27. Prat M, Derumeaux H, Sailer L, et al (2018) Positive predictive values of peripheral  
491 arterial and venous thrombosis codes in French hospital database. *Fundam Clin Pharmacol*  
492 32:108–113. <https://doi.org/10.1111/fcp.12326>
- 493 28. Fuentes S, Cosson E, Mandereau-Bruno L, et al (2019) Identifying diabetes cases in  
494 health administrative databases: a validation study based on a large French cohort. *Int J Public*

- 495 Health 64:441–450. <https://doi.org/10.1007/s00038-018-1186-3>
- 496 29. Walker AM (2010) Orthogonal predictions: follow-up questions for suggestive data.  
497 *Pharmacoepidemiol Drug Saf* 19:529–532. <https://doi.org/10.1002/pds.1929>
- 498 30. Wang SV, Kulldorff M, Glynn RJ, et al (2018) Reuse of data sources to evaluate drug  
499 safety signals: When is it appropriate? *Pharmacoepidemiol Drug Saf* 27:567–569.  
500 <https://doi.org/10.1002/pds.4442>
- 501 31. Gould AL (2010) Generating and confirming hypotheses. *Pharmacoepidemiol Drug*  
502 *Saf* 19:533–536. <https://doi.org/10.1002/pds.1928>
- 503 32. Berger ML, Sox H, Willke RJ, et al (2017) Good practices for real-world data studies  
504 of treatment and/or comparative effectiveness: Recommendations from the joint ISPOR-ISPE  
505 Special Task Force on real-world evidence in health care decision making.  
506 *Pharmacoepidemiol Drug Saf* 26:1033–1039. <https://doi.org/10.1002/pds.4297>
- 507 33. Ibanez B, James S, Agewall S, et al (2018) 2017 ESC Guidelines for the management  
508 of acute myocardial infarction in patients presenting with ST-segment elevationThe Task  
509 Force for the management of acute myocardial infarction in patients presenting with ST-  
510 segment elevation of the European Society of Cardiology (ESC). *Eur Heart J* 39:119–177.  
511 <https://doi.org/10.1093/eurheartj/ehx393>
- 512 34. Pelucchi C, Grigoryan L, Galeone C, et al (2012) Guideline for the management of  
513 acute sore throat: ESCMID Sore Throat Guideline Group. *Clin Microbiol Infect* 18:1–27.  
514 <https://doi.org/10.1111/j.1469-0691.2012.03766.x>
- 515 35. Pavan S, Rommel K, Mateo Marquina ME, et al (2017) Clinical Practice Guidelines  
516 for Rare Diseases: The Orphanet Database. *PLoS ONE* 12:.  
517 <https://doi.org/10.1371/journal.pone.0170365>
- 518 36. Bradley CJ, Penberthy L, Devers KJ, Holden DJ (2010) Health Services Research and  
519 Data Linkages: Issues, Methods, and Directions for the Future. *Health Serv Res* 45:1468.

520 <https://doi.org/10.1111/j.1475-6773.2010.01142.x>

521 37. Scailteux L-M, Droitcourt C, Balusson F, et al (2019) French administrative health

522 care database (SNDS): The value of its enrichment. *Therapie* 74:215–223.

523 <https://doi.org/10.1016/j.therap.2018.09.072>

524

525

## 526 TABLES

527 Table 1. Positive (PPV) and negative (NPV) predictive values of the final algorithm for the  
528 identification of relapse in multiple sclerosis

		Validation committee		
		Relapse +	Relapse -	Total
Algo- ithm	Relapse +	99	5	104
	Relapse -	0	96	96
	Total	99	101	200

PPV = 95% (95%CI = [91; 99])  
NPV = 100%

529

530 Table 2. Positive (PPV) and negative (NPV) predictive values of the final algorithm for the  
531 identification of metastatic castration-resistant prostate cancer (mCRPC), adapted from  
532 Thurin NH, et al. 2020

		Validation committee		
		mCRPC +	mCRPC -	Total
Algo- ithm	mCRPC +	90	3	93
	mCRPC -	1.23*	105.77*	107
	Total	91.23	108.77	200

PPV = 97% (95%CI= [93; 100])  
NPV = 99% (95%CI= [97; 100])

533 \*After weighting

534

535

536

537

538

539

## 540 FIGURE CAPTIONS

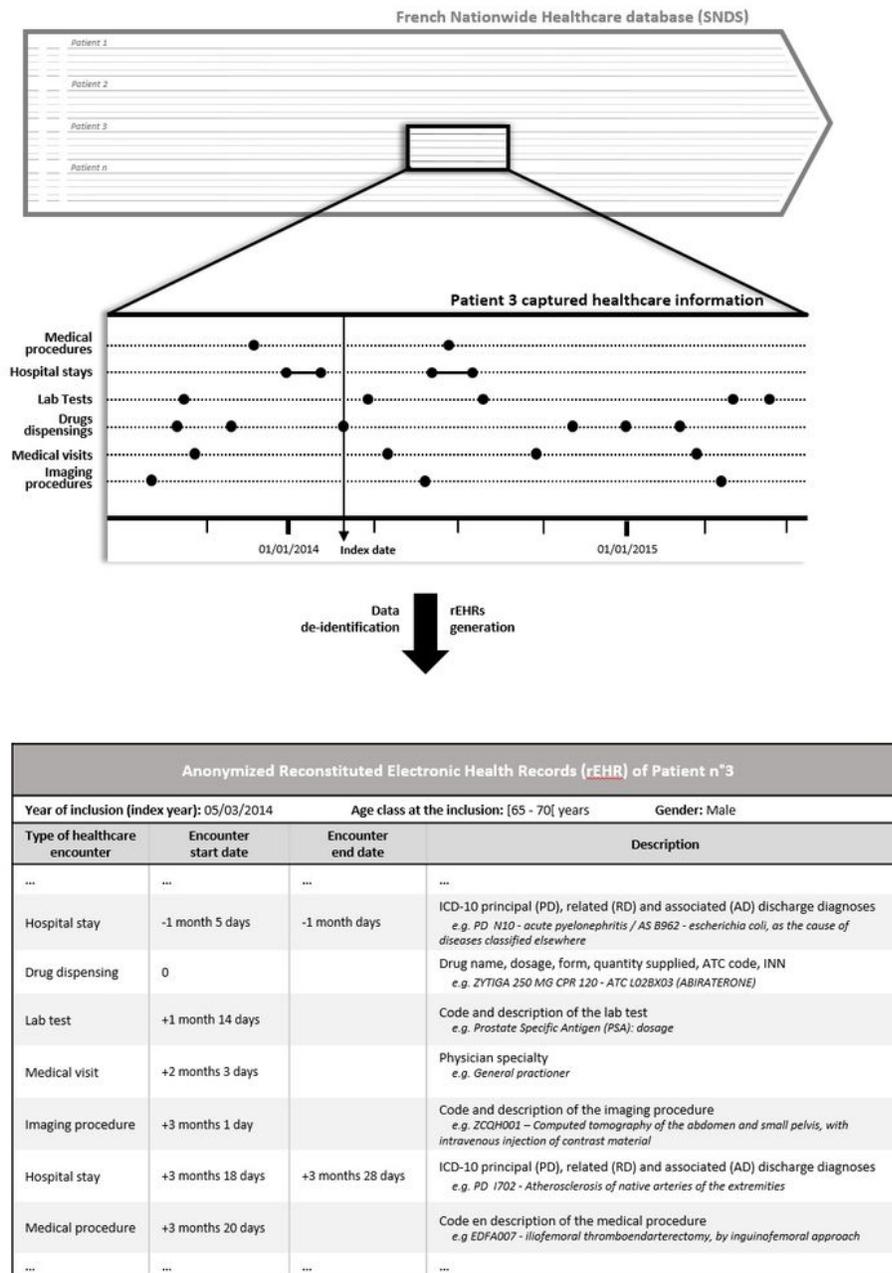
541 Fig. 1 Generation of anonymized reconstituted Electronic Health Records from the French

542 Nationwide Healthcare database (SNDS) to conduct intra-database case-identifying algorithm

543 validation

544

# Figures



**Figure 1**

Generation of anonymized reconstituted Electronic Health Records from the French Nationwide Healthcare database (SNDS) to conduct intra-database case-identifying algorithm validation

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [APPENDIX1.pdf](#)