

1 Target journal: BMJ /Article type: RESEARCH METHODS AND REPORTING

2

3 **Validation of algorithms in studies based on routinely collected health data: general**

4 **principles**

5 Vera Ehrenstein, MPH, DSc,<sup>1</sup> Maja Hellfritzsich, MD, PhD,<sup>2</sup> Johnny Kahlert, PhD,<sup>1</sup> Sinéad M

6 Langan, FRCP, MSc, PhD,<sup>3</sup> Hisashi Urushihara, MSc, DrPH,<sup>4</sup> Danica Marinac-Dabic, MD, PhD,

7 MMSc, FISPE,<sup>5</sup> Jennifer L Lund, PhD,<sup>6</sup> Henrik Toft Sørensen, MD, PhD, DMSc, DSc,<sup>1</sup> Eric I

8 Benchimol, MD, PhD<sup>7,8,9</sup>

9 **Corresponding author:** Dr. Vera Ehrenstein, Department of Clinical Epidemiology,

10 Department of Clinical Medicine, Aarhus University, Olof Palmes Allé 43-45, Aarhus N,

11 Denmark, [ve@clin.au.dk](mailto:ve@clin.au.dk)

12 Author affiliations:

13 <sup>1</sup> Department of Clinical Epidemiology, Aarhus University, Aarhus, Denmark

14 <sup>2</sup> Clinical Pharmacology, Pharmacy and Environmental Medicine, Department of Public Health,  
15 University of Southern Denmark, Odense, Denmark

16 <sup>3</sup> Department of Non-communicable Disease Epidemiology, Faculty of Epidemiology and  
17 Population Health, London School of Hygiene and Tropical Medicine, London, United Kingdom

18 <sup>4</sup> Division of Drug Development & Regulatory Science Faculty of Pharmacy, Keio University,  
19 Tokyo, Japan

20 <sup>5</sup> Office of Clinical Evidence and Analysis, Center for Devices and Radiological, United States  
21 Food and Drug Administration, Silver Spring MD, USA

22 <sup>6</sup> Department of Epidemiology, Gillings School of Global Public Health, University of North  
23 Carolina at Chapel Hill, Chapel Hill NC, USA

24 <sup>7</sup> Division of Gastroenterology, Hepatology and Nutrition and Child Health Evaluative Sciences,  
25 SickKids Research Institute, The Hospital for Sick Children, Toronto, ON, Canada

26 <sup>8</sup> Department of Paediatrics and Institute of Health Policy, Management and Evaluation,  
27 University of Toronto, Toronto, ON, Canada

28 <sup>9</sup> ICES, Toronto, ON, Canada

29 **Word count**, excluding title page, abstract, references, boxes and figures: **approx. 6,136**

30

31 **Abstract (140 words)**

32 Clinicians, researchers, regulators and other decision-makers are increasingly relying on  
33 evidence based on routinely collected health data (RCD). This paper aims to systematise  
34 terminology, methods, and practical considerations relevant to the conduct of validation studies  
35 of RCD-based algorithms. First, we define algorithms in the context of RCD-based research and  
36 summarise measures of algorithm accuracy. Second, we discuss the concepts of gold standard  
37 and reference standard. Third, we offer considerations for determining study size, prioritising  
38 accuracy measures, and algorithm portability. Finally, we briefly discuss the use of validity data  
39 in interpreting results. We use published studies to illustrate all points. In all studies, some  
40 degree of information and other biases is inevitable, thus decisions regarding prioritising  
41 measures of algorithm accuracy and basis for those decisions should be transparently reported.  
42 Validation work should be an ongoing routine task of RCD source maintenance.

43 **Key words**

44 Accuracy; algorithm; data quality, epidemiology; information bias; measurement error;  
45 misclassification; observational studies; routinely collected health data; real-world data, real-  
46 world evidence, validity

47

48 **Introduction**

49 *Real-world data* (RWD) are data originating in the course of routine clinical practice, and the  
50 term is used by contrast with data collected in phase I-III interventional trials.<sup>1-3</sup> *Real-world*  
51 *evidence* (RWE) is evidence generated by analysing RWD. RWD is represented to a large extent  
52 by *routinely collected health data* (RCD), i.e., data accruing as a by-product of health care  
53 delivery, encompassing electronic health records; health administrative data; claims data, records  
54 of treatments, procedures, devices, or diseases;<sup>4</sup> and patient-reported events.<sup>5</sup> RCD-based studies  
55 are playing an increasingly important role in decisions by regulators, payers, and in health  
56 technology assessment,<sup>1-3 7-14</sup> as corroborated by the 2021 publication of the CONSORT-  
57 ROUTINE - extension for the reporting of randomised controlled trials conducted using cohorts  
58 and routinely collected data.<sup>15</sup> Rapid accumulation of RCD in diverse health care systems,  
59 coupled with expanding computing capabilities, predict the ever-strengthening role of RCD in  
60 generating evidence on disease prevention, aetiology, epidemiology, clinical course and quality  
61 of care, and on utilisation, benefits, risks, and costs of treatments and devices.<sup>18-21</sup>

62  
63 The advantages of RCD/RWE –include accrual independent of research (and thus free of  
64 investigator’s preconceptions), comparatively low cost, large number of observations, rapid  
65 availability for analysis, analysis-friendly structure, standard coding for many data types – afford  
66 precise results from diverse patient populations in routine care settings. An important  
67 disadvantage of the RCD is susceptibility to measurement error, introduced along the path from  
68 the point of care to a data point in an analytic dataset.<sup>22</sup> In RCD-based studies this causes  
69 *information bias*, i.e., systematic error resulting from imperfect correspondence between the true  
70 (health) event and its representation in a RCD record, e.g. misclassification. *Random*  
71 measurement errors on the level of data accrual translate into *systematic* error on the level of

72 study result,<sup>23-27</sup> potentially leading to wrong conclusions, and, ultimately, suboptimal clinical  
73 practice.<sup>28</sup> Although misclassification is inevitable in any epidemiologic study, regardless of data  
74 collection methods,<sup>28</sup> the extent, the nature, and the impact of misclassification for any given  
75 study using RCD may be difficult to quantify.<sup>29</sup> Large-scale RCD-based studies, yielding results  
76 that are precise but biased are especially troublesome, because large size and associated precision  
77 may create a false sense of trust among RWE consumers, if precision is mistaken for validity.  
78 Thus, “without taking data quality into account [...] the more the data, the surer we fool  
79 ourselves”.<sup>33</sup> For RCD-based studies to provide high-quality evidence for clinical practice, data  
80 validity must be quantified and the quantitative information incorporated in the analysis and  
81 interpretation of study results.<sup>34-36</sup> There have been several comprehensive efforts to validate  
82 sources of RCD,<sup>37-43</sup> while guidelines for assessing and reporting validation studies<sup>30</sup> have  
83 helped improve reporting transparency.<sup>44-51</sup> Nevertheless, literature on validating RCD remains  
84 fragmented,<sup>31 52-55</sup> and scientists have repeatedly raised awareness about the need for better  
85 education in the realm of validation studies.<sup>28 56</sup>

86  
87 This paper aims to systematise terminology, methods, and practical considerations relevant to  
88 validation studies of RCD-based algorithms. First, we define algorithms in the context of RCD  
89 and summarise measures of algorithm validity. Second, we discuss the concepts of gold standard  
90 and reference standard. Third, we consider prioritising validity measures, discuss portability of  
91 algorithms, and offer considerations for determining the size of a validation study. Finally, we  
92 discuss the role of validation information in interpreting study results. We use published RCD-  
93 based studies for illustration. Our intended audience includes RWE originators and consumers in  
94 academia, regulatory science, industry and other decision-makers involved in planning,  
95 conducting and interpreting RCD-based studies. This effort aligns with the Guidelines for Good

96 Pharmacoepidemiology Practices (GPP) maintained by the International Society for  
97 Pharmacoepidemiology (ISPE);<sup>57</sup> with the data-quality-focussed regulatory strategy of the  
98 European Medicines Agency,<sup>58</sup> and with the RWE transparency initiative.<sup>59</sup>

99

## 100 **Manuscript development**

101 This paper was conceptualised as a contribution of the ongoing efforts to improve conduct and  
102 reporting of observational studies,<sup>3 44-51 60 61</sup> specifically, validation studies in RCD sources.<sup>32 34 52</sup>

103 <sup>54 62</sup> The authors are ISPE members with clinical, academic, or regulatory affiliations in Europe,  
104 North America, and Asia and with collective expertise in clinical research, epidemiology, RCD,  
105 validation studies, and development of reporting guidelines.<sup>30 47 48 50-52 63</sup> Manuscript drafts were  
106 circulated among the authors and discussed via teleconferences or e-mail. At the face-to-face  
107 meeting at the ISPE Annual Meeting in Philadelphia, USA (25 August 2019), the authors  
108 provided input on the scope, structure, and terminology of the manuscript. After several  
109 subsequent drafting rounds, the manuscript was circulated to the ISPE membership for feedback  
110 and subsequently revised accordingly. The ISPE membership includes experts from clinical,  
111 academic, and industry sectors from the Americas, Africa, Asia, Oceania, and Europe.

112

## 113 **Terminology**

114 Throughout this paper, by validity we mean internal validity. External validity (generalisability)  
115 of study results is not a topic of this paper, except when portability of algorithms is discussed.

116 Box 1 provides the glossary of the main terms used in this paper. Terminology may differ  
117 slightly from study to study; our aim was to use general, all-inclusive terms, while maintaining  
118 consistency with previous studies and the current guideline documents.<sup>57</sup> Following an earlier  
119 study on RWD validation,<sup>30</sup> we use the term *health state* for any event of interest measured in a

120 study. Depending on study aims, a given health state will take on a role of exposure (e.g.,  
121 medicinal or surgical treatment or device whose safety or effectiveness is of interest),  
122 outcome/endpoint (e.g., treatment adverse event, disease relapse, mortality), or a  
123 covariate/confounder/subgroup (e.g., comorbidity).<sup>30</sup> For example, the RCD record of an  
124 antidiabetics dispensing may define exposure groups “antidiabetic treatment” in one study, the  
125 endpoint “diabetes onset” in another, and a “comorbid diabetes” indicator in the third.

#### 126 BOX 1 HERE

#### 128 **Algorithms in routinely collected data**

129 An *algorithm* is “a completely defined set of operations that will produce the desired outcome”.<sup>64</sup>  
130 In RCD-based studies, algorithms are operational definitions of health states, used to classify  
131 persons with respect to presence/absence and attributes of those health states (e.g., illness and its  
132 severity or treatment and its duration). RCD-based research relies on algorithms to define  
133 inclusion and exclusion criteria; treatment/exposure status, duration, intensity; date of onset of a  
134 disease or a behaviour; and covariates for inclusion in the analyses. An individual’s RCD-based  
135 record, and, by extension, any algorithm based on that record, is a product of the unique chain of  
136 objective and subjective events and practices. These include patient-level events and behaviours,  
137 such as symptoms and care-seeking thresholds (which may vary by sex, age, education,  
138 income);<sup>65 66</sup> characteristics of the health care system (e.g., universal vs. insurance-based  
139 coverage); referral patterns (e.g., whether or not specialist care is general-practitioner-triaged);  
140 clinical aspects (e.g., diagnostic process and treatment decisions); and disease outcome (e.g.,  
141 fatality).<sup>29</sup> RCD-based algorithms depend on record-keeping systems and the specific purpose of  
142 RCD generation (e.g., health administrative data vs. discharge summaries vs. general-practitioner  
143 clinical notes); health sector feeding to a RCD source (e.g., primary vs. hospital care); level of

144 detail allowable by RCD (e.g., type and version of classification used to code diagnoses  
145 (SNOMED, International Classification of Diseases [ICD], 9th Revision vs. ICD-10 vs. -10 vs.  
146 ICD-10-CM [clinical modification] vs. Read codes); and database maintenance routines (e.g.,  
147 frequency of updates, plausibility checks).<sup>30 31</sup> Each link in the chain of events from point of care  
148 to a research data point harbours a potential source of error (Figure 1).<sup>22</sup>

149 **FIGURE 1 HERE**

150  
151 RCD-based algorithms are analogous to diagnostic tests in patient care: both must accurately  
152 classify – “diagnose” – individuals with respect to presence or absence of a given health state at  
153 a given point in time.<sup>30</sup> Similar to the diagnostic process, RCD-based algorithms may be based  
154 on rules (e.g., biomarker value cut-offs), on formal guidelines (e.g., result of a given diagnostic  
155 test), or on clinical features (e.g., a set of signs and symptoms).<sup>67</sup> An algorithm for a given health  
156 state may be as simple as a search for a single diagnostic code during a hospitalisation or as  
157 complex as a decision tree with time-dependent variables measured across multiple data  
158 sources.<sup>31 67 68</sup> A hospital discharge summary carrying the ICD-10 code I21 “Acute myocardial  
159 infarction” is an example of a simple algorithm that can be used to identify patients with acute  
160 myocardial infarction, albeit recorded posthoc. The subcodes may be used to classify patients  
161 according to the affected area of the heart muscle (I21.1 “Acute transmural myocardial infarction  
162 of inferior wall”). Figure 2 shows a comparison of incidence rates of pertussis, identified in three  
163 European countries using various candidate RWD-based algorithms and as defined by the  
164 European surveillance based on a standard case definition.

165 **FIGURE 2 HERE**

166



167 An example of a complex algorithm is a decision-tree based semi-automated procedure that  
168 establishes drug treatment start, duration and continuation based on patients' age, sex,  
169 purchasing history, standard packaging information, and expert input.<sup>69</sup> Machine learning,  
170 increasingly used to aid development of RCD algorithms, is valuable as a screening tool for  
171 identifying candidate algorithm components for subsequent validation.<sup>31 67 68 70 71</sup> The raised  
172 concepts apply regardless of algorithm-generating mechanism, be they static or dynamic,  
173 machine-learning or human expertise based. Box 2 provides examples of RCD-based algorithms  
174 and of studies assessing their validity.

## 175 BOX 2 HERE

176  
177 Researchers planning an RCD-based study should consider whether candidate algorithms are  
178 suitable for use in a given RCD source, based on the data source completeness and algorithm  
179 validity (discussed below). Meaningful development and application of RCD-based algorithms  
180 may require input from clinicians, database administrators, epidemiologists, statisticians, and  
181 data custodians, collectively familiar with local coding, referral and record-keeping practices,  
182 and overall data flow. Figure 3 exemplifies considerations regarding whether and to what extent  
183 to rely on RCD in a given study. For an RCD source to be deemed suitable for a given study  
184 question, the investigator needs to ensure validity and reliability of the algorithms that will be  
185 used to identify study-specific health states. The investigator needs to consider whether a  
186 validated algorithm is available for the RCD source and if so, whether it can be used given the  
187 study question (see Prioritising validity measures); whether it applies to a given population and  
188 the study period (e.g., whether validity estimated for ICD-9 codes holds for ICD-10 codes). If a  
189 validated algorithm is unavailable, should the validation be performed on all or a sample of  
190 health states observed in the study at hand (internal validation) or should a separate – external –

191 validation effort be undertaken. Further decisions relate to the choice of gold/reference standards  
192 (described below), logistics of the validation, and incorporating algorithm validity metrics in  
193 study results and interpretation (Figure 3).

194 [FIGURE 3 HERE](#)

195

### 196 **Completeness of a data source**

197 Completeness of a data source is the proportion of all events of the study-relevant health states in  
198 the study target population captured in that data source.<sup>32</sup> For example, before planning a RCD-  
199 based study of safety of oral non-steroidal anti-inflammatory drugs (NSAIDs) using data on  
200 outpatient dispensings from the Danish National Prescription Registry, we need to understand  
201 the registry's completeness in capturing the NSAIDs dispensings in the population. In Denmark,  
202 the proportion of by-prescription dispensings is 66% for ibuprofen and 100% for diclofenac,<sup>72</sup>  
203 corresponding to the registry's completeness with respect to each drug, as it captures only by-  
204 prescription dispensings.<sup>73</sup> Completeness of an RCD source may be affected by changes in health  
205 policy (e.g., from <100% to 100% for diclofenac dispensings in the Swedish Prescribed Drug  
206 Register<sup>74</sup> as Sweden banned over-the-counter sales following a cardiovascular safety concern<sup>75</sup>).  
207 Completeness is also affected by changes in diagnostic process (e.g., introduction of troponin as  
208 the main diagnostic biomarker of myocardial infarction<sup>65</sup>); or by introduction of screening (e.g.,  
209 screening for colorectal cancer will lead to an increase over time in the number of persons with a  
210 diagnosed colorectal cancer recorded in cancer registries). When planning RCD-based studies,  
211 completeness of candidate data sources needs to be assessed with respect to health states  
212 representing all study variables – exposures, outcomes, covariates, and subgroup identifiers.

213

## 214 **Measures of algorithm validity**

215 The RCD algorithms/diagnostic test analogy extends to the measures of validity.<sup>30 76</sup> These  
216 measures, covered in standard texts and tutorial articles,<sup>20 21 77 78</sup> are summarised in Box 3, for  
217 reference. Algorithm validity measures include sensitivity and specificity, and their derivatives  
218 (receiver operating characteristic (ROC)-curve, area under the curve (AUC)<sup>79</sup>, diagnostic odds  
219 ratio (DOR)<sup>80</sup>); indicators of performance (positive and negative predictive values (PPV and  
220 NPV)); and measures of agreement (kappa statistic<sup>77 81</sup>). (Other sources have used, the term ‘bias  
221 parameters’<sup>56</sup> or validity indices<sup>78</sup> to collectively describe the measures of algorithm validity).  
222 When a health state is rare, both specificity and NPV of an algorithm are less useful because they  
223 are expected to be close to 100%. Sensitivity and specificity of algorithms are generally,  
224 although not always, independent of prevalence of the health state in the underlying population  
225 (equivalent of pre-test probability for a diagnostic test). Crucially, both PPV (equivalent of post-  
226 test probability for a diagnostic test) and NPV depend on the prevalence of the health state of  
227 interest in the underlying population (PPV directly proportional, NPV inversely proportional to  
228 the prevalence).<sup>79</sup> A PPV estimated in a validation cohort in which the health state of interest is  
229 more prevalent than in the underlying source population, will be overly ‘optimistic’, meaning  
230 that the PPV applied to the data source will be lower. Simple calculations allow derivation of  
231 PPV and NPV from known sensitivity, specificity, and prevalence of the health state in a  
232 population.<sup>30 56</sup> Many validation studies are not designed to estimate all validity measures.  
233 Sensitivity and PPV are the most commonly reported measures<sup>30 82</sup> because estimation of  
234 specificity and NPV requires identification of a representative group of true-negative individuals,  
235 which may be challenging with routinely collected data due to requirement of exhausting labour.

236 **BOX 3 HERE**

237

238 **Gold standard and reference standard**

239 Sensitivity, specificity, PPV or NPV of an algorithm are estimated against a gold standard.

240 Conceptually, a gold standard is a method that classifies individuals with respect to presence and  
241 absence of a given health state without errors, i.e., there are no false-positives and no false-  
242 negatives. As a true gold standard rarely exists, in practice validation studies rely on a “reference  
243 standard” or an “alloyed gold standard”, i.e., a classification method inferior to the theoretical  
244 (but often non-existent) gold standard, but superior to the RCD algorithm being validated (Box 2  
245 provides examples of gold and reference standards).<sup>77</sup> For the health state under validation, the  
246 reference standard provides case and non-case definitions, against which to compute the  
247 measures of validity and performance (Box 3). Figure 4 shows relations between the gold  
248 standard, the reference standard, and the algorithm. Setoguchi *et al.*, in a validation of  
249 haematological malignancies defined from health administrative data against a cancer registry as  
250 the reference standard illustrated dependence of validity measures on the completeness of the  
251 reference standard.<sup>83</sup>

252 **FIGURE 4 HERE**

253  
254 Medical records are commonly used as a gold/reference standard in validation studies of RCD-  
255 based algorithms. When planning a validation study based on review of medical records, points  
256 to consider include:

- 257 • Look-back and look-forward period for record review relative to the date of the potential  
258 health state identified by the RCD algorithm, e.g., only information in 3 months before  
259 and after the RCD-recorded event will be considered in validation, to allow for data flow  
260 artefacts, e.g., delay of diagnosis recording relative to a true diagnosis date in  
261 administrative data;

- 262 • What types of key records should be sought to confirm the health state (e.g., glycated  
263 haemoglobin for diabetes, liver enzyme measurements for hepatotoxicity, bone  
264 scintigraphy for bone metastases etc.);
- 265 • Whether only records from specific contacts are of interest for validation, e.g., at specific  
266 clinical departments (see the osteonecrosis of the jaw example, in the description to  
267 Figure 3);
- 268 • Whether to flag prevalent vs. incident status of the health events under validation;
- 269 • Whether chart abstraction and case adjudication processes should or can be combined –  
270 e.g., conducted by an expert adjudicator or should a nurse conduct the abstraction and  
271 should the health state be subsequently adjudicated by an expert? De-coupling chart  
272 abstraction and adjudication may be necessary, for example, if adjudicators need to be  
273 blinded to certain information, such as treatment status in validating endpoints in safety  
274 assessments;
- 275 • How to quantify to resolve disagreements and quantify interrater variability, if chart  
276 abstraction/adjudication is done by more than one person.<sup>30 81</sup>
- 277 • Design of data collection instruments such as chart abstraction forms and software: in  
278 what order should the information be abstracted?
- 279 • What are the mechanisms for correcting entries if needed or unclear? Is it possible and if  
280 so, for how long, to revisit the abstraction source to edit information?
- 281 • Staff skills and potential need for study-specific training.

282 A pilot medical record review ahead of a large-scale effort will inform the process and help use  
283 resources efficiently.

284

285 Medical records must be viewed especially critically when used as the reference standard for  
286 adjudicating RCD-identified endpoints in safety studies of medicines, where the initial safety  
287 data originate from preapproval trials. Justifiably, a researcher may want to have the same case  
288 definitions for trial- and RCD-based events. However, the components of case definition  
289 available for closely monitored participants in clinical trials may be under-recorded or not  
290 recorded in the course of routine clinical care, especially if events are asymptomatic or  
291 insufficiently severe to trigger diagnostic activity or care seeking. For example, if QT  
292 prolongation is a safety concern in a prospective trial, patients' electrocardiogram (ECG) would  
293 be taken at baseline and monitored during the entire follow-up, regardless of severity or  
294 symptoms. In routine clinical practice, physicians will detect and document ECG results  
295 predominantly on symptomatic care-seeking patients, and generally not have information on the  
296 equivalent of the baseline ECG to gauge whether QT- prolongation is pre-existing or newly  
297 occurring (Berkson's bias<sup>84</sup>). In this case, the medical record is likely to have low sensitivity, and  
298 trial case definitions may need to be modified to the type of information typically available from  
299 medical records.<sup>85</sup> Furthermore, if physicians are on a higher alert to adverse events with new  
300 treatments, differential misclassification may ensue.

301  
302 Another validation option in the absence of a suitable gold/reference standard is use of latent  
303 class modelling, whereby different RCD-based algorithms for a given health state are assigned  
304 probabilities of correctly classifying patients with respect to the health state status, through data-  
305 driven modelling. Prosser and colleagues used latent class modelling to validate three RCD-  
306 based algorithms to identify treated asthma in administrative data in Canada, with each algorithm  
307 based on a combination of frequency of physician visits and hospitalisation carrying qualifying  
308 diagnostic codes.<sup>86</sup> In addition to medical records, patient or physician interviews, or another

309 database are examples of gold/reference standards used in RCD validation studies (examples in  
310 Box 2).<sup>30 31 77</sup>

311

### 312 **What measures of validity can be estimated in a validation study?**

313 Types of validity measures estimable in a given validation study depend on the strategy used to  
314 assemble the study population, which typically falls into three scenarios.<sup>56</sup> In the first scenario,  
315 study population is assembled based on the RCD-algorithm defined health status – algorithm-  
316 positive and algorithm-negative are sampled from an RCD source and their RCD-record is  
317 compared with the gold-standard based case definition (e.g., medical records). This scenario  
318 allows estimation of PPV and NPV. In the second scenario, the study population consists of  
319 persons who are health-status positive and health-status negative according to the gold standard.  
320 Their RCD records are searched for the elements of the proposed RCD-based algorithm. This  
321 scenario enables estimation of sensitivity and specificity. In the third scenario, study population  
322 is assembled independent of either RCD- or reference-standard based status and allows  
323 estimation of sensitivity, specificity, PPV, and NPV (see Fox et al for a detailed discussion<sup>56</sup>). If  
324 RCD algorithms are assessed in data sources, none of which is superior, kappa statistic may be  
325 the only estimable metric. Bollaerts et al. derived analytical expressions connecting the observed  
326 (algorithm-based) prevalence of a health state and its four validity measures (sensitivity,  
327 specificity, PPV, and NPV). These expressions, and an associated web-based application allows  
328 derivation of unknown validity measures based on the observed prevalence and any two other  
329 measures.<sup>78</sup> A worked example of this application is provided in a study of Morkem et al.<sup>87</sup>

330

331 **Prioritising validity measures**

332 Ideally, we want to maximise all measures of algorithm validity, however, practically we must  
333 accept trade-offs between the measures: increasing algorithm sensitivity (allowing it to capturing  
334 more true-positive cases) is accompanied by a decreasing specificity (by capturing more false-  
335 positive cases) and vice versa.<sup>64 83</sup> The impact of imperfect validity measures on study findings  
336 depends on the role of an RCD-defined health state in a given analysis: inclusion/exclusion  
337 criterion, exposure, outcome, covariate, as illustrated in several examples, below.

338 *RCD algorithms for eligibility criteria*

339 Consider an observational study of comparative effectiveness and safety of oral anticoagulants  
340 among patients with atrial fibrillation (AF) as an example. Suppose we want to conduct this  
341 study using RCD from Scandinavian nationwide registries, whereby data on diagnoses originate  
342 from hospital contacts, and data on treatment, from outpatient dispensings.<sup>88</sup> In addition to AF,  
343 indications of oral anticoagulants include post-arthroplasty thromboprophylaxis and treatment of  
344 recurrent venous thromboembolism.<sup>66</sup> Patients' characteristics, treatment dosing and duration  
345 vary by indication. Therefore, we wish to ensure that the study population includes only patients  
346 treated for AF. As routine dispensing records rarely contain information on indication or  
347 prescribed dose, the RCD-based eligibility criteria must indirectly identify patients treated for  
348 AF. A drug utilisation study showed that AF is recorded among nearly 80% of patients who  
349 initiate anticoagulant.<sup>66</sup> Thus, we could first identify drug initiators and exclude patients with a  
350 record of arthroplasty or venous thromboembolism and assume that all remaining patients have  
351 AF.<sup>89</sup> A limitation of this approach is that some remaining patients will not have AF. An  
352 approach that alleviates, but not fully resolves this limitation, would require a hospital diagnosis  
353 of AF as an inclusion criterion: this will ensure that all patients have AF, but not that they are  
354 treated for it with anticoagulants. To further increase specificity of identifying the eligible



355 population, we could exclude from that population patients with a hospital diagnosis of venous  
356 thromboembolism or a record of arthroplasty. Although this exclusion criterion will remove  
357 some eligible patients, we may decide to pay this price for the sake of obtaining a valid  
358 estimation of comparative effectiveness and safety. The measures we take to increase the  
359 specificity of the eligibility criteria algorithm will come at a price of potential selection bias, if  
360 patients with hospital-diagnosed AF systematically differ from patients with AF initiating  
361 anticoagulants seen in all medical care settings.<sup>90</sup> Patients seen in hospitals are likely to represent  
362 the most severe spectrum of the disease; this bias would be avoided by the sensitive eligibility  
363 algorithm, described above. Selection bias due to inclusion of patients with other indications  
364 could be avoided by keeping such patients in the study population and conducting a stratified  
365 analysis. With rare diseases, we may prioritise eligibility criteria algorithms with high PPV and  
366 specificity, even at a price of reduced precision, to avoid “contamination” of the study cohort  
367 with false-positives; alternatively one may opt for a hybrid approach by identifying potentially  
368 eligible patients based on RCD algorithms and subsequently eliminate false-positive cases by  
369 reviewing medical records.<sup>64</sup>

370

### 371 *RCD algorithms for exposure*

372 RCD-based algorithms establishing patients’ treatment/drug exposure rely on routine records of  
373 prescriptions (e.g., in EHR data), dispensing (e.g., in pharmacy administrative data), or  
374 administrations (e.g., in hospitalization data) to measure treatment initiation, duration, dosage,  
375 discontinuation. In identifying new users of a drug, the length of the washout (treatment-free)  
376 period may produce severe misclassification between the status of new and prevalent user, which  
377 may lead to underestimation of risks of early side effects.<sup>91</sup> Incorrect allocation of on-/off-  
378 treatment person-time will likewise bias studies of comparative effectiveness. Although routine

379 prescription and dispensing records are generally considered to be high quality, their  
380 correspondence to the true treatment status/adherence is difficult to validate. Therefore RWE on  
381 medication safety and effectiveness often include multiple sensitivity analyses to assess  
382 robustness of studies against different RCD-based algorithms and their assumptions.<sup>51</sup>  
383 Information on treatment absence may be of interest in studies aiming to explore whether some  
384 patients who may benefit from treatment remain untreated. If an RCD source has incomplete  
385 treatment records, there is likely a mixture of true and false negatives amongst the cohort of  
386 patients without an RCD-based treatment record. If under-treatment is of interest, the NPV of the  
387 RCD-based treatment record should be prioritised. In addition, whether the effect of exposure is  
388 prolonged for a predetermined period after one-time exposure to a drug or limited during the  
389 period when a patient is actually exposed to a drug should be determined considering the  
390 characteristics of the drug for each outcome health state.

391

#### 392 *RCD algorithms for outcomes*

393 When RCD-based algorithms are used to define events to estimate their absolute risk or risk  
394 difference, such as adverse treatment effects, an algorithm with low sensitivity will lead to  
395 underestimation of the true risk, potentially providing false reassurance. However, the  
396 corresponding relative effect (risk ratio) will be unbiased no matter how low the sensitivity,  
397 provided a near-perfect specificity. This consideration is only true for a binary outcome. The  
398 same consideration will apply in studies of disease epidemiology for the purposes of planning  
399 health services and in studies comparing disease populations with different severity or among  
400 different medical-care settings. For example, an RCD-algorithm based study on epidemiology of  
401 diabetes when only secondary-care diagnoses are available (e.g., in Sweden or Denmark<sup>92,93</sup>) will  
402 underestimate the incidence and prevalence of diabetes (diagnostic delay) and overestimate the

403 average disease severity. Moreover, the time of disease onset will be incorrect for patients with  
404 diabetes seen in hospitals after being diagnosed in primary care. In the case of diabetes or other  
405 diseases with specific treatments, sensitivity of identification can be increased by use of  
406 treatment proxies or laboratory values.

407

#### 408 *RCD algorithms for confounders*

409 Accurate RCD-based algorithms are important for confounding control in observational studies.  
410 Poorly measured confounding variables perform poorly in controlling for confounding and yield  
411 biased results. For example, obesity is an important confounder in many associations, yet many  
412 RCD sources do not capture obesity or body mass index with sufficient accuracy or  
413 completeness. For example, the ICD-10 based algorithm for overweight or obesity in the Danish  
414 National Patient Registry has a PPV close to 90%, but the registry is merely 11% complete in  
415 capturing patients with obesity.<sup>94</sup> If obesity is an important confounder, analyses adjusted for  
416 obesity using that algorithm will remain largely confounded in the analyses based on that RWD  
417 source, and researcher should consider using an external adjustment to evaluate whether some or  
418 all of the observed association is likely to be explained by the unmeasured confounding.<sup>95</sup>

419

420 In summary, decisions about prioritising validity measures imply balancing the risks and benefits  
421 of including false-positives vs. false-negatives that should be considered given the study aims.

422

#### 423 **Portability of algorithms**

424 Variability of algorithm validity within and between RCD sources is likely to depend on several  
425 factors including patients' age (paediatric vs adult population), calendar time (changes of  
426 guidelines, diagnostics, classifications, referral patterns), disease severity (hospitalisation

427 patterns), health care funding scheme (single-payer vs. multiple-payers), accessibility (universal  
428 vs. choice-based), health care sector contributing to RCD (general practice vs. hospital care vs.  
429 both), and coding practices rooted from insurance system and policy, and changes in coding  
430 dictionaries and supporting IT infrastructure.<sup>67</sup> Therefore, algorithms validated in one setting  
431 (country, RCD source) cannot be assumed to have the same validity in a different setting  
432 (country, RCD source). As illustrated in the pertussis infection example (Figure 2), in multi-  
433 database studies, assuming that an algorithm is transportable may not be tenable, and must be  
434 critically assessed given the differences of the underlying health care systems, guidelines, and  
435 mechanism of RCD record generation.<sup>96-100</sup> In preparation for studies of effectiveness of the  
436 vaccines against the severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) in Europe,  
437 background incidence rates of multiple potential adverse events have been estimated in using  
438 multiple and heterogeneous RWD sources in several European countries.<sup>101</sup> The variability of the  
439 observed background rates estimated based on common code lists underscores the importance of  
440 algorithm calibration and benchmarking using available evidence, and features of a given RWD  
441 source. These background rates should remain country-specific and are meant for detecting  
442 signals for subsequent formal evaluation. Differences in observed rates across multiple  
443 databases, in a process similar to triangulation,<sup>102</sup> may be used to glean the direction and the  
444 magnitude of measurement error associated with a given RWD source.

445

446 On the other hand, because of cost and effort associated with conducting a validation study,  
447 before embarking on one, researchers should search the relevant literature on whether a validated  
448 algorithms are already available for their given study aim (Figure 3). At the same time, with  
449 strategic investments in interoperability, creation of searchable and updated algorithm libraries,  
450 alignments and partnership consortia between the countries and relevant data sources the burden

451 associated with conducting validation studies can be dramatically reduced. Approaches to  
452 ensuring applicability range from calibration-refinement exercises<sup>98</sup> to full-scale re-validation in  
453 the new setting (Figure 3).<sup>103</sup>

454

455 Another aspect of algorithm portability is whether incident and/or prevalent events are of  
456 interest. This distinction is important in safety studies of association between a treatment and a  
457 side effect that is a chronic disease, whereby only incident events, i.e., events with onset  
458 following the medication initiation should be counted as a case. Such case definition would  
459 require the gold/reference standard to afford a sufficiently long lookback period to rule out  
460 prevalent conditions. For example, some patients with inflammatory bowel disease may have  
461 more than 8 years without disease-related records in administrative data,<sup>104</sup> suggesting that a  
462 lookback period of 1 year, commonly used in studies based on insurance claims, may not be  
463 sufficient to exclude prevalent cases.

464

#### 465 **Size of a validation study**

466 If the aim of a validation study is to estimate a dichotomous measure of validity (Box 3), the  
467 number of the validated records determines precision of the resulting estimates. When the target  
468 validity measure is dichotomous, size of the study becomes the consideration of precision of a  
469 dichotomous proportion.<sup>79 105</sup> Precision considerations must be balanced against considerations  
470 of available human and financial resources and must account for the expected nonresponse rate  
471 (e.g., if medical records are unavailable). Further considerations include whether variability of  
472 the validity measures is expected to vary by subgroups, such as calendar period (e.g., record  
473 recency may determine classification version, diagnostic validity, treatment guidelines, referral  
474 or reimbursement practices; classifications); disease subtype or severity (e.g., subcodes); priority

475 and setting of the diagnosis (e.g., primary vs secondary, inpatient vs outpatient); or patients’  
476 characteristics (e.g., sex or comorbidity). If estimation of the validity measures is important in a  
477 given subgroup, the number of validated records should allow for the desired subgroup precision.  
478 For planning validation studies of complex algorithms, researchers should consider whether  
479 validity measures would be estimated for the overall algorithm or for each algorithm component  
480 (e.g., a group of diagnostic codes vs. a single diagnostic code), with potential subsequent  
481 algorithm refinement, such as removal of diagnostic codes with low validity.<sup>106</sup> The study size  
482 should include sufficient number of observations with a given algorithm component. For internal  
483 validation, the proportion of records selected for validation may depend on whether the condition  
484 is rare or common.<sup>107</sup> An adaptive validation study design using a Bayesian approach has been  
485 recently proposed for internal validation studies aiming to estimate PPV and NPV. The  
486 validation information accrues as the study is being conducted. The design defines decision rules  
487 for when the amount of validation accrued is sufficient to provide reliable inference, whereupon  
488 validation can be stopped.<sup>108</sup> The increasing availability of electronic medical charts in  
489 combination with machine learning and text mining may help undertake validation studies on  
490 larger scale in the future.

491

#### 492 **Interpreting study results in the light of information bias**

493 Researchers may be tempted to dismiss misclassification bias (e.g., of exposure status) on the  
494 grounds that the obtained estimates are “conservative” i.e., that the true association is speculated  
495 to be (even) stronger than the one observed. In a safety study of a treatment, procedure, or  
496 device, nondifferential misclassification of exposure may mask an important safety signal. A true  
497 association masked by misclassification may provide false assurance about treatment safety,  
498 while a spurious association when none exists may cause unnecessary treatment withdrawal.

499 More generally the “nondifferential misclassification mantra”<sup>109</sup> refers to a statistical expectation  
500 under a narrow set of conditions (in which the exposure has an effect, the health state is  
501 dichotomous, and misclassification is non-differential).<sup>24 25</sup> For polytomous variables created by  
502 categorising continuous variables, the direction of measurement error depends on a category.<sup>110</sup>  
503 Finally, in RCD, all variables with rare exceptions, are measured with error, and estimation of  
504 joint effect of errors in all study variables cannot be assumed to be in a given direction without  
505 quantitative assessment.<sup>2828 109 111</sup>

506  
507 If a validation study reveals that an RCD-based algorithm’s validity measures are unsuitable  
508 given study aims and its role in the analysis, an alternative approach to identifying that health  
509 state should be considered (e.g., de-novo data collection) (Figure 3). If one proceeds with the use  
510 of an RCD-based algorithm, the potential impact of information bias on the interpretation of  
511 study findings should be considered and, ideally, quantified. Methods and software for such  
512 quantification<sup>82 109 112-116</sup> range from simple calculations examining the impact of  
513 misclassification of one variable at a time<sup>113</sup> to simultaneous assessment of the impact of errors  
514 of multiple study variables (bias analysis)<sup>82 112</sup> to using imputation methods to correct for  
515 misclassification.<sup>117</sup> The latter methods can be used to estimate the range of study results that  
516 could be observed under plausible ranges of validity measures or allow re-computation of what  
517 study result would have been expected had the measurement been perfect or better than the one  
518 available.<sup>109</sup> If sensitivity and specificity parameters are unknown, these methods can be used to  
519 evaluate the impact of a plausible range of different sets of sensitivities and specificities on study  
520 results. Predictive values can be used to conduct quantitative bias analysis to correct for  
521 misclassification of outcome variables and a research has shown that imperfect predictive values  
522 impact study results ranging from negligible to significant extent, producing incorrect

523 conclusions (e.g., misclassified data show presence of an exposure – outcome association, but  
524 corrected data do not).<sup>118</sup> Other techniques for correcting study results for measurement error  
525 include hierarchical semi-Bayes methods,<sup>114</sup> or bootstrap imputation.<sup>115</sup>

526

527 Interpretation of study findings warrants quantifying the amount and the potential impact of  
528 information bias due to imperfect algorithm validity and especially discuss and quantify  
529 plausible scenarios that are likely to invalidate study conclusions. Despite extensive theoretical  
530 developments and practical applications and examples<sup>109 111 119</sup> measurement error, especially  
531 differential and non-independent errors, remain neglected in interpreting study findings, and  
532 formal analysis of potential impact of measurement error in interpreting results remains an  
533 exception, rather than a rule.<sup>28</sup>

534

### 535 **Summary**

536 As the RCD landscape evolves, ensuring, quantifying, and updating algorithm validity  
537 information should become a standard maintenance task of all RCD-based research efforts.  
538 Information on algorithm validity needs to reflect up-to-date clinical knowledge, diagnostic  
539 accuracy, clinical guidelines, regulatory practice, insurance policies, coding conventions,  
540 granularity and completeness of recording, and data flow.<sup>34 120-124</sup> As a starting point for  
541 developing protocols and implementation of validation studies of RCD-based algorithms, we  
542 recommend use of the published checklist of reporting criteria for studies validating RCD-based  
543 algorithms, spelling out practical elements of conducting the validation study.<sup>30</sup> Formalising the  
544 conduct of validation studies will help improve quality of evidence of RCD-based research and  
545 foster trust between RCD researchers and research consumers, such as policy-makers and  
546 clinicians. We posit that if RCD-based research delivers valid evidence, accurate clinical and



547 administrative record-keeping at the point of health care delivery will be viewed by record-  
548 keepers as an investment in the quality of care rather than a bureaucratic task. In the words of the  
549 FDA Commissioner, “To enable greater adoption of RWE in clinical and regulatory decisions,  
550 we’ll need to work with the healthcare system to change the way clinical information is  
551 collected. Ideally, we’d like to have a system where providers have the right incentives to enter  
552 clinically relevant information into EMRs at the point of care”.<sup>125</sup> Despite their advantages, RCD  
553 cannot be a default solution to any research question. For example, for studies of rare diseases  
554 with centralised clinical management at a few treatment centres, a study that uses primary data  
555 collection, alone or in combination with RCD, may be more efficient and valid than a purely  
556 RCD-based study, as using RCD may be like seeking a needle in a haystack because of the  
557 dependence of some validity measures on diseases prevalence. Limited information on  
558 confounders is another potentially important limitation of RCD.

559  
560 The principles described here are broadly applicable for all studies relying on routinely collected  
561 health data, including observational studies of diseases, benefit-risk assessments of interventions,  
562 as well as pragmatic randomised trials harnessing routinely collected health data to measure  
563 patients’ characteristics, risk factors, and/or outcomes in routine clinical practice. As RCD are  
564 being increasingly used to make important clinical, regulatory and policy decisions about drugs,  
565 procedures, devices and other healthcare interventions, we hope to have raised awareness about  
566 the importance of quantifying measurement error around RCD-based algorithms when  
567 generating real-world evidence.

568

569 **Funding**

570 Funding to support this manuscript development was provided by the International Society for  
571 Pharmacoepidemiology (ISPE). A draft of the manuscript was made available for review to the  
572 ISPE members and subsequently revised in response to that review.

573 **Acknowledgements**

574 This manuscript is endorsed by the International Society for Pharmacoepidemiology (ISPE). We  
575 gratefully acknowledge review of the manuscript by ISPE members. We thank Ms. Helle Vester,  
576 MSc, for excellent administrative support.

577

578 **Disclosures**

579 Dr. Ehrenstein is a salaried employee of Aarhus University, which receives institutional research  
580 funding from various pharmaceutical companies and regulatory agencies, administered by  
581 Aarhus University.

582 Dr. Hellfritsch has nothing to disclose.

583 Dr. Kahlert is a salaried employee of Aarhus University; institutional research funding from  
584 various pharmaceutical companies and regulatory agencies to and administered by Aarhus  
585 University.

586 Dr. Langan was supported by a Wellcome Senior Research Fellowship in Clinical Science  
587 (205039/Z/16/Z). Dr. Langan was also supported by Health Data Research UK (grant No.  
588 LOND1), which is funded by the UK Medical Research Council, Engineering and Physical  
589 Sciences Research Council, Economic and Social Research Council, Department of Health and  
590 Social Care (England), Chief Scientist Office of the Scottish Government Health and Social Care  
591 Directorates, Health and Social Care Research and Development Division (Welsh Government),  
592 Public Health Agency (Northern Ireland), British Heart Foundation and Wellcome Trust.

593 Dr. Urushihara is a salaried employee of Keio University, and has received research support  
594 from CAC Croit Corporation, Shionogi & Co., Ltd., and Senju Pharmaceutical Co., Ltd.

595 Dr. Marinac-Dabic has nothing to disclose.

596 Dr. Lund receives research support from AbbVie, Inc. and her spouse is a full-time, paid  
597 employee of GlaxoSmithKline.

598 Dr. Sørensen is a salaried employee of Aarhus University; institutional research funding from  
599 various pharmaceutical companies and regulatory agencies to and administered by Aarhus  
600 University.

601 Dr. Benchimol was supported by a New Investigator Award from the Canadian Institutes of  
602 Health Research, Canadian Association of Gastroenterology and Crohn's and Colitis Canada. Dr.  
603 Benchimol was also supported by the Career Enhancement Program of the Canadian Child  
604 Health Clinician Scientist Program. Dr. Benchimol has acted as a legal consultant for Hoffman  
605 La-Roche Limited and Peabody & Arnold LLP, has received consulting fees from McKesson  
606 Canada, and receives research funding from Crohn's and Colitis Canada for matters unrelated to  
607 this article.

608

**Box 1. Glossary of key terms**

<b>Term</b>	<b>Definition/elaboration</b>
Algorithm	An operational definition of a <i>health state</i> , used to classify persons with respect to presence/absence and attributes of those health states. An algorithm may be as simple as a single diagnostic code or as complex as a set of decision rules. Implicitly, persons not meeting algorithm criteria are classified in the analysis as not having a given health state.
Alloyed gold standard	See <i>Reference standard</i>
Area under the curve (AUC)	Area under the ROC curve, used as a single measure of the discriminative performance of an algorithm based on cut-off values of a continuous variable. See Box 3 for details.
Completeness	Completeness of a data source is the proportion of all events of the study-relevant health states in the study target population captured in that data source.
Confounder	A confounder can be conceptualized as a common cause of the exposure status (eg treatment decisions) and the outcome development. Thus, a confounder cannot be a causal intermediate between the exposure and the outcome. In a given study, a characteristic that fulfils the above criteria will exert confounding if it is 1) unequally distributed across exposure categories AND b) predicts the study outcome. Confounder is a special case of a covariate.
Covariate	A characteristic of the members of the study population besides the exposure and the outcome that may be of interest. This term is often used for predictors of the outcome of interest. By contrast with a confounder, covariate is need not be associated the exposure of interest.
External validation	Validation of events of a health state identified with an RWD-based algorithm in a population external to a given study.
Diagnostic odds ratio (DOR)	The ratio of the odds of algorithm positivity in those with a health state of interest to the odds of algorithm positivity in those without the health state of interest.
Gold standard	In an ideal sense, a method that classifies individuals with respect to a health state of interest without errors. True gold standards rarely exist (see <i>Reference standard</i> ), as no method of measurement is error-free. In practice, the quality of gold standards may be time-dependent. For example, new diagnostic tests may, paradoxically, perform better than methods previously considered gold standards.
Health state	A generic term used in this paper for any event of interest measured in a study. Depending on study aims, the same health state may play a role of exposure, outcome, covariate, or a subgroup indicator.
Information bias	Discrepancy between the true value of a given trait or characteristic (height, weight, disease status) and its measured value. Errors may be inherent in a measurement instrument (eg biased scale used to measure weight) or be human (eg data entry errors). <i>Information bias</i> is often used synonymously with <i>Measurement error</i> and <i>Misclassification</i> .
Internal validation	Validation of events of a health state identified with an RWD-based algorithm in a given study.
Measurement error	See Information bias.
Misclassification	Incorrect classification by an <i>algorithm</i> of an individual's health state, resulting from <i>Measurement error</i> (eg, a patient with a body mass index in a

	normal range may be classified as obese as a result of a measurement error of their height or weight).
Negative predictive value (NPV)	Proportion of patients who truly negative for a health state among all those who are classified as negative by an algorithm. See Box 3 for details.
Positive predictive value (PPV)	Proportion of patients who truly positive for a health state among all those who are classified as positive by an algorithm. See Box 3 for details.
Real-world data (RWD)	Data that originate in the course of routine clinical practice, usually by contrast with data originating in phase I-III clinical trials.
Real-world evidence (RWE)	Evidence generated using RWD.
Reference standard	A method or a data source that is expected to classify individuals with respect to a health state of interest better than the algorithm being validated. Alternative term: alloyed gold standard.
Receiver operating characteristic (ROC) curve	For algorithms that classify patients into health states (eg., presence/absence/severity) a plot of sensitivity (or true positive proportion) against 1-specificity (or false-positive proportion), Box 2.
Routinely collected health data (RCD)	Data that accrue routinely as a by-product of health care delivery or administration of the health care system. Routinely collected data may also include non-health characteristics such as income, education, and employment, if important in studying health states. RCD are a type of RWD.
Sensitivity	Proportion of persons with a given health state according to gold/reference standard who are classified as positive by an (RWD) algorithm. See Box 2 for details.
Specificity	Proportion of persons without a given health state according to gold/reference standard who are classified as negative by an (RWD) algorithm. See Box 2 for details.
Exposure	A health state that, in a given study, explicitly or implicitly, defines patients' characteristic that is equivalent to treatment allocation in randomized trials. In an RWD -based study, exposure status may correspond to treatment status (medicinal, surgical, device), but may also correspond to a patient's characteristic (comorbidity, age, sex, socioeconomic status). Conceptually, exposure is the independent variable.
Outcome	A health state that, in a given study, plays a role of an outcome of interest. In an RWD -based study, it could be any outcome of treatment (benefit or risk), disease recurrence (e.g., recurrent malignancy, readmission for a myocardial infarction), or death. Conceptually, exposure is the dependent variable.
Validation cohort	Study population used to evaluate measures of algorithm validity.
Validity (internal)	Correspondence between an estimate and the true value of a parameter. In epidemiologic studies, a valid (unbiased) estimate may estimate occurrence of a health event, or an association between exposure and outcome.

**Box 2. Examples of RCD-based algorithms and validation studies of RCD-based algorithms conducted with different types of gold/reference standard**

Health state or event being validated	Country	RCD source(s)	RCD-based algorithm	Gold/reference standard	Validation study
In-hospital chemotherapy treatment	Denmark	Danish National Patient Registry	Agent-specific treatment codes	In-hospital pharmacy production system and chart review	By comparing RCD data with the reference standards, and vice-versa, for 50 randomly selected patients with colorectal cancer having nodal involvement, the kappa, sensitivity, specificity, positive and negative predictive values were estimated for the RCD-based algorithm for receipt of any chemotherapy and for specific treatments (Lund et al 2013 <sup>126</sup> )
Start, duration and end of prescribed medicine	Finland	Finnish Prescription Registry	Different approaches to estimation of drug use periods based on prescription data. Three fixed methods were included and one data-driven method (PRE2DUP)	Expert review of purchase history	For 200 randomly selected purchase histories of warfarin, bisoprolol, simvastatin, risperidone and mirtazapine in patients with Alzheimer's disease, two experts reviewed purchase histories and evaluated the RCD-based algorithms with regards to their ability to provide purchases and duration. (Tanskanen A et al 2017 <sup>127</sup> )
Cardiac interventions	Denmark	Danish National Patient Registry	Surgical and procedure codes according to the NSCP classification	Chart review	Charts from randomly sampled patients identified using the RCD-based algorithm were reviewed and positive predictive values for each intervention were calculated (Adelborg et al 2016 <sup>128</sup> ).
Osteonecrosis of the jaw	Denmark	Danish National Patient Registry	ICD-10 codes recorded at departments of oral and maxillofacial surgery in cancer patients	Chart review and an independent sample of confirmed cases	By assessing the occurrence of osteonecrosis of the jaw according to the chart among patients with a positive algorithm, the positive predictive value of the RCD-based algorithm was obtained. Sensitivity of the algorithm was estimated as the proportion of confirmed cases with a positive algorithm (Ehrenstein et al 2015 <sup>129</sup> )
Colorectal cancer recurrence	Denmark	Danish National Patient Registry and the Danish Pathology Registry	Combination of codes compatible with metastatic disease, chemotherapy administration and cancer recurrence in patients previously diagnosed with non-metastatic colorectal cancer	Independent sample of patients with known recurrence status	By comparing the algorithm-based and the true recurrence status, sensitivity, specificity, positive and negative predictive values were calculated (Lash et al 2015 <sup>130</sup> ).
Frailty in the elderly	United States	Administrative claims	Medicare claims-based algorithm of dependency in activities of daily living (or dependency) developed as a proxy for frailty	A reference standard measure of phenotypic frailty	We examined the discriminative ability of the claims-based algorithm to predict phenotypic frailty using a receiver operating characteristic (ROC) curve and estimated the sensitivity, specificity, positive predictive value (PPV), and negative predictive value (NPV) of the claims-based algorithm relative to the reference standard using a range of dichotomous cut-off values (Cuthbertson et al 2018 <sup>131</sup> ).

Hepatic encephalopathy	United States	Administrative claims	ICD-10 codes and treatment proxies after switch to ICD-10 in 2015 rendered algorithms validated using ICD-9 invalid	Independent prospectively acquired sample containing both true positive and true negative cases, with case status established by diagnostic criteria.	Data from the independent sample were cross-tabulated with various ICD-10 and treatment-proxy based algorithms, computing sensitivity, specificity, PPV, NPV, and area under the curve (AUC) (Tapper et al 2019 <sup>132</sup> ).
Malignant melanoma	Denmark	Danish Cancer Registry and Danish Pathology Registry	ICD-10 code of malignant melanoma and ICDO-10 for topography and morphology as registered in the Cancer Registry	Another RCD source: Danish Pathology Registry	In a sample identified and characterized by the algorithm, the positive predictive values of malignant melanoma diagnosis, topography and morphology were calculated by comparing with data registered in the Pathology registry. Using a sample from the Pathology Registry, the completeness and the sensitivity of the melanoma diagnosis in the Cancer Registry was estimated (Pedersen et al 2018 <sup>133</sup> ).
Childhood-onset inflammatory bowel disease	Canada	Health administrative databases in Ontario, Canada collectively including hospital discharge data, billing claims for all physician services and demographic data including region of residence).	Various combinations of physician office and procedure billings with hospital records	(1) Another RCD source: clinical database of patients with childhood-onset inflammatory bowel disease (2) Chart review and an independent sample of confirmed cases and controls	A population-based clinical database of patients with IBD diagnosed aged <15 years was used to define cases, and patient information was linked to health administrative data to compare the sample of patients with known childhood-onset inflammatory bowel disease was used to test the diagnostic accuracy (sensitivity, specificity, positive predictive value and negative predictive value) of various patterns of healthcare use as registered in health administrative registries. The most accurate algorithm was further validated with chart data of children medical practices. The algorithm was validated in a sample of patients with IBD diagnosed <18 years derived from chart reviews at various medical practices. (Benchimol et al 2009 <sup>134</sup> ).
Use of oral anticoagulants due to atrial fibrillation	France	French national health insurance system database linked to the French hospital discharge database	A logistic model based on 14 independent covariates included based on their ability to discriminate between atrial fibrillation and other treatment indications. Developed in a training sample of oral anticoagulant users with known treatment indication	Sample of oral anticoagulant users with known treatment indication (validation sample)	Patients with a known treatment indication were randomly divided into a training sample (50% of patients) and a validation sample (the remaining 50%). The training sample was used to develop the algorithm and estimate logistic regression coefficients, while the validation sample was used to select the best model based on an area under ROC curve (c-index). The predictive accuracy of the final model was assessed by determining discrimination (c-statistics) and calibration, both evaluated on the validation sample (Billionnet et al 2017 <sup>135</sup> ).
Five events deemed to be important in studying risks and benefits of treatments:	Eight European countries	Multi-sector, multi-lingual RCD from eight European countries with	The Unified Medical Language System was used to identify concepts and corresponding codes	Benchmarking using published rates	An iterative process of computing incidence rates of the five events, and obtaining feedback from database holders to explore any disparities/difference from expected/published rates (Coloma et al 2013 <sup>98</sup> ).

acute myocardial infarction; acute renal failure; anaphylactic shock; bullous eruption; rhabdomyolysis		heterogeneous classifications			
Hyponatremia	Denmark	Danish National Patient Registry	ICD-10 codes recorded at discharge	Another RCD: database collecting laboratory tests performed in hospital-based diagnoses. One sodium value <135 mmol/L measured at any time during hospitalisation confirmed the diagnosis.	Example of a validation study on a large scale, possible when one database can be plausibly assumed to be a reference standard (i.e., lab values are superior to ICD-10 codes). Study population included 819 701 patients of all ages admitted to hospital during the study period (2 186 642 hospitalisations). Estimated sensitivity, specificity, positive predictive value and negative predictive value for ICD-10 codes for hyponatraemia overall and for cut-off points for increasing hyponatraemia severity. An important advantage (Holland-Bill et al 2014 <sup>136</sup> ).
Asthma	Canada	Provincial health administrative data in British Columbia, Canada, containing billing records of physician visits, hospital discharge summaries, and prescription drug purchases.	Three algorithms based on combinations of ICD-9 codes and purchases of asthma medications	Latent class modelling based definition of treated asthma	Each algorithm component is associated with an estimable probability of assigning a given patient to a treatment asthma category, allowing for a range of different asthma definitions, depending on a given study aim, to prioritise high sensitivity for health services planning, or a need to restrict analysis to asthma cases requiring hospitalization (Prosser et al 2008 <sup>86</sup> ).
Serious infection	Denmark	Danish National Patient Registry	ICD-10 diagnoses for infections identified in primary position of inpatient hospital encounters	Chart review	Case definition was imported by adopting/mapping an ICD-9 codes validated in the US to ICD-10 codes as well as criteria for presence of a serious infection. Charts were reviewed by medical doctors. Data on presence of criteria of any or preselected specific infections (pneumonia, sepsis) were collected as well as information on a chart reviewer's clinical judgement, whether an infection was present and if so, what specific infection it was (Holland-Bill et al 2014 <sup>103</sup> ).
Cardiovascular risk scores previously developed and validated in other sources of data: CHADS <sub>2</sub> (congestive	Denmark	Linked data from Danish population registries	ICD-8/ICD-10 codes for congestive heart failure and stroke; ICD-8/ICD-10 codes and/or treatment proxies (ATC codes) for	Expected 1-, 5-, and 10-year risk of thromboembolism	Used time-to event analysis – Cox regression – to estimate c-statistic to assess the predictive capability of each score and its components to predict the recorded occurrence of thromboembolism (Olesen et al 2011 <sup>137</sup> )

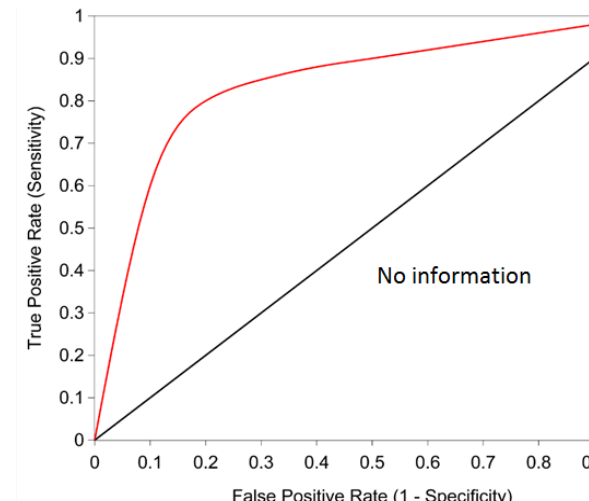


heart failure, hypertension, age $\geq$ 75 years, diabetes mellitus, previous stroke/transient ischaemic attack)  CHA2DS2-VASc (congestive heart failure, hypertension, age category $\geq$ 75 years, diabetes mellitus, previous, stroke/transient ischaemic attack, vascular disease)			diabetes and hypertension from outpatient dispensing data		
Registered cancer staging according to TNM at diagnosis	Denmark	Danish Cancer Registry	TNM tumour stage at diagnosis recorded in a cancer registry	Expectation that the data source must have complete data on TNM staging on all primary tumours captured in a national cancer registry	Computed proportion of records with missing TNM stage based on the expectation that the data source must have complete data on TNM staging (Deleruan et al 2012 <sup>138</sup> )
Multiple sclerosis (MS)	Canada	Ontario Electronic Medical Record Administrative data Linked Database (EMRALD)	Multiple algorithms based on combinations of cumulative patient profile, prescription and physicians billing codes were tested.	EMR search followed by chart review.	To establish reference standards, the 943 potential MS cases were identified among 73,003 enrollees of the database based on only information in the EMR (cumulative patient profile, free text entries for multiple sclerosis-related phrases, and physicians billing codes) but not the linked administrative data. Random sampling was not feasible since it was likely to yield a very small number of cases. Chart review of potential cases was conducted by trained chart abstractors (Krysko et al 2015 <sup>139</sup> ).
Stevens-Johnson Syndrome (SJS)/Toxic Epidermal Necrolysis (TEN)	Japan	A single educational hospital electronic medical records	Combination of ICD-10 codes, treatment/test/procedural proxies for SJS/TEN and other similar diseases for exclusion in a single	Another data source: disease registry at division of dermatology	This study is to develop and optimise the algorithm based on claims data to make differential diagnosis between SJS/TEN and other similar diseases. Expert dermatologist confirmed cases in the registry. The patients who had the ICD-10 codes of differential diagnoses were identified in the EHR of division of dermatology. They were tested against 1214 algorithms composed of all combinations of six

			hospital administrative data.		algorithm components to calculate validity and performance measures including sensitivity, specificity, PPV, NPV and diagnostic odds ratio (Fukasawa et al 2019 <sup>106</sup> )
Lung cancer	United States	The HealthCore Integrated Research Environment (HIRE)-Oncology clinical database that were linked with the HealthCore Integrated Research Database (HIRD).	Modified Duh Algorithm based on the first-line chemotherapy regimens for NSCLC in Healthcare Common Procedure Coding System and ICD-10, with exclusion criteria included procedures, surgeries, and chemotherapies administered to patients with SCLC as recommended by the 2006 NCCN Guidelines.	Another data source: Lung cancer patients identified in the HIRE-Oncology database based on a diagnosis of lung cancer with a molecular classification of NSCLC or SCLC.	Performance of NSCLC case-finding algorithm to distinguish NSCLC and SCLC cases was compared with the control algorithm using sensitivity, specificity, PPV, NPV, accuracy, diagnostic odds ratio (DOR), and area under the curve (AUC) for ROC. The control algorithm included all tests and first-line treatments based on the complete list of the American Cancer Society <sup>140</sup> and 2016 NCCN <sup>141</sup> recommended treatments for NSCLC and SCLC were combined and no exclusion criteria (Turner et al, 2017 <sup>142</sup> ).
Multiple malignancies	Taiwan	The National Health Insurance (NHI) database and the National Cancer Registry (NCR) in Taiwan	Cancer patients in the NHI database were identified based on the catastrophic illness certificates for medical co-payment waivers, which were recorded in the Registry of Catastrophic Illness Patients. Cancer types including all cancers and top 10 cancers were based on ICD-9.	Another data source: the National Cancer Registry (NCR), a mandatory population-based cancer registry%, respectively.	Deterministic linkage by the unique personal identification assigned to each Taiwanese resident was conducted, which also allows the linkage to civil registration, birth and death registries, and the NCR. Sensitivity, specificity, PPV, and NPV for all cancer and each cancer were estimated using the NCR as reference standard and survival recorded in the NHI were calculated. (Kao et al, 2017 <sup>143</sup> )
Multiple malignancies	United States	Medicare claims	Four candidate algorithms validated for cancers of lung, colorectal, breast, lymphoma and leukaemia: Definition 1: Combination of diagnosis and procedures on the same day or within the same hospitalisation Definition 2: Two diagnoses of a specific cancer within 2 months Definition 3: Definition 1	Another RCD source: Pennsylvania State Cancer Registry, with certified high-quality and completeness of data on malignancies	Cross-tabulation of the claims-based algorithms with the cancer registry data to assess sensitivity, specificity and PPV among 157,310 Medicare beneficiaries in 1997-2000. In addition, the study evaluated algorithm's ability to distinguish between incident and recurrent malignancies and evaluated the impact of imperfect completeness of the references standard in capturing haematological malignancies (Setoguchi et al, 2007 <sup>83</sup> ).

			or Definition 2 Definition 4: One diagnosis of cancer		
--	--	--	---	--	--

611 Abbreviations not explained in the table: ICD-9 International Classification of Diseases, 9<sup>th</sup> Revision; ICD-10 International Classification of Diseases, 10<sup>th</sup> Revision; ATC Anatomical  
612 Therapeutic Chemical; TNM Tumour, Nodes, Metastases; SCLC small cell lung cancer; NSCLC non-small cell lung cancer; EHR electronic health records; NSCP Nomesco  
613 Classification of Surgical Procedures  
614

Box 3. Measures of algorithm validity																		
Measure	Range	Description	Formulas															
Sensitivity	[0, 1]	Proportion of persons with a given health state according to gold/reference standard who are classified as such by an algorithm	<table border="1"> <thead> <tr> <th colspan="2"></th> <th colspan="2">Health state status based on gold/reference standard</th> </tr> <tr> <th colspan="2"></th> <th>Positive</th> <th>Negative</th> </tr> </thead> <tbody> <tr> <th rowspan="2">Health state status based on algorithm</th> <th>Positive</th> <td>A (True positive)</td> <td>C (False negative)</td> </tr> <tr> <th>Negative</th> <td>B (False positive)</td> <td>D (True negative)</td> </tr> </tbody> </table> $Sensitivity = \frac{A}{(A + B)} =$ $Specificity = \frac{D}{(C + D)} =$ $PPV = \frac{A}{(A + C)} =$ $NPV = \frac{D}{(B + D)} =$			Health state status based on gold/reference standard				Positive	Negative	Health state status based on algorithm	Positive	A (True positive)	C (False negative)	Negative	B (False positive)	D (True negative)
		Health state status based on gold/reference standard																
		Positive		Negative														
Health state status based on algorithm	Positive	A (True positive)		C (False negative)														
	Negative	B (False positive)	D (True negative)															
Specificity	[0, 1]	Proportion of persons without a given health state according to gold/reference standard who are classified as such by an algorithm																
Positive predictive value (PPV)	[0, 1]	Proportion of patients who truly have the health state among all those who are classified as positive by the algorithm																
Negative predictive value (NPV)	[0, 1]	Proportion of patients who truly do not have the health state among all those who are classified as negative by the algorithm. Chuback et al and Benchimol et al provide formulae connecting sensitivity and specificity, PPV and NPV with prevalence of a health state. <sup>30,64</sup>																
Receiver operating characteristic (ROC) curve		Definition of an RCD-based algorithm may be based on a cut-off value of a measured continuous variable (eg, hyponatremia based on serum potassium levels in the study by Holland-Bill et al <sup>156</sup> ) Generally, more extreme cut-off values of a continuous variable leads to increase in proportion of both false-positives, but also true-positives. For algorithms that classify patients into health states (eg., presence/absence/severity) ROC is a plot of sensitivity (or true positive proportion) on the x-axis against (1-specificity) or false-positive proportion on the y-axis. For a (hypothetical) perfect algorithm, sensitivity=specificity=1. <sup>79</sup>	 <p>Example of a receiver operating characteristic (ROC) curve and area under the curve (AUC). (Sørensen and Vandembroucke (in press).<sup>144</sup>)</p>															
Area under the curve (AUC)		Area under the curve (AUC) is derived from ROC curve and is a single measure, frequently used to indicate performance of diagnostic tests, or algorithms, as they are used to 'diagnose' presence of a health state in the study population. With varying thresholds for case definition, the points (1-Sensitivity, Specificity) are plotted on the plain with 1-Sensitivity on the x-axis and Specificity on the y-axis to construct the curve. An algorithm with sensitivity = specificity=1 ROC=1; an algorithm that would classify patients not better than a coin toss would have a AUC=0.5. <sup>79</sup>																

Diagnostic odds ratio (DOR)	[0, ∞)	An odds ratio of dichotomous tests in diagnostic application and frequently used for meta-analysis of diagnostic tests. For the purposes of validation studies, DOR can be calculated as the ratio of the odds of positivity in those with a health state of interest relative to the odds of positivity in those without the health state of interest. <sup>80</sup> DOR does not depend on prevalence of health state, with higher values indicating better discriminatory test performance and used in combination with sensitivity and specificity. A value of 1 means that the test does not discriminate between cases and non-cases. Higher DOR correspond to higher probability of an algorithm to be positive among true cases than in non-cases of a given health state.	$DOR = \frac{A (True\ Positive) / C (False\ Positive)}{B (False\ Negative) / D (True\ Negative)}$																																				
Kappa statistic	[-1, 1]	<p>Kappa statistic is used to quantify interrater variability.<sup>144</sup> In evaluating RCD-based algorithm, it can be used to quantify agreement between two algorithms or data sources, none of which can be considered a better (a reference standard) relative to the other.</p> <p>Based on the value of kappa statistic, agreement may be qualified on a scale 'less than chance' to 'almost perfect' (Viera and Garrett 2005<sup>81</sup>.)</p> <table border="0" data-bbox="560 1033 787 1316"> <tr> <td>Kappa</td> <td>Agreement</td> </tr> <tr> <td>&lt; 0</td> <td>Less than chance</td> </tr> <tr> <td>0.01–0.20</td> <td>Slight</td> </tr> <tr> <td>0.21–0.40</td> <td>Fair</td> </tr> <tr> <td>0.41–0.60</td> <td>Moderate</td> </tr> <tr> <td>0.61–0.80</td> <td>Substantial</td> </tr> <tr> <td>0.81–0.99</td> <td>Almost perfect</td> </tr> </table>	Kappa	Agreement	< 0	Less than chance	0.01–0.20	Slight	0.21–0.40	Fair	0.41–0.60	Moderate	0.61–0.80	Substantial	0.81–0.99	Almost perfect	<table border="1" data-bbox="930 685 1531 1131"> <thead> <tr> <th colspan="2" rowspan="2"></th> <th colspan="2">Health state status according to data source 1</th> <th rowspan="2"></th> </tr> <tr> <th>Positive</th> <th>Negative</th> </tr> </thead> <tbody> <tr> <th rowspan="2">Health state status according to data source 2</th> <th>Positive</th> <td>a</td> <td>b</td> <td>m<sub>1</sub></td> </tr> <tr> <th>Negative</th> <td>c</td> <td>d</td> <td>m<sub>0</sub></td> </tr> <tr> <td colspan="2"></td> <td>n<sub>1</sub></td> <td>n<sub>0</sub></td> <td>n</td> </tr> </tbody> </table> <p>Expected agreement between two data sources/algorithm, none of which is superior to the other in classifying a given health state into present/absent</p> $Observed\ agreement\ p_{obs} = \frac{(a + d)}{n}$ $Expected\ agreement\ p_{exp} = \left[ \left( \frac{n_1}{n} \right) \times \left( \frac{m_1}{n} \right) \right] + \left[ \left( \frac{n_0}{n} \right) \times \left( \frac{m_0}{n} \right) \right]$ $Kappa\ statistic\ \kappa = \frac{(p_{obs} - p_{exp})}{(1 - p_{exp})}$				Health state status according to data source 1			Positive	Negative	Health state status according to data source 2	Positive	a	b	m <sub>1</sub>	Negative	c	d	m <sub>0</sub>			n <sub>1</sub>	n <sub>0</sub>	n
Kappa	Agreement																																						
< 0	Less than chance																																						
0.01–0.20	Slight																																						
0.21–0.40	Fair																																						
0.41–0.60	Moderate																																						
0.61–0.80	Substantial																																						
0.81–0.99	Almost perfect																																						
		Health state status according to data source 1																																					
		Positive	Negative																																				
Health state status according to data source 2	Positive	a	b	m <sub>1</sub>																																			
	Negative	c	d	m <sub>0</sub>																																			
		n <sub>1</sub>	n <sub>0</sub>	n																																			

616 **References**

- 617 1. Makady A, de Boer A, Hillege H, et al. What Is Real-World Data? A Review of Definitions Based on Literature  
618 and Stakeholder Interviews. *Value Health* 2017;20(7):858-65. doi:  
619 <https://doi.org/10.1016/j.jval.2017.03.008>
- 620 2. Corrigan-Curay J, Sacks L, Woodcock J. Real-World Evidence and Real-World Data for Evaluating Drug Safety  
621 and Effectiveness. *JAMA* 2018;320(9):867-68. doi: 10.1001/jama.2018.10136 [published Online First:  
622 2018/08/15]
- 623 3. Berger ML, Sox H, Willke RJ, et al. Good practices for real-world data studies of treatment and/or comparative  
624 effectiveness: Recommendations from the joint ISPOR-ISPE Special Task Force on real-world evidence in  
625 health care decision making. *Pharmacoepidemiol Drug Saf* 2017;26(9):1033-39. doi: 10.1002/pds.4297  
626 [published Online First: 2017/09/16]
- 627 4. Sherman RE, Anderson SA, Dal Pan GJ, et al. Real-World Evidence - What Is It and What Can It Tell Us? *N*  
628 *Engl J Med* 2016;375(23):2293-97. doi: 10.1056/NEJMsb1609216 [published Online First: 2016/12/14]
- 629 5. Administration UFaD. Real-World Evidence 2019 [Available from: [https://www.fda.gov/science-](https://www.fda.gov/science-research/science-and-research-special-topics/real-world-evidence)  
630 [research/science-and-research-special-topics/real-world-evidence](https://www.fda.gov/science-research/science-and-research-special-topics/real-world-evidence) accessed 16 October 2019.
- 631 6. Booth CM, Karim S, Mackillop WJ. Real-world data: towards achieving the achievable in cancer care. *Nature*  
632 *Reviews Clinical Oncology* 2019;16(5):312-25. doi: 10.1038/s41571-019-0167-7
- 633 7. Administration UFaD. Use of Real-World Evidence to Support Regulatory Decision-Making for Medical Devices  
634 2017 [Available from: [https://www.fda.gov/regulatory-information/search-fda-guidance-documents/use-](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/use-real-world-evidence-support-regulatory-decision-making-medical-devices)  
635 [real-world-evidence-support-regulatory-decision-making-medical-devices](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/use-real-world-evidence-support-regulatory-decision-making-medical-devices) accessed 16 October 2019.
- 636 8. Administration UFaD. Submitting Documents Using Real-World Data and Real-World Evidence to FDA for  
637 Drugs and Biologics 2019 [Available from: [https://www.fda.gov/regulatory-information/search-fda-](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/submitting-documents-using-real-world-data-and-real-world-evidence-fda-drugs-and-biologics-guidance)  
638 [guidance-documents/submitting-documents-using-real-world-data-and-real-world-evidence-fda-drugs-and-](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/submitting-documents-using-real-world-data-and-real-world-evidence-fda-drugs-and-biologics-guidance)  
639 [biologics-guidance](https://www.fda.gov/regulatory-information/search-fda-guidance-documents/submitting-documents-using-real-world-data-and-real-world-evidence-fda-drugs-and-biologics-guidance) accessed 16 October 2019.
- 640 9. Plueschke K, McGettigan P, Pacurariu A, et al. EU-funded initiatives for real world evidence: descriptive analysis  
641 of their characteristics and relevance for regulatory decision-making. *BMJ Open* 2018;8(6):e021864. doi:  
642 10.1136/bmjopen-2018-021864
- 643 10. Pharmaceuticals and Medical Devices Agency. The guidelines for the conduct of pharmacoepidemiological  
644 studies using medical information databases, etc., in evaluation of drug safety. 1 ed. Tokyo2014 [Available  
645 from: <https://www.pmda.go.jp/files/000147250.pdf> accessed 31 January 2020.
- 646 11. Pharmaceutical Evaluation Division and Pharmaceutical Safety Division. Basic Principles on the utilization of  
647 health information databases for Post-Marketing Surveillance of Medical Products Tokyo2017 [Available  
648 from: <https://www.pmda.go.jp/files/000218531.pdf> accessed 31 January 2020.
- 649 12. Ministry of Health Labour and Welfare. Amendment of Ministerial Ordinance Concerning Good Post-Marketing  
650 Study Practice [in Japanese]  
651 Tokyo2017 [Available from: [https://elaws.e-](https://elaws.e-gov.go.jp/search/elawsSearch/elaws_search/lsg0500/detail?lawId=416M60000100171)  
652 [gov.go.jp/search/elawsSearch/elaws\\_search/lsg0500/detail?lawId=416M60000100171](https://elaws.e-gov.go.jp/search/elawsSearch/elaws_search/lsg0500/detail?lawId=416M60000100171) accessed 31 January  
653 2020.
- 654 13. European Medicines Agency. Medical Devices [Available from: [https://www.ema.europa.eu/en/human-](https://www.ema.europa.eu/en/human-regulatory/overview/medical-devices)  
655 [regulatory/overview/medical-devices](https://www.ema.europa.eu/en/human-regulatory/overview/medical-devices) accessed 24 February 2020.
- 656 14. Strom BL, Kimmel SE, Hennessy S. Pharmacoepidemiology. Sixth edition. ed. Hoboken, NJ: Wiley-Blackwell  
657 2020.
- 658 15. Kwakkenbos L, Imran M, McCall SJ, et al. CONSORT extension for the reporting of randomised controlled  
659 trials conducted using cohorts and routinely collected data (CONSORT-ROUTINE): checklist with  
660 explanation and elaboration. *BMJ* 2021;373:n857. doi: 10.1136/bmj.n857 [published Online First:  
661 2021/05/01]
- 662 16. About the Food and Drug Administration (FDA) Sentinel Initiative [Available from:  
663 <https://www.sentinelinitiative.org/about#section-1592851590618> accessed 11 May 2021.
- 664 17. (EMA) EMA. Technical workshop on real-world metadata for regulatory purposes 2021 [Available from:  
665 [https://www.ema.europa.eu/en/events/technical-workshop-real-world-metadata-regulatory-](https://www.ema.europa.eu/en/events/technical-workshop-real-world-metadata-regulatory-purposes#documents-section)  
666 [purposes#documents-section](https://www.ema.europa.eu/en/events/technical-workshop-real-world-metadata-regulatory-purposes#documents-section) accessed 19 May 2021.
- 667 18. Rajkomar A, Dean J, Kohane I. Machine Learning in Medicine. *N Engl J Med* 2019;380(14):1347-58. doi:  
668 10.1056/NEJMra1814259 [published Online First: 2019/04/04]
- 669 19. Eichler HG, Bloechl-Daum B, Broich K, et al. Data Rich, Information Poor: Can We Use Electronic Health  
670 Records to Create a Learning Healthcare System for Pharmaceuticals? *Clin Pharmacol Ther*  
671 2019;105(4):912-22. doi: 10.1002/cpt.1226 [published Online First: 2018/09/05]

- 672 20. Fletcher RH, Fletcher SW. Clinical epidemiology : the essentials. 4th ed. Philadelphia: Lippincott Williams &  
673 Wilkins 2005.
- 674 21. Weiss NS. Clinical epidemiology : the study of the outcome of illness. 3rd ed. Oxford ; New York: Oxford  
675 University Press 2006.
- 676 22. Maclure M, Schneeweiss S. Causation of bias: the episcopo. *Epidemiology* 2001;12(1):114-22. [published  
677 Online First: 2001/01/04]
- 678 23. Bosco JL, Silliman RA, Thwin SS, et al. A most stubborn bias: no adjustment method fully resolves  
679 confounding by indication in observational studies. *J Clin Epidemiol* 2010;63(1):64-74. doi:  
680 10.1016/j.jclinepi.2009.03.001 [published Online First: 2009/05/22]
- 681 24. Jurek AM, Greenland S, Maldonado G, et al. Proper interpretation of non-differential misclassification effects:  
682 expectations vs observations. *Int J Epidemiol* 2005;34(3):680-7. doi: 10.1093/ije/dyi060 [published Online  
683 First: 2005/04/02]
- 684 25. Rothman KJ, Greenland S, Lash TL. Chapter 9. Validity in Epidemiologic Studies. In: Rothman KJ, Greenland  
685 S, Lash TL, eds. Modern epidemiology. 3rd ed. Philadelphia: Wolters Kluwer Health/Lippincott Williams  
686 & Wilkins 2008:x, 758 p.
- 687 26. Manuel DG, Rosella LC, Stukel TA. Importance of accurately identifying disease in studies using electronic  
688 health records. *BMJ* 2010;341:c4226. doi: 10.1136/bmj.c4226 [published Online First: 2010/08/21]
- 689 27. Greenland S. Basic methods for sensitivity analysis of biases. *Int J Epidemiol* 1996;25(6):1107-16. [published  
690 Online First: 1996/12/01]
- 691 28. van Smeden M, Lash TL, Groenwold RHH. Reflection on modern methods: five myths about measurement error  
692 in epidemiological research. *Int J Epidemiol* 2020;49(1):338-47. doi: 10.1093/ije/dy251 [published Online  
693 First: 2019/12/11]
- 694 29. Agniel D, Kohane IS, Weber GM. Biases in electronic health record data due to processes within the healthcare  
695 system: retrospective observational study. *BMJ* 2018;361:k1479. doi: 10.1136/bmj.k1479 [published  
696 Online First: 2018/05/02]
- 697 30. Benchimol EI, Manuel DG, To T, et al. Development and use of reporting guidelines for assessing the quality of  
698 validation studies of health administrative data. *J Clin Epidemiol* 2011;64(8):821-9. doi:  
699 10.1016/j.jclinepi.2010.10.006 [published Online First: 2011/01/05]
- 700 31. Lanes S, Brown JS, Haynes K, et al. Identifying health outcomes in healthcare databases. *Pharmacoepidemiol*  
701 *Drug Saf* 2015;24(10):1009-16. doi: 10.1002/pds.3856 [published Online First: 2015/08/19]
- 702 32. Sørensen HT, Baron JA. Registries and medical databases. In: Olsen J, Greene N, Saracci R, et al., eds. Teaching  
703 epidemiology : a guide for teachers in epidemiology, public health and clinical medicine  
704 Fourth edition. ed. Oxford, United Kingdom ; New York, NY, United States of America: Oxford University Press  
705 2015.
- 706 33. Meng XL. Statistical Paradises and Paradoxes in Big Data (I): Law of Large Populations, Big Data Paradox, and  
707 the 2016 Us Presidential Election. *Ann Appl Stat* 2018;12(2):685-726. doi: 10.1214/18-Aoas1161sf
- 708 34. Sorensen HT, Sabroe S, Olsen J. A framework for evaluation of secondary data sources for epidemiological  
709 research. *Int J Epidemiol* 1996;25(2):435-42. [published Online First: 1996/04/01]
- 710 35. Schneeweiss S, Eichler HG, Garcia-Altes A, et al. Real World Data in Adaptive Biomedical Innovation: A  
711 Framework for Generating Evidence Fit for Decision-Making. *Clin Pharmacol Ther* 2016;100(6):633-46.  
712 doi: 10.1002/cpt.512 [published Online First: 2016/10/21]
- 713 36. Franklin JM, Glynn RJ, Martin D, et al. Evaluating the Use of Nonrandomized Real-World Data Analyses for  
714 Regulatory Decision Making. *Clin Pharmacol Ther* 2019;105(4):867-77. doi: 10.1002/cpt.1351 [published  
715 Online First: 2019/01/14]
- 716 37. Jick SS, Kaye JA, Vasilakis-Scaramozza C, et al. Validity of the general practice research database.  
717 *Pharmacotherapy* 2003;23(5):686-9. [published Online First: 2003/05/14]
- 718 38. Ludvigsson JF, Andersson E, Ekbom A, et al. External review and validation of the Swedish national inpatient  
719 register. *BMC Public Health* 2011;11(1):450. doi: 10.1186/1471-2458-11-450 [published Online First:  
720 2011/06/11]
- 721 39. Carnahan RM. Mini-Sentinel's systematic reviews of validated methods for identifying health outcomes using  
722 administrative data: summary of findings and suggestions for future research. *Pharmacoepidemiol Drug*  
723 *Saf* 2012;21 Suppl 1:90-9. doi: 10.1002/pds.2318
- 724 40. Springate DA, Kontopantelis E, Ashcroft DM, et al. ClinicalCodes: an online clinical codes repository to  
725 improve the validity and reproducibility of research using electronic medical records. *PLoS One*  
726 2014;9(6):e99825. doi: 10.1371/journal.pone.0099825 [published Online First: 2014/06/19]
- 727 41. Schmidt M, Schmidt SA, Sandegaard JL, et al. The Danish National Patient Registry: a review of content, data  
728 quality, and research potential. *Clin Epidemiol* 2015;7:449-90. doi: 10.2147/CLEP.S91125 [published  
729 Online First: 2015/11/26]

- 730 42. Mini-Sentinel. [Available from: <http://www.mini-sentinel.org/> accessed 7 December 2015.
- 731 43. Koram N, Delgado M, Stark JH, et al. Validation studies of claims data in the Asia-Pacific region: A  
732 comprehensive review. *Pharmacoepidemiol Drug Saf* 2019;28(2):156-70. doi: 10.1002/pds.4616
- 733 44. von Elm E, Altman DG, Egger M, et al. Strengthening the Reporting of Observational Studies in Epidemiology  
734 (STROBE) statement: guidelines for reporting observational studies. *BMJ* 2007;335(7624):806-8. doi:  
735 10.1136/bmj.39335.541782.AD [published Online First: 2007/10/20]
- 736 45. Vandembroucke JP, von Elm E, Altman DG, et al. Strengthening the Reporting of Observational Studies in  
737 Epidemiology (STROBE): explanation and elaboration. *Ann Intern Med* 2007;147(8):W163-94. doi:  
738 10.7326/0003-4819-147-8-200710160-00010-w1 [published Online First: 2007/10/17]
- 739 46. Nicholls SG, Langan SM, Sorensen HT, et al. The RECORD reporting guidelines: meeting the methodological  
740 and ethical demands of transparency in research using routinely-collected health data. *Clin Epidemiol*  
741 2016;8:389-92. doi: 10.2147/CLEP.S110528
- 742 47. Benchimol EI, Langan S, Guttman A, et al. Call to RECORD: the need for complete reporting of research using  
743 routinely collected health data. *J Clin Epidemiol* 2013;66(7):703-5. doi: 10.1016/j.jclinepi.2012.09.006
- 744 48. Langan SM, Benchimol EI, Guttman A, et al. Setting the RECORD straight: developing a guideline for the  
745 REporting of studies Conducted using Observational Routinely collected Data. *Clin Epidemiol* 2013;5:29-  
746 31. doi: 10.2147/clep.s36885 [published Online First: 2013/02/16]
- 747 49. Nicholls SG, Quach P, von Elm E, et al. The REporting of Studies Conducted Using Observational Routinely-  
748 Collected Health Data (RECORD) Statement: Methods for Arriving at Consensus and Developing  
749 Reporting Guidelines. *PLoS One* 2015;10(5):e0125620. doi: 10.1371/journal.pone.0125620
- 750 50. Benchimol EI, Smeeth L, Guttman A, et al. The REporting of studies Conducted using Observational  
751 Routinely-collected health Data (RECORD) statement. *PLoS Med* 2015;12(10):e1001885. doi:  
752 10.1371/journal.pmed.1001885 [published Online First: 2015/10/07]
- 753 51. Langan SM, Schmidt SA, Wing K, et al. The reporting of studies conducted using observational routinely  
754 collected health data statement for pharmacoepidemiology (RECORD-PE). *BMJ* 2018;363:k3532. doi:  
755 10.1136/bmj.k3532 [published Online First: 2018/11/16]
- 756 52. Ehrenstein V, Petersen I, Smeeth L, et al. Helping everyone do better: a call for validation studies of routinely  
757 recorded health data. *Clin Epidemiol* 2016;8:49-51. doi: 10.2147/CLEP.S104448 [published Online First:  
758 2016/04/26]
- 759 53. Chun DS, Lund JL, Sturmer T. Pharmacoepidemiology and Drug Safety's special issue on validation studies.  
760 *Pharmacoepidemiol Drug Saf* 2019;28(2):123-25. doi: 10.1002/pds.4694 [published Online First:  
761 2019/02/05]
- 762 54. Lash TL, Olshan AF. EPIDEMIOLOGY Announces the "Validation Study" Submission Category.  
763 *Epidemiology* 2016;27(5):613-4. doi: 10.1097/EDE.0000000000000532 [published Online First:  
764 2016/07/09]
- 765 55. Nissen F, Quint JK, Morales DR, et al. How to validate a diagnosis recorded in electronic health records.  
766 *Breathe (Sheff)* 2019;15(1):64-68. doi: 10.1183/20734735.0344-2018 [published Online First: 2019/03/07]
- 767 56. Fox MP, Lash TL, Bodnar LM. Common misconceptions about validation studies. *Int J Epidemiol* 2020 doi:  
768 10.1093/ije/dyaa090
- 769 57. International Society for Pharmacoepidemiology. Guidelines for Good Pharmacoepidemiology Practices (GPP)  
770 2015 [Available from: <https://www.pharmacoepi.org/resources/policies/guidelines-08027> accessed 23 May  
771 2021.
- 772 58. Hines PA, Janssens R, Gonzalez-Quevedo R, et al. A future for regulatory science in the European Union: the  
773 European Medicines Agency's strategy. *Nat Rev Drug Discov* 2020;19(5):293-94. doi: 10.1038/d41573-  
774 020-00032-0 [published Online First: 2020/04/03]
- 775 59. Orsini LS, Monz B, Mullins CD, et al. Improving transparency to build trust in real-world secondary data studies  
776 for hypothesis testing-Why, what, and how: recommendations and a road map from the real-world evidence  
777 transparency initiative. *Pharmacoepidemiol Drug Saf* 2020;29(11):1504-13. doi: 10.1002/pds.5079  
778 [published Online First: 2020/09/15]
- 779 60. Wang SV, Schneeweiss S, Berger ML, et al. Reporting to Improve Reproducibility and Facilitate Validity  
780 Assessment for Healthcare Database Studies V1.0. *Pharmacoepidemiol Drug Saf* 2017;26(9):1018-32. doi:  
781 10.1002/pds.4295 [published Online First: 2017/09/16]
- 782 61. Greenfield S. Making Real-World Evidence More Useful for Decision Making. *Value Health* 2017;20(8):1023-  
783 24. doi: 10.1016/j.jval.2017.08.3012 [published Online First: 2017/10/02]
- 784 62. Nicholls SG, Langan SM, Benchimol EI. Routinely collected data: the importance of high-quality diagnostic  
785 coding to research. *CMAJ* 2017;189(33):E1054-e55. doi: 10.1503/cmaj.170807 [published Online First:  
786 2017/08/23]



- 787 63. Nicholls SG, Langan SM, Benchimol EI, et al. Reporting transparency: making the ethical mandate explicit.  
788 *BMC Med* 2016;14:44. doi: 10.1186/s12916-016-0587-5
- 789 64. Chubak J, Pocobelli G, Weiss NS. Tradeoffs between accuracy measures for electronic health care data  
790 algorithms. *J Clin Epidemiol* 2012;65(3):343-49 e2. doi: 10.1016/j.jclinepi.2011.09.002
- 791 65. Schmidt M, Jacobsen JB, Lash TL, et al. 25 year trends in first time hospitalisation for acute myocardial  
792 infarction, subsequent short and long term mortality, and the prognostic impact of sex and comorbidity: a  
793 Danish nationwide cohort study. *BMJ* 2012;344:e356. doi: 10.1136/bmj.e356 [published Online First:  
794 2012/01/27]
- 795 66. Vinter N, Linder M, Andersen M, et al. Classification and characteristics of on-label and off-label apixaban use  
796 in Denmark and Sweden. *Pharmacoepidemiol Drug Saf* 2019;28(6):867-78. doi: 10.1002/pds.4778  
797 [published Online First: 2019/04/18]
- 798 67. Shivade C, Raghavan P, Fosler-Lussier E, et al. A review of approaches to identifying patient phenotype cohorts  
799 using electronic health records. *J Am Med Inform Assoc* 2014;21(2):221-30. doi: 10.1136/amiajnl-2013-  
800 001935 [published Online First: 2013/11/10]
- 801 68. Wong J, Horwitz MM, Zhou L, et al. Using machine learning to identify health outcomes from electronic health  
802 record data. *Curr Epidemiol Rep* 2018;5(4):331-42. doi: 10.1007/s40471-018-0165-9 [published Online  
803 First: 2018/12/18]
- 804 69. Tanskanen A, Taipale H, Koponen M, et al. From prescription drug purchases to drug use periods - a second  
805 generation method (PRE2DUP). *BMC Med Inform Decis Mak* 2015;15:21. doi: 10.1186/s12911-015-0140-  
806 z [published Online First: 2015/04/19]
- 807 70. Kirby JC, Speltz P, Rasmussen LV, et al. PheKB: a catalog and workflow for creating electronic phenotype  
808 algorithms for transportability. *J Am Med Inform Assoc* 2016;23(6):1046-52. doi: 10.1093/jamia/ocv202  
809 [published Online First: 2016/03/31]
- 810 71. Beam AL, Kohane IS. Big Data and Machine Learning in Health Care. *JAMA* 2018;319(13):1317-18. doi:  
811 10.1001/jama.2017.18391 [published Online First: 2018/03/14]
- 812 72. Pottgard A, Schmidt SAJ, Wallach-Kildemoes H, et al. Data Resource Profile: The Danish National  
813 Prescription Registry. *Int J Epidemiol* 2017;46(3):798-98f. doi: 10.1093/ije/dyw213 [published Online  
814 First: 2016/10/30]
- 815 73. Schmidt M, Hallas J, Friis S. Potential of prescription registries to capture individual-level use of aspirin and  
816 other nonsteroidal anti-inflammatory drugs in Denmark: trends in utilization 1999-2012. *Clin Epidemiol*  
817 2014;6:155-68. doi: 10.2147/CLEP.S59156 [published Online First: 2014/05/30]
- 818 74. Sweden bans over-the-counter sales of painkiller over heart failure fears 2019 [Available from:  
819 [https://www.thelocal.se/20191105/sweden-bans-over-the-counter-sales-of-painkiller-over-heart-failure-](https://www.thelocal.se/20191105/sweden-bans-over-the-counter-sales-of-painkiller-over-heart-failure-fears)  
820 [fears](https://www.thelocal.se/20191105/sweden-bans-over-the-counter-sales-of-painkiller-over-heart-failure-fears) accessed 28 August 2020.
- 821 75. Schmidt M, Sorensen HT, Pedersen L. Diclofenac use and cardiovascular risks: series of nationwide cohort  
822 studies. *BMJ* 2018;362:k3426. doi: 10.1136/bmj.k3426 [published Online First: 2018/09/06]
- 823 76. Deeks JJ. Systematic reviews in health care: Systematic reviews of evaluations of diagnostic and screening tests.  
824 *BMJ* 2001;323(7305):157-62. doi: 10.1136/bmj.323.7305.157 [published Online First: 2001/07/21]
- 825 77. West SL, Ritchey ME, Poole C. Chapter 12. Validity of Pharmacoepidemiologic Drug and Diagnosis Data. In:  
826 Strom BL, Kimmel SE, Hennessy S, eds. Textbook of Pharmacoepidemiology. 5th ed. Chichester, West  
827 Sussex, UK: Wiley-Blackwell 2013:p.
- 828 78. Bollaerts K, Rekkas A, De Smedt T, et al. Disease misclassification in electronic healthcare database studies:  
829 Deriving validity indices-A contribution from the ADVANCE project. *PLoS One* 2020;15(4):e0231333.  
830 doi: 10.1371/journal.pone.0231333 [published Online First: 2020/04/23]
- 831 79. Bewick V, Cheek L, Ball J. Statistics review 13: receiver operating characteristic curves. *Crit Care*  
832 2004;8(6):508-12. doi: 10.1186/cc3000 [published Online First: 2004/11/30]
- 833 80. Glas AS, Lijmer JG, Prins MH, et al. The diagnostic odds ratio: a single indicator of test performance. *J Clin*  
834 *Epidemiol* 2003;56(11):1129-35. doi: 10.1016/s0895-4356(03)00177-x [published Online First:  
835 2003/11/15]
- 836 81. Viera AJ, Garrett JM. Understanding interobserver agreement: the kappa statistic. *Fam Med* 2005;37(5):360-3.  
837 [published Online First: 2005/05/11]
- 838 82. Fox MP, Lash TL, Greenland S. A method to automate probabilistic sensitivity analyses of misclassified binary  
839 variables. *Int J Epidemiol* 2005;34(6):1370-6. doi: 10.1093/ije/dyi184 [published Online First: 2005/09/21]
- 840 83. Setoguchi S, Solomon DH, Glynn RJ, et al. Agreement of diagnosis and its date for hematologic malignancies  
841 and solid tumors between medicare claims and cancer registry data. *Cancer Causes Control*  
842 2007;18(5):561-9. doi: 10.1007/s10552-007-0131-1 [published Online First: 2007/04/21]
- 843 84. Berkson J. Limitations of the application of fourfold table analysis to hospital data. *Biometrics* 1946;2(3):47-53.  
844 [published Online First: 1946/06/01]

- 845 85. Ording AG, Cronin-Fenton D, Ehrenstein V, et al. Challenges in translating endpoints from trials to  
846 observational cohort studies in oncology. *Clin Epidemiol* 2016;8:195-200. doi: 10.2147/CLEP.S97874  
847 [published Online First: 2016/06/30]
- 848 86. Prosser RJ, Carleton BC, Smith MA. Identifying persons with treated asthma using administrative data via latent  
849 class modelling. *Health Serv Res* 2008;43(2):733-54. doi: 10.1111/j.1475-6773.2007.00775.x [published  
850 Online First: 2008/03/29]
- 851 87. Morkem R, Handelman K, Queenan JA, et al. Validation of an EMR algorithm to measure the prevalence of  
852 ADHD in the Canadian Primary Care Sentinel Surveillance Network (CPCSSN). *BMC Med Inform Decis  
853 Mak* 2020;20(1):166. doi: 10.1186/s12911-020-01182-2 [published Online First: 2020/07/22]
- 854 88. Comparative effectiveness and safety of non-vitamin K oral anticoagulants and warfarin in non-valvular atrial  
855 fibrillation - a cohort study in 3 Nordic countries. ESC Congress; 2019; Paris, France.
- 856 89. Nielsen PB, Skjoth F, Sogaard M, et al. Effectiveness and safety of reduced dose non-vitamin K antagonist oral  
857 anticoagulants and warfarin in patients with atrial fibrillation: propensity weighted nationwide cohort study.  
858 *BMJ* 2017;356:j510. doi: 10.1136/bmj.j510
- 859 90. Hellfritsch M, Pottegard A, Haastrup SB, et al. Cohort selection in register-based studies of direct oral  
860 anticoagulant users with atrial fibrillation: An inevitable trade-off between selection bias and  
861 misclassification. *Basic Clin Pharmacol Toxicol* 2020;127(1):3-5. doi: 10.1111/bcpt.13423 [published  
862 Online First: 2020/05/05]
- 863 91. Riis AH, Johansen MB, Jacobsen JB, et al. Short look-back periods in pharmacoepidemiologic studies of new  
864 users of antibiotics and asthma medications introduce severe misclassification. *Pharmacoepidemiol Drug  
865 Saf* 2015;24(5):478-85. doi: 10.1002/pds.3738 [published Online First: 2015/01/21]
- 866 92. Schmidt M, Schmidt SAJ, Adelborg K, et al. The Danish health care system and epidemiological research: from  
867 health care contacts to database records. *Clin Epidemiol* 2019;11:563-91. doi: 10.2147/CLEP.S179083  
868 [published Online First: 2019/08/03]
- 869 93. Bakken IJ, Ariansen AMS, Knudsen GP, et al. The Norwegian Patient Registry and the Norwegian Registry for  
870 Primary Health Care: Research potential of two nationwide health-care registries. *Scand J Public Health*  
871 2019;0(0):1403494819859737. doi: 10.1177/1403494819859737 [published Online First: 2019/07/11]
- 872 94. Gribsholt SB, Pedersen L, Richelsen B, et al. Validity of ICD-10 diagnoses of overweight and obesity in Danish  
873 hospitals. *Clin Epidemiol* 2019;11:845-54. doi: 10.2147/CLEP.S214909 [published Online First:  
874 2019/10/02]
- 875 95. Sturmer T, Glynn RJ, Rothman KJ, et al. Adjustments for unmeasured confounders in pharmacoepidemiologic  
876 database studies using external information. *Med Care* 2007;45(10 Supl 2):S158-65. doi:  
877 10.1097/MLR.0b013e318070c045 [published Online First: 2007/10/25]
- 878 96. Valkhoff VE, Coloma PM, Masclee GM, et al. Validation study in four health-care databases: upper  
879 gastrointestinal bleeding misclassification affects precision but not magnitude of drug-related upper  
880 gastrointestinal bleeding risk. *J Clin Epidemiol* 2014;67(8):921-31. doi: 10.1016/j.jclinepi.2014.02.020  
881 [published Online First: 2014/05/06]
- 882 97. Gini R, Schuemie M, Brown J, et al. Data Extraction and Management in Networks of Observational Health  
883 Care Databases for Scientific Research: A Comparison of EU-ADR, OMOP, Mini-Sentinel and MATRICE  
884 Strategies. *EGEMS (Wash DC)* 2016;4(1):1189. doi: 10.13063/2327-9214.1189 [published Online First:  
885 2016/03/26]
- 886 98. Avillach P, Coloma PM, Gini R, et al. Harmonization process for the identification of medical events in eight  
887 European healthcare databases: the experience from the EU-ADR project. *J Am Med Inform Assoc*  
888 2013;20(1):184-92. doi: 10.1136/amiajnl-2012-000933 [published Online First: 2012/09/08]
- 889 99. Ehrenstein V, Huang K, Kahlert J, et al. Abstract. *Pharmacoepidemiol Drug Saf* 2017;26(S2):3-636. doi:  
890 10.1002/pds.4275
- 891 100. Forns J, Cainzos-Achirica M, Hellfritsch M, et al. Validity of ICD-9 and ICD-10 codes used to identify acute  
892 liver injury: A study in three European data sources. *Pharmacoepidemiol Drug Saf* 2019;28(7):965-75. doi:  
893 10.1002/pds.4803 [published Online First: 2019/06/07]
- 894 101. ACCESS STUDY PLACEHOLDER - AVAILABLE IN JUNE 2021 FOR CITATION. PLACEHOLDER
- 895 102. Lawlor DA, Tilling K, Davey Smith G. Triangulation in aetiological epidemiology. *Int J Epidemiol*  
896 2016;45(6):1866-86. doi: 10.1093/ije/dyw314 [published Online First: 2017/01/22]
- 897 103. Holland-Bill L, Xu H, Sorensen HT, et al. Positive predictive value of primary inpatient discharge diagnoses of  
898 infection among cancer patients in the Danish National Registry of Patients. *Ann Epidemiol*  
899 2014;24(8):593-7, 97 e1-18. doi: 10.1016/j.annepidem.2014.05.011 [published Online First: 2014/08/03]
- 900 104. Benchimol EI, Guttman A, Mack DR, et al. Validation of international algorithms to identify adults with  
901 inflammatory bowel disease in health administrative data from Ontario, Canada. *J Clin Epidemiol*  
902 2014;67(8):887-96. doi: 10.1016/j.jclinepi.2014.02.019 [published Online First: 2014/04/30]

903 105. Flahault A, Cadilhac M, Thomas G. Sample size calculation should be performed for design accuracy in  
904 diagnostic test studies. *J Clin Epidemiol* 2005;58(8):859-62. doi: 10.1016/j.jclinepi.2004.12.009 [published  
905 Online First: 2005/07/16]

906 106. Fukasawa T, Takahashi H, Kameyama N, et al. Development of an electronic medical record-based algorithm  
907 to identify patients with Stevens-Johnson syndrome and toxic epidermal necrolysis in Japan. *PLoS One*  
908 2019;14(8):e0221130. doi: 10.1371/journal.pone.0221130 [published Online First: 2019/08/14]

909 107. Holcroft CA, Spiegelman D. Design of validation studies for estimating the odds ratio of exposure-disease  
910 relationships when exposure is misclassified. *Biometrics* 1999;55(4):1193-201. doi: 10.1111/j.0006-  
911 341x.1999.01193.x [published Online First: 2001/04/21]

912 108. Collin LJ, MacLehose RF, Ahern TP, et al. Adaptive Validation Design: A Bayesian Approach to Validation  
913 Substudy Design With Prospective Data Collection. *Epidemiology* 2020;31(4):509-16. doi:  
914 10.1097/EDE.0000000000001209 [published Online First: 2020/06/03]

915 109. Lash TL, Fox MP, Fink AK. Applying quantitative bias analysis to epidemiologic data  
916 (<https://sites.google.com/site/biasanalysis/>). Dordrecht ; New York: Springer 2009.

917 110. van Walraven C, Hart RG. Leave 'em alone - why continuous variables should be analyzed as such.  
918 *Neuroepidemiology* 2008;30(3):138-9. doi: 10.1159/000126908 [published Online First: 2008/04/19]

919 111. Hall GC, Lanes S, Bollaerts K, et al. Outcome misclassification: Impact, usual practice in  
920 pharmacoepidemiology database studies and an online aid to correct biased estimates of risk ratio or  
921 cumulative incidence. *Pharmacoepidemiol Drug Saf* 2020;29(11):1450-55. doi: 10.1002/pds.5109  
922 [published Online First: 2020/08/30]

923 112. Lash TL, Fox MP, MacLehose RF, et al. Good practices for quantitative bias analysis. *Int J Epidemiol*  
924 2014;43(6):1969-85. doi: 10.1093/ije/dyu149 [published Online First: 2014/08/01]

925 113. Rothman KJ. Episheet [Available from: <http://krothman.hostbyet2.com/episheet.xls> accessed 26 October 2019.

926 114. MacLehose RF, Bodnar LM, Meyer CS, et al. Hierarchical Semi-Bayes Methods for Misclassification in  
927 Perinatal Epidemiology. *Epidemiology* 2018;29(2):183-90. doi: 10.1097/ede.0000000000000789  
928 [published Online First: 2017/11/23]

929 115. Walraven CV. A comparison of methods to correct for misclassification bias from administrative database  
930 diagnostic codes. *Int J Epidemiol* 2018;47(2):605-16. doi: 10.1093/ije/dyx253 [published Online First:  
931 2017/12/19]

932 116. Chu H, Wang Z, Cole SR, et al. Sensitivity analysis of misclassification: a graphical and a Bayesian approach.  
933 *Ann Epidemiol* 2006;16(11):834-41. doi: 10.1016/j.annepidem.2006.04.001 [published Online First:  
934 2006/07/18]

935 117. van Walraven C. Bootstrap imputation with a disease probability model minimized bias from misclassification  
936 due to administrative database codes. *J Clin Epidemiol* 2017;84:114-20. doi: 10.1016/j.jclinepi.2017.01.007  
937 [published Online First: 2017/02/09]

938 118. Newcomer SR, Xu S, Kulldorff M, et al. A primer on quantitative bias analysis with positive predictive values  
939 in research using electronic health data. *J Am Med Inform Assoc* 2019;26(12):1664-74. doi:  
940 10.1093/jamia/ocz094 [published Online First: 2019/08/01]

941 119. Funk MJ, Landi NS. Misclassification in administrative claims data: quantifying the impact on treatment effect  
942 estimates. *Curr Epidemiol Rep* 2014;1(4):175-85. doi: 10.1007/s40471-014-0027-z [published Online First:  
943 2015/06/19]

944 120. Herrett E, Thomas SL, Schoonen WM, et al. Validation and validity of diagnoses in the General Practice  
945 Research Database: a systematic review. *Br J Clin Pharmacol* 2010;69(1):4-14. doi: 10.1111/j.1365-  
946 2125.2009.03537.x [published Online First: 2010/01/19]

947 121. Centers for Medicare and Medicaid Services. [Available from:  
948 [https://www.cms.gov/Newsroom/MediaReleaseDatabase/Press-releases/2014-Press-releases-items/2014-  
949 07-31.html](https://www.cms.gov/Newsroom/MediaReleaseDatabase/Press-releases/2014-Press-releases-items/2014-07-31.html)] accessed 9 December 2015.

950 122. ICD-11 [Available from: <https://icd.who.int/en/> accessed 06 May 2019.

951 123. Miksad RA, Abernethy AP. Harnessing the Power of Real-World Evidence (RWE): A Checklist to Ensure  
952 Regulatory-Grade Data Quality. *Clin Pharmacol Ther* 2018;103(2):202-05. doi: 10.1002/cpt.946  
953 [published Online First: 2017/12/08]

954 124. WHO. ICD-11: Classifying disease to map the way we live and die 2019 [Available from:  
955 <https://www.who.int/health-topics/international-classification-of-diseases> accessed 06 May 2019.

956 125. Gottlieb S. Remarks to the National Academy of Sciences on the Impact of Real World Evidence on Medical  
957 Product Development 2017 [Available from: <https://www.fda.gov/newsevents/speeches/ucm576519.htm>  
958 accessed 02 July 2021.

- 959 126. Lund JL, Froslev T, Deleuran T, et al. Validity of the Danish National Registry of Patients for chemotherapy  
960 reporting among colorectal cancer patients is high. *Clin Epidemiol* 2013;5:327-34. doi:  
961 10.2147/CLEP.S49773 [published Online First: 2013/09/17]
- 962 127. Tanskanen A, Taipale H, Koponen M, et al. Drug exposure in register-based research-An expert-opinion based  
963 evaluation of methods. *PLoS One* 2017;12(9):e0184070. doi: 10.1371/journal.pone.0184070 [published  
964 Online First: 2017/09/09]
- 965 128. Adelborg K, Sundboll J, Munch T, et al. Positive predictive value of cardiac examination, procedure and  
966 surgery codes in the Danish National Patient Registry: a population-based validation study. *BMJ Open*  
967 2016;6(12):e012817. doi: 10.1136/bmjopen-2016-012817 [published Online First: 2016/12/13]
- 968 129. Ehrenstein V, Gammelager H, Schiodt M, et al. Evaluation of an ICD-10 algorithm to detect osteonecrosis of  
969 the jaw among cancer patients in the Danish National Registry of Patients. *Pharmacoepidemiol Drug Saf*  
970 2015;24(7):693-700. doi: 10.1002/pds.3786 [published Online First: 2015/05/15]
- 971 130. Lash TL, Riis AH, Ostefeld EB, et al. A validated algorithm to ascertain colorectal cancer recurrence using  
972 registry resources in Denmark. *Int J Cancer* 2015;136(9):2210-15. doi: 10.1002/ijc.29267
- 973 131. Cuthbertson CC, Kucharska-Newton A, Faurot KR, et al. Controlling for Frailty in Pharmacoepidemiologic  
974 Studies of Older Adults: Validation of an Existing Medicare Claims-based Algorithm. *Epidemiology*  
975 2018;29(4):556-61. doi: 10.1097/EDE.0000000000000833 [published Online First: 2018/04/06]
- 976 132. Tapper EB, Korovaichuk S, Baki J, et al. Identifying Patients With Hepatic Encephalopathy Using  
977 Administrative Data in the ICD-10 Era. *Clin Gastroenterol Hepatol* 2019 doi:  
978 <https://doi.org/10.1016/j.cgh.2019.12.017>
- 979 133. Pedersen SA, Schmidt SAJ, Klausen S, et al. Melanoma of the Skin in the Danish Cancer Registry and the  
980 Danish Melanoma Database: A Validation Study. *Epidemiology* 2018;29(3):442-47. doi:  
981 10.1097/ede.0000000000000802 [published Online First: 2018/01/18]
- 982 134. Benchimol EI, Guttman A, Griffiths AM, et al. Increasing incidence of paediatric inflammatory bowel disease  
983 in Ontario, Canada: evidence from health administrative data. *Gut* 2009;58(11):1490-7. doi:  
984 10.1136/gut.2009.188383 [published Online First: 2009/08/05]
- 985 135. Billionnet C, Alla F, Berigaud E, et al. Identifying atrial fibrillation in outpatients initiating oral anticoagulants  
986 based on medico-administrative data: results from the French national healthcare databases.  
987 *Pharmacoepidemiol Drug Saf* 2017;26(5):535-43. doi: 10.1002/pds.4192 [published Online First:  
988 2017/03/16]
- 989 136. Holland-Bill L, Christiansen CF, Ulrichsen SP, et al. Validity of the International Classification of Diseases,  
990 10th revision discharge diagnosis codes for hyponatraemia in the Danish National Registry of Patients.  
991 *BMJ Open* 2014;4(4):e004956. doi: 10.1136/bmjopen-2014-004956 [published Online First: 2014/04/25]
- 992 137. Olesen JB, Lip GY, Hansen ML, et al. Validation of risk stratification schemes for predicting stroke and  
993 thromboembolism in patients with atrial fibrillation: nationwide cohort study. *BMJ* 2011;342:d124. doi:  
994 10.1136/bmj.d124
- 995 138. Deleuran T, Sogaard M, Froslev T, et al. Completeness of TNM staging of small-cell and non-small-cell lung  
996 cancer in the Danish cancer registry, 2004-2009. *Clin Epidemiol* 2012;4:39-44. doi: 10.2147/CLEP.S33315  
997 [published Online First: 2012/09/01]
- 998 139. Krysko KM, Ivers NM, Young J, et al. Identifying individuals with multiple sclerosis in an electronic medical  
999 record. *Mult Scler* 2015;21(2):217-24. doi: 10.1177/1352458514538334 [published Online First:  
1000 2014/06/21]
- 1001 140. Sacher AG, Le LW, Lau A, et al. Real-world chemotherapy treatment patterns in metastatic non-small cell  
1002 lung cancer: Are patients undertreated? *Cancer* 2015;121(15):2562-69. doi: 10.1002/ncr.29386
- 1003 141. Ettinger DS, Wood DE, Akerley W, et al. NCCN Guidelines Insights: Non-Small Cell Lung Cancer, Version  
1004 4.2016. *J Natl Compr Canc Netw* 2016;14(3):255-64. doi: 10.6004/jnccn.2016.0031 [published Online  
1005 First: 2016/03/10]
- 1006 142. Turner RM, Chen Y-W, Fernandes AW. Validation of a Case-Finding Algorithm for Identifying Patients with  
1007 Non-small Cell Lung Cancer (NSCLC) in Administrative Claims Databases. *Frontiers in Pharmacology*  
1008 2017;8(883) doi: 10.3389/fphar.2017.00883
- 1009 143. Kao W-H, Hong J-H, See L-C, et al. Validity of cancer diagnosis in the National Health Insurance database  
1010 compared with the linked National Cancer Registry in Taiwan. *Pharmacoepidemiology and drug safety*  
1011 2018;27(10):1060-66. doi: 10.1002/pds.4267
- 1012 144. Vandenbroucke JP, Sorensen HT. Chapter \*\*. Clinical Epidemiology. In: Rothman KJ, Greenland S, Lash TL,  
1013 eds. *Modern Epidemiology*. 4th ed. Philadelphia: Wolters Kluwer Health/Lippincott Williams & Wilkins  
1014 (in press).

- 1015 145. Gini R, Dodd CN, Bollaerts K, et al. Quantifying outcome misclassification in multi-database studies: The case  
1016 study of pertussis in the ADVANCE project. *Vaccine* 2020;38 Suppl 2:B56-B64. doi:  
1017 10.1016/j.vaccine.2019.07.045 [published Online First: 2019/11/05]  
1018 146. Lohse SR, Farkas DK, Lohse N, et al. Validation of spontaneous abortion diagnoses in the Danish National  
1019 Registry of Patients. *Clin Epidemiol* 2010;2:247-50. doi: 10.2147/CLEP.S13815 [published Online First:  
1020 2010/12/15]  
1021 147. Acquavella J, Ehrenstein V, Schiodt M, et al. Design and methods for a Scandinavian pharmacovigilance study  
1022 of osteonecrosis of the jaw and serious infections among cancer patients treated with antiresorptive agents  
1023 for the prevention of skeletal-related events. *Clin Epidemiol* 2016;8:267-72. doi: 10.2147/CLEP.S107270  
1024 [published Online First: 2016/08/09]  
1025 148. Xue F, Ma H, Stehman-Breen C, et al. Design and methods of a postmarketing pharmacoepidemiology study  
1026 assessing long-term safety of Prolia(R) (denosumab) for the treatment of postmenopausal osteoporosis.  
1027 *Pharmacoepidemiol Drug Saf* 2013;22(10):1107-14. doi: 10.1002/pds.3477 [published Online First:  
1028 2013/07/17]  
1029 149. Gammelager H, Svaerke C, Noerholt SE, et al. Validity of an algorithm to identify osteonecrosis of the jaw in  
1030 women with postmenopausal osteoporosis in the Danish National Registry of Patients. *Clin Epidemiol*  
1031 2013;5:263-7. doi: 10.2147/CLEP.S45226 [published Online First: 2013/08/16]  
1032 150. Schiodt M, Larsson Wexell C, Herlofson BB, et al. Existing data sources for clinical epidemiology:  
1033 Scandinavian Cohort for osteonecrosis of the jaw - work in progress and challenges. *Clin Epidemiol*  
1034 2015;7:107-16. doi: 10.2147/CLEP.S71796 [published Online First: 2015/02/07]  
1035 151. Abstract 2436. Pharmacovigilance cohort study of osteonecrosis of the jaw and serious infection among cancer  
1036 patients treated with denosumab or zoledronic acid in Denmark, Norway, and Sweden. ICPE All Access  
1037 Virtual Event, September 2020; 2020; Virtual.  
1038 152. Abstract 3835. A multinational European study of anaphylaxis among recipients of intravenous Iron. ICPE All  
1039 Access Virtual Event, September 2020; 2020; Virtual.  
1040

1041 **Figure legends**

1042 *Figure 1 Point of care to RWD data point and back*

1043 *Figure 1 Incidence rates (per 100,000 person-years) of pertussis infection identified in RWD in*

1044 *three countries using various RWD-based algorithms (horizontal bars) and using a gold*

1045 *standard (the dashed lines), i.e., national surveillance data reported to the European Centre for*

1046 *Disease Prevention and Control. For RWD-based composite algorithms, the grey bar represents*

1047 *cases detected only by the left-hand component (indicated in the label before the key Boolean*

1048 *operator ‘OR’); the black bar represents cases detected by both components; the white bar*

1049 *represents cases detected by the right-hand component (indicated in the label after the key*

1050 *Boolean operator word ‘OR’). Reproduced from Gini et al. 2020 (permitted given full*

1051 *attribution).*<sup>145</sup>

1052 *Figure 2 An example of a decision process when considering validation studies/algorithm*  
1053 *selection. As an example of a ready-to-use algorithm, consider ICD codes for spontaneous*  
1054 *abortion recorded in the Danish National Patient Registry. International classifications of*  
1055 *diseases have specific diagnostic codes for spontaneous abortion, and a validation study*  
1056 *revealed that >95% of ICD records corresponded to a spontaneous abortion record in medical*  
1057 *chart(reference standard), with slight variation by period, and type of hospital.<sup>146</sup> Therefore, the*  
1058 *ICD-10 based algorithm is suitable for constructing cohorts of women with spontaneous*  
1059 *abortions or for studying relative associations between an exposure and the outcome of*  
1060 *spontaneous abortion. Use of the algorithm for studies evaluating absolute risks or risk*  
1061 *differences for the outcome of spontaneous abortions will need evidence (or assumption) of high*  
1062 *sensitivity of the algorithm and high completeness of the data source. The two were not examined*  
1063 *in the validations study, but the assumptions may be defensible using the knowledge that most*  
1064 *clinically apparent spontaneous abortions in Denmark are seen in hospital settings. (Very early*  
1065 *events may not be clinically apparent even to the affected woman and are therefore generally*  
1066 *difficult to quantify.) On the other end of the spectrum is osteonecrosis of the jaw – a side effect*  
1067 *of antiresorptive therapy, whose risk is an important parameter of the therapy safety profile –*  
1068 *exemplifies a condition without a clear-cut algorithm. Antiresorptive therapy is used in*  
1069 *osteoporosis (low dose) in bone malignancies (high dose), and risk of osteonecrosis of the jaw is*  
1070 *dose dependent. Although ICD-10 has a potentially useful code M87.1 “Osteonecrosis due to*  
1071 *drugs”, consistency of its use was unclear, including attribution to the specific agent.*  
1072 *Furthermore, majority of the osteonecrosis cases affect bones other than the jaw. In evaluating*  
1073 *suitability of the Danish National Patient Registry’s for estimating risks of osteonecrosis of the*  
1074 *jaw in postauthorisation safety studies of antiresorptive agents,<sup>147,148</sup> a candidate algorithm*  
1075 *consisting of a list ICD-10 codes originating from departments of oral and maxillofacial surgery*

1076 (to ensure jaw localization) was validated (external validation). The PPV of the candidate  
1077 algorithm ranged from 20% in patients with osteoporosis to 42% in patients with cancer, and  
1078 had a sensitivity of 73%.<sup>129 149</sup> Because the studies' main objective was estimation of absolute  
1079 risks of osteonecrosis of the jaw, the PPV and the sensitivity were deemed unsuitable for this  
1080 purpose, cases of osteonecrosis of the jaw for the study were identified directly by at the  
1081 departments of oral and maxillofacial surgery.<sup>150 151</sup> In a multinational European  
1082 postauthorisation safety study of anaphylaxis following the use of intravenous iron preparations,  
1083 the initial RWD-based algorithm for anaphylaxis originated from earlier US studies, whose PPV  
1084 was estimated in the European setting using medical chart review; notably, application of a US-  
1085 based algorithm in the European setting yielded a lower-than-expected risk of anaphylaxis,  
1086 potentially indicative of differences in diagnostic or recording practices leading to record  
1087 generation in the two types of setting.<sup>152</sup>

1088 *Figure 3 Gold standard and reference standard*

1089 *A: True positive, B: False negative, C: False positive, D: True negative*

1090 *Inner circle represents positive reference standard.*

1091 *a\*, b\*, c\*, d\* represent misclassifications with reference standard*

1092 *Unbiased sensitivity =  $(a+a^*)/(a+a^*+b+b^*)$*

1093 *Sensitivity against reference standard =  $(a+c^*)/(a+b+c^*+d^*)$*

1094 *c\* and d\* are consequently eliminated from sensitivity calculation through chart review against  
1095 reference standard cases.*

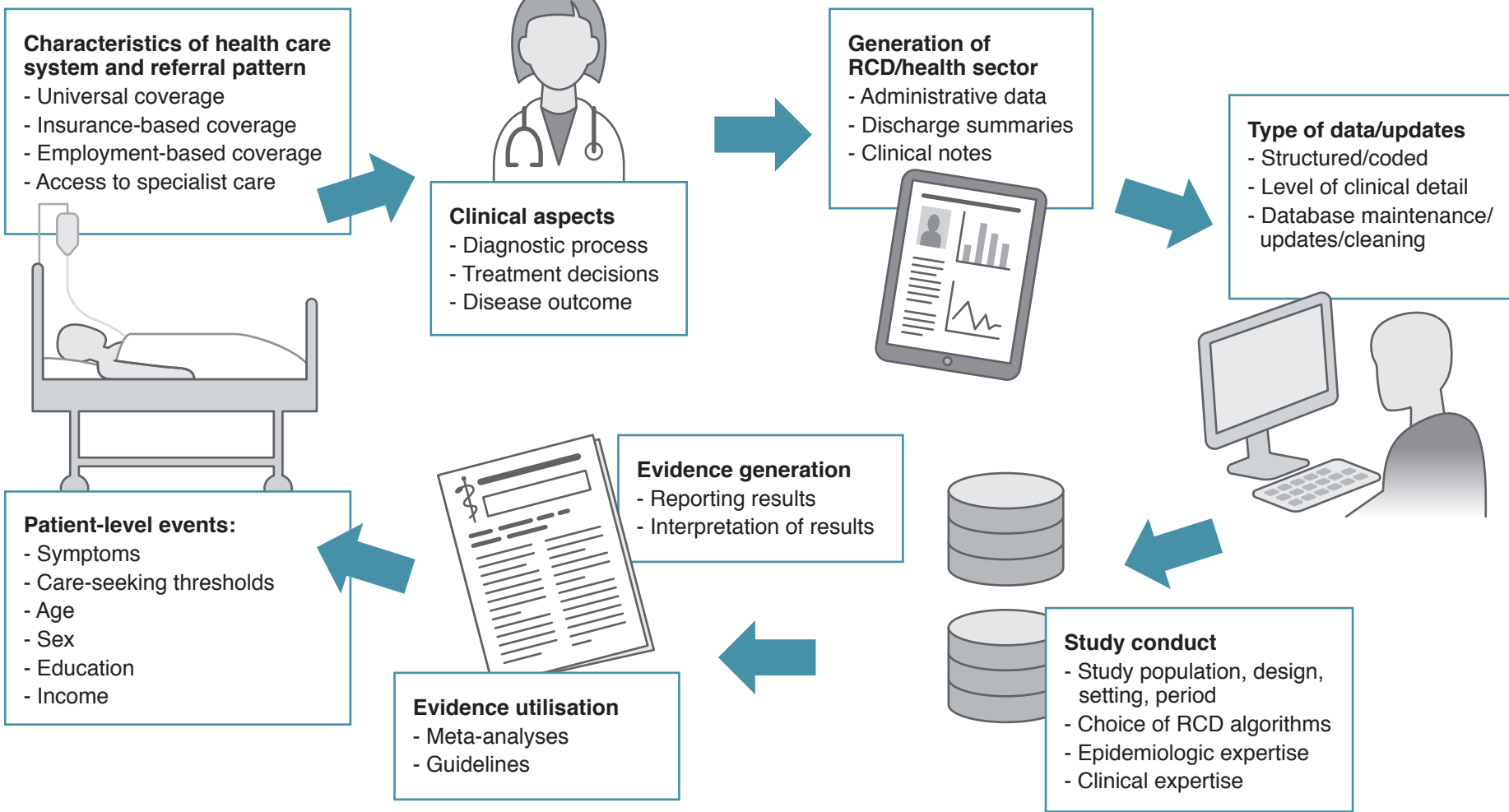
1096 *if  $(a:b=a^*:b^*)$  holds, no bias. Otherwise, bias in sensitivity unavoidable.*

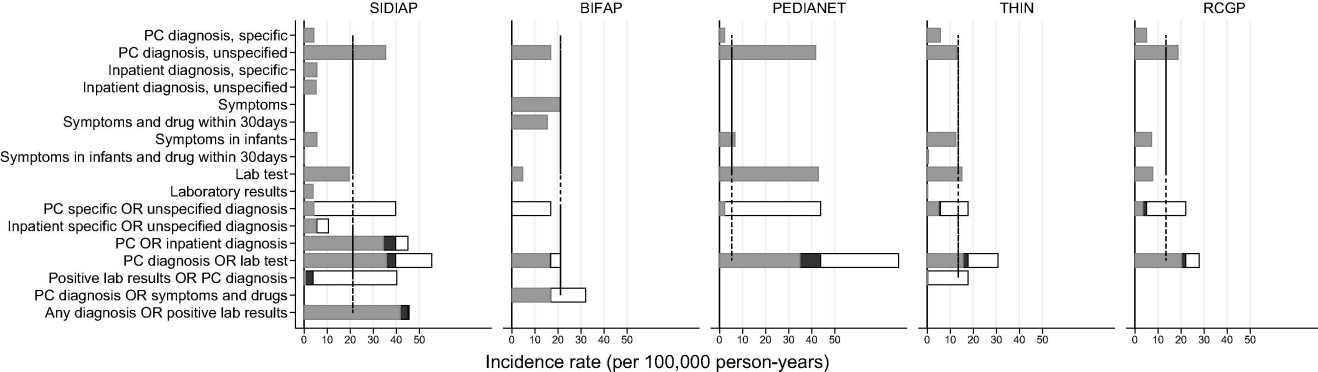
1097 *It is necessary to review the charts of all non-cases against reference standards  $(a^*+b^*+c+d)$ ,*

1098 *Remarks to the National Academy of Sciences on the Impact of Real World Evidence on Medical*

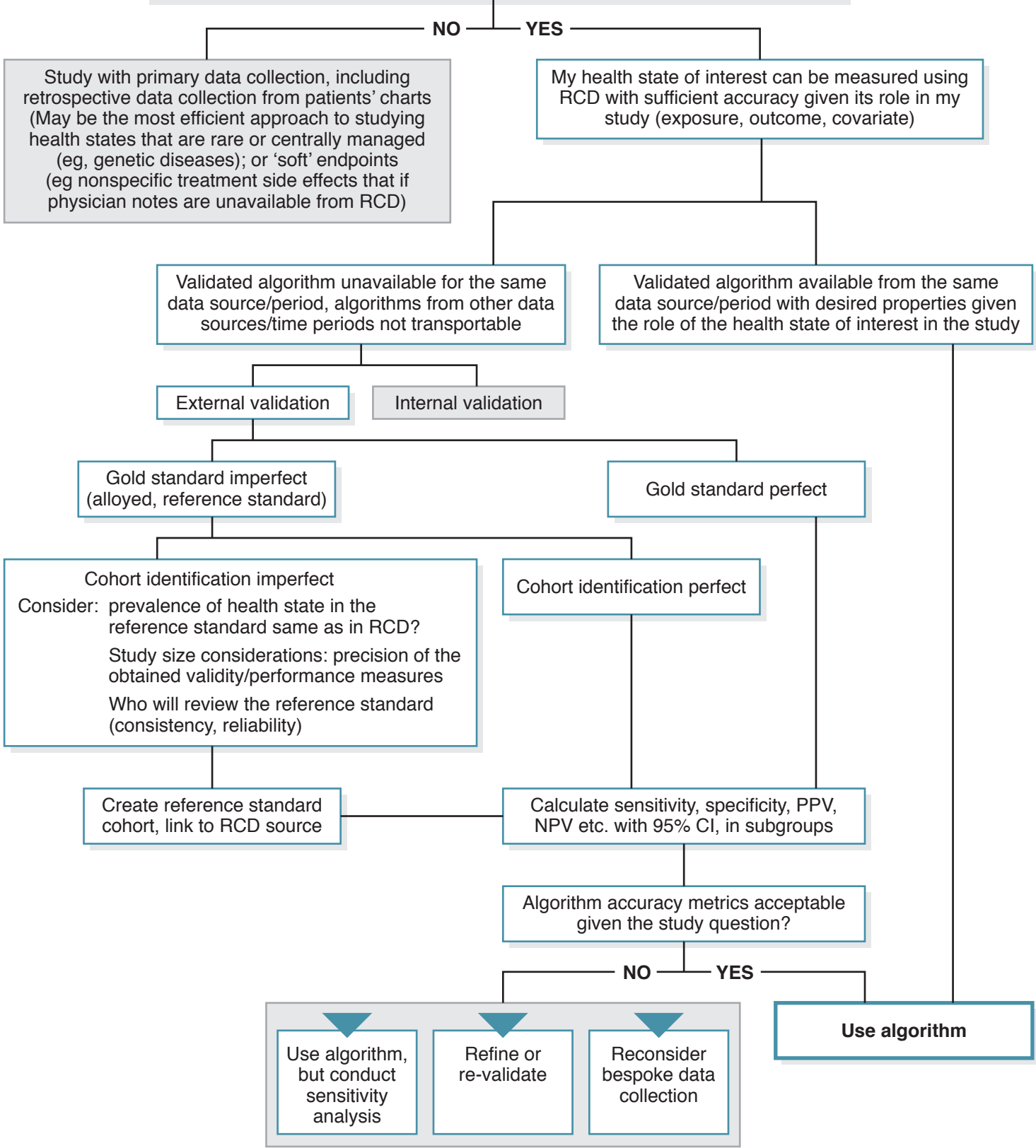


1099 *Product Development to determine the numbers of  $a^*$  and  $b^*$ , possibly taking a tremendous*  
1100 *effort.*

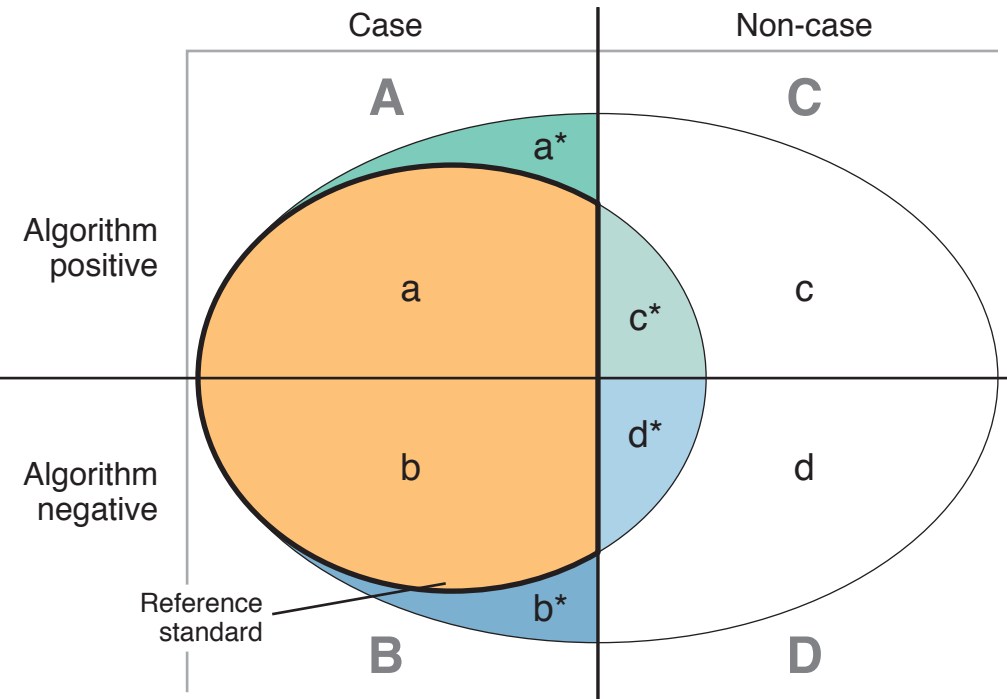




**Can RCD be considered to measure my health state of interest?**  
Consider: How prevalent is the health state?  
How complete is the RCD source (sector, health system/access)?  
How rich are the data (structured data/free text; specific/nonspecific codes/severity; subtype)?



# Gold standard



A: True positive

B: False negative

C: False positive

D: True negative

Inner circle represents positive reference standard.

Misclassifications with reference standard:  $a^*$ ,  $b^*$ ,  $c^*$ ,  $d^*$

Unbiased sensitivity =  $(a+a^*)/(a+a^*+b+b^*)$

Sensitivity against reference standard  
=  $(a+c^*)/(a+b+c^*+d^*)$

$c^*$  and  $d^*$  are consequently eliminated from sensitivity calculation through chart review against reference standard cases.

if  $(a:b=a^*:b^*)$  holds, no bias. Otherwise, bias in sensitivity unavoidable.

However, it necessitates reviewing all non-reference standards ( $a^*+b^*+c+d$ ) against chart to determine the numbers of  $a^*$  and  $b^*$ , possibly taking a tremendous effort.