

VALIDATION STUDIES: ISPE GUIDELINE



Misclassified present-day Bahamas as Spain



Christopher Columbus



Misclassified Cuba as mainland China



Amerigo Vespucci



Christopher Columbus

MISCLASSIFICATION IN EPIDEMIOLOGY

Wrong classification of study subjects:

- Classifying treated as untreated and vice versa
- Classifying cases as noncases and vice versa
- Classifying smokers as nonsmokers and vice versa
- More complex scenarios with > 2 levels

a.k.a.: measurement error

a.k.a.: information bias

SOURCES OF ERROR

At patient level

- Coding errors



ICD-10 G40 Epilepsy	69	100.0%
Epilepsy	52	75.4%
Seizures	5	7.2%
Suspected seizures	5	7.2%
Coding error	2	2.9%
Asphyxia	1	1.4%
Mental retardation	1	1.4%
Unspecified neurologic problem	1	1.4%
Heart failure	2	2.9%

MEDICAL CHART REVIEW



ISPE MANUSCRIPT: THE MOTIVATION

Received: 15 May 2020 | Accepted: 22 December 2020





DOI: 10.1002/pds.5193



ORIGINAL ARTICLE

WILEY

Outcomes in patients with lung cancer treated with crizotinib and erlotinib in routine clinical practice: A post-authorization safety cohort study conducted in Europe and in the United States

Vera Ehrenstein¹  | Kui Huang²  | Johnny Kahlert¹ | Shahram Bahmanyar³ | Pär Karlsson³ | Lukas Löfling³  | Anthony P. Nunes⁴ | Cheryl Enger⁴ | Irene D. Bezemer⁵ | Josephina G. Kuiper⁵  | Fabian Hoti⁶ | Rosa Juuti⁶ | Pasi Korhonen⁶ | Jingping Mo² | Stephen E. Schachterle^{2,7} | Keith D. Wilner² | Mikael Rørth¹ | Henrik T. Sørensen¹

QTc prolongation (ECG records reviewed by Adjudication Committee on request)

If the abstraction form includes evidence of any of the following, then event is definite:

Severe/symptomatic

- Clinical diagnosis of polymorphic ventricular tachycardia rate lasting for more than 5 seconds or ventricular fibrillation, documented by an ECG recording;
- Syncope recorded in case notes;
- Clinical diagnosis of torsade de pointes or TdP documented by an ECG recording;
- Clinical diagnosis of sudden death or sudden cardiac death;

Mild/asymptomatic

- ECG showing Fridericia rate-corrected QT interval ($QT/RR^{1/3}$) >450 ms in men or >470 ms in women;
- QTcF (Fridericia correction) change ≥ 60 msec from baseline;
- Absolute QT or QTc of >500 msec.

Exclude: Known coronary syndromes or known congenital long QT syndromes

Vision disorders

If the abstraction form includes evidence of any of the following then event is definite:

- Clinical diagnosis of vision disturbances including blurred vision, photophobia, photopsia, palinopsia (ask about multiple sclerosis), visual illusion, reduced visual acuity, diplopia, visual impairment, visual field defect, and vitreous floaters, maculopathy, retinal edema, retinal hemorrhage.

Exclude:

- Refractive error, amblyopia, corneal disorder (abnormal sensation in eye, anterior chamber collapse, anterior chamber opacity, aqueous humour leakage, asthenopia, chemical burns of eye, chemical eye injury, contact lens intolerance, corneal suture, corneal sutures removal, deposit eye, dry eye, eye burns, eye inflammation, eye injury, eye irritation, eye laser surgery, eye operation complication, eye penetration, flat anterior chamber of eye, foreign body in eye, foreign body sensation in eyes, hypoaesthesia eye, ocular toxicity, slit-lamp tests abnormal, superficial injury of eye, thermal burns of eye, vitamin A deficiency eye disorder, xerophthalmia, acquired corneal dystrophy, allergic keratitis, arcus lipoides, atopic keratoconjunctivitis, benign neoplasm of cornea, biopsy cornea, biopsy cornea abnormal, bowman's membrane disorder, corneal abrasion, corneal bleeding, corneal cyst, corneal decompensation, corneal defect, corneal degeneration, corneal deposits, corneal diameter decreased, corneal diameter increased, corneal disorder, corneal endothelial cell loss, corneal endotheliitis, corneal epithelial microcysts, corneal epithelium defect, corneal erosion, corneal exfoliation, corneal flap complication, corneal graft rejection, corneal hypertrophy, corneal implant, corneal infiltrates, corneal lesion, corneal lesion removal, corneal light reflex test abnormal, corneal oedema, corneal opacity, corneal operation, corneal perforation, corneal pigmentation, corneal reflex decreased, corneal scar, corneal staining, corneal striae, corneal thickening, corneal thinning, corneal touch, corneal transplant, corneal intraepithelial neoplasia, dellen, detached Descemet's membrane, diffuse lamellar keratitis, injury corneal, iridocorneal endothelial syndrome, Kayser-Fleischer ring, keratectomy, keratitis, keratitis interstitial, keratitis sclerosing, keratoconus, keratomalacia, keratometry, keratomileusis, keratopathy, keratorhexis, keratotomy, limbal hyperaemia, limbal swelling, macrocornea, malignant neoplasm of cornea, microcornea, neoplasm of cornea unspecified malignancy, neurotrophic keratopathy, photokeratitis, photorefractive keratectomy, punctate keratitis, Terrien's marginal degeneration, topography corneal abnormal, ulcerative keratitis, vital dye staining cornea present, vitamin A deficiency related corneal disorder), visual impairing cataracts, uncontrolled diabetes, brain tumor, age-related macular degeneration, toxic maculopathy (eg. Chloroquine or Tamoxifen)

COMMENTARY

WILEY

Validation of safety outcomes in routinely collected data: Lessons learned from a multinational postapproval safety study

Patrick J. Arena^{1,2}  | Kui Huang¹  | Lukas Löfling^{3,4}  | Shahram Bahmanyar³ |
Jingping Mo⁵ | Stephen E. Schachterle¹ | Anthony P. Nunes^{6,7} |
Elisabeth Smits⁸ | Rosa Juuti⁹ | Fabian Hoti⁹ | Pasi Korhonen⁹ |
Kasper Adelborg¹⁰ | Jens Sundbøll¹⁰ | Torben Riis Rasmussen¹¹ |
Anders Løkke^{12,13} | Vera Ehrenstein⁹ 

Key Lessons Learned



Trial-based definitions may be unattainable when using routinely collected data



It is crucial to compare clinical trials criteria to medical records data



Investigators should ensure they have permission to fully access relevant medical records



Investigators should consider key differences between geographies



Multiple methods to evaluate outcomes should be considered



Clear distinctions are needed to differentiate cases from noncases



Trial-based definitions are often unrealistic for use as “gold standards” when validating real-world data algorithms















American Journal of Epidemiology, 2024, 00, 1–13

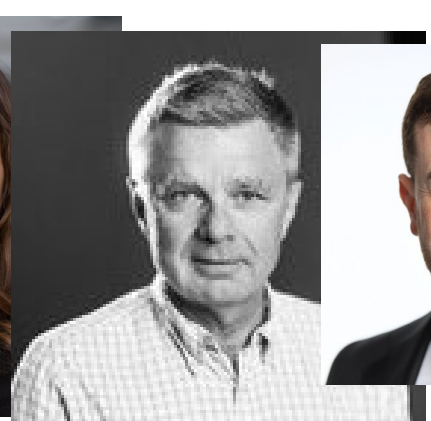
<https://doi.org/10.1093/aje/kwae071>

Advance access publication date May 17, 2024

Practice of Epidemiology

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

Vera Ehrenstein^{*,1} , Maja Hellfritsch^{2,3} , Johnny Kahlert¹, Sinéad M. Langan⁴ , Hisashi Urushihara⁵ ,
Danica Marinac-Dabic⁶ , Jennifer L. Lund⁷ , Henrik Toft Sørensen¹ , Eric I. Benchimol^{8,9,10} 



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

Vera Ehrenstein^{1*}, Maja Hellfritsch^{2,3}, Johnny Kahlert¹, Sinéad M. Langan⁴, Hisashi Urushihara⁵, Danica Marinac-Dabic⁶, Jennifer L. Lund⁷, Henrik Toft Sørensen¹, Eric I. Benchimol^{8,9,10}

¹Department of Clinical Epidemiology, Department of Clinical Medicine, Aarhus University and Aarhus University Hospital, 8200 Aarhus N, Denmark

²Research Unit of Clinical Pharmacology, Pharmacy, and Environmental Medicine, University of Southern Denmark, 5230 Odense M, Denmark

³Department of Cardiology, Gødstrup Hospital, 7400 Herning, Denmark

⁴Department of Non-communicable Disease Epidemiology, Faculty of Epidemiology and Population Health, London School of Hygiene and Tropical Medicine, London WC1E 7HT, United Kingdom

⁵Division of Drug Development and Regulatory Science, Faculty of Pharmacy, Keio University, Tokyo 105-8512, Japan

⁶Office of Clinical Evidence and Analysis, Center for Devices and Radiological Health, US Food and Drug Administration, Silver Spring, MD 20993, United States

⁷Department of Epidemiology, Gillings School of Global Public Health, University of North Carolina at Chapel Hill, Chapel Hill, NC 27599, United States

⁸Division of Gastroenterology, Hepatology and Nutrition and Child Health Evaluative Sciences, SickKids Research Institute, The Hospital for Sick Children, Toronto, ON M5G 1X8, Canada

⁹Department of Paediatrics and Institute of Health Policy, Management and Evaluation, University of Toronto, Toronto, ON M5T 3M6, Canada

¹⁰Institute for Clinical Evaluative Sciences, Toronto, ON M4N 3M5, Canada

*Corresponding author: Vera Ehrenstein, Department of Clinical Epidemiology, Department of Clinical Medicine, Aarhus University and Aarhus University Hospital, Olof Palmes Allé 43-45, 8200 Aarhus N, Denmark (ve@clin.au.dk)

Abstract

Clinicians, researchers, regulators, and other decision-makers increasingly rely on evidence from real-world data (RWD), including data routinely accumulating in health and administrative databases. RWD studies often rely on algorithms to operationalize variable definitions. An algorithm is a combination of codes or concepts used to identify persons with a specific health condition or characteristic. Establishing the validity of algorithms is a prerequisite for generating valid study findings that can ultimately inform evidence-based health care. In this paper, we aim to systematize terminology, methods, and practical considerations relevant to the conduct of validation studies of RWD-based algorithms. We discuss measures of algorithm accuracy, gold/reference standards, study size, prioritization of accuracy measures, algorithm portability, and implications for interpretation. Information bias is common in epidemiologic studies, underscoring the importance of transparency in decisions regarding choice and prioritizing measures of algorithm validity. The validity of an algorithm should be judged in the context of a data source, and one size does not fit all. Prioritizing validity measures within a given data source depends on the role of a given variable in the analysis (eligibility criterion, exposure, outcome, or covariate). Validation work should be part of routine maintenance of RWD sources.

This article is part of a Special Collection on Pharmacoepidemiology.

Key words: algorithms; data quality; information bias; measurement error; misclassification; routinely collected health data; real-world data; validity.

Introduction

Real-world data (RWD) and real-world evidence (RWE) originate outside phase I–III (preapproval) interventional trials.^{1,2} Routinely collected health data (RCD) accrue as a by-product of health-care delivery and encompass electronic health records; health administrative data; claims; electronic records of diagnoses, treatments, procedures, or devices; and patient-reported events.³ The role of RCD-based RWE is increasing in decision-making by clinicians, regulators, and health-care payers.^{1,2,4–11} This development predicts an ever-strengthening role of RWE in research on disease prevention, etiology, and clinical course.¹²

Advantages of RWD include accrual independent of research questions, population-wide coverage, sustainability, structured format, and standard coding.^{1,12} Measurement error, inevitably introduced along the path from the point of care to an analytical data point, may translate into systematic errors (or information

bias) on the level of study results and, ultimately, to misleading clinical practice.¹³ Large-scale RWD-based studies yielding precise but potentially biased results may exacerbate the fallacy of mistaking precision for validity.¹⁴ Validation studies of RWD and reporting guidelines have helped improve transparency^{15–19}; however, there remains an unmet need for education and awareness about the conduct of validation studies in RWD-based research.^{12,20–27} A validation study aims to quantify the validity of a measurement instrument. In (pharmaco)epidemiologic studies, operational definitions of the study variables—exposure, outcomes, covariates—are the measurement instruments (commonly called algorithms) whose validity is estimated against a chosen standard.²⁶

This paper outlines universal principles of validation “best practice,” including the basic terminology, the underlying theory, and the practical approaches. After defining an algorithm,

× Bookmarks

✓ Validation of algorithms in studies based on routinely collected health data: general principles

Introduction

Algorithms in routinely collected data

Completeness of a data source

Measures of algorithm validity

Gold standard and reference standard

Which measures can be estimated in a validation study?

Prioritizing validity measures

Portability of algorithms

Size of a validation study

Interpretation of study results and consideration of further approaches

Summary

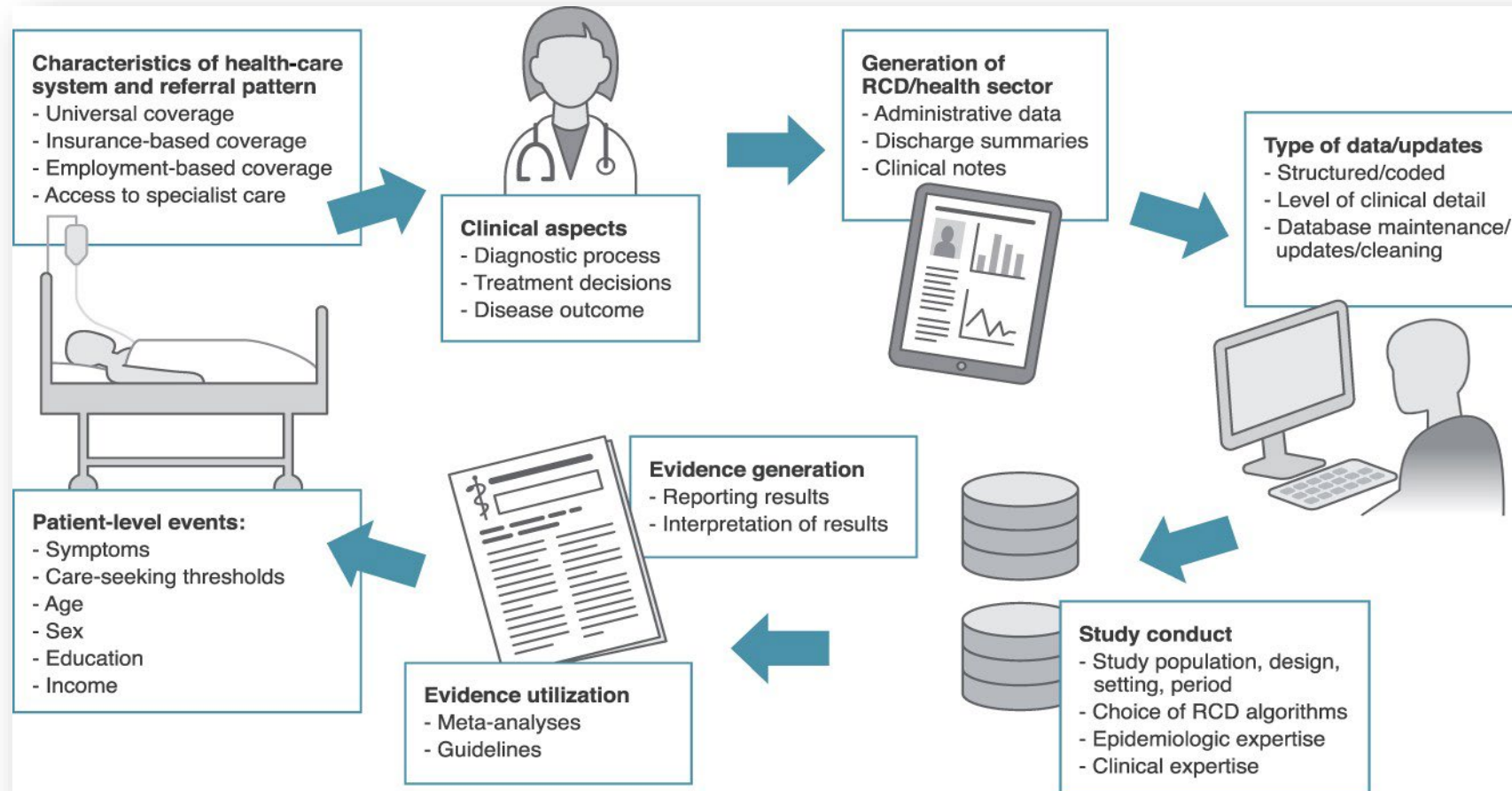
Acknowledgments

Supplementary material

Funding

Conflict of interest

Data availability



THE DIAGNOSIS PARADIGM



ALGORITHM/PHENOTYPE: BACTERIAL INFECTION

Algorithm	Health state of interest
Dispensing with ATC J01 (antibacterials for systemic use)	Bacterial infection
Dispensing with ATC J01 (antibacterials for systemic use)	Exposure to a systemic antibacterial agent

International Journal of Obesity (2015) 39, 1450–1455
© 2015 Macmillan Publishers Limited All rights reserved 0307-0565/15
www.nature.com/ijob

PEDIATRIC ORIGINAL ARTICLE

Prenatal exposure to systemic antibacterials and overweight and obesity in Danish schoolchildren: a prevalence study

A Mor¹, S Antonsen¹, J Kahlert¹, V Holsteen², S Jørgensen², J Holm-Pedersen², HT Sørensen¹, O Pedersen³ and V Ehrenstein¹

BACKGROUND/OBJECTIVE: Prenatal exposure to antibacterials may permanently dysregulate fetal metabolic patterns via epigenetic pathways or by altering maternal microbiota. We examined the association of prenatal exposure to systemic antibacterials with overweight and obesity in schoolchildren.

SUBJECTS/METHODS: We conducted a prevalence study among Danish schoolchildren aged 7–16 years using data from routine school anthropometric evaluations conducted during 2002–2013. Prenatal exposure to antibacterials was ascertained by using maternal prescription dispensations and infection-related hospital admissions during pregnancy. We defined overweight and obesity among the children using standard age- and sex-specific cutoffs. We computed sex-specific adjusted prevalence ratios (aPRs) of overweight and obesity associated with exposure to prenatal antibacterials, adjusting for maternal age at delivery, marital status, smoking in pregnancy and multiple gestation; we also stratified the analyses by birth weight.

RESULTS: Among 9886 schoolchildren, 3280 (33%) had prenatal exposure to antibacterials. aPRs associated with the exposure were 1.26 (95% confidence interval (CI): 1.10–1.45) for overweight and 1.29 (95% CI: 1.03–1.62) for obesity. Among girls, aPRs were 1.16 (95% CI: 0.95–1.42) for overweight and 1.27 (95% CI: 0.89 to 1.82) for obesity. Among boys, aPRs were 1.37 (95% CI: 1.13–1.66) for overweight and 1.29 (95% CI: 0.96–1.73) for obesity. The aPR for overweight was higher among schoolchildren with birth weight < 3500 g (aPR: 1.30, 95% CI: 1.05–1.61) than in schoolchildren with birth weight ≥ 3500 g (aPR: 1.18, 95% CI: 0.95–1.46). Inversely, the association for obesity was higher among schoolchildren with birth weight ≥ 3500 g (aPR: 1.35, 95% CI: 1.00–1.81) than among those who were < 3500 g at birth (aPR: 1.16, 95% CI: 0.82–1.65).

CONCLUSIONS: Prenatal exposure to systemic antibacterials is associated with an increased risk of overweight and obesity at school age, and this association varies by birth weight.

International Journal of Obesity (2015) 39, 1450–1455; doi:10.1038/ijob.2015.129

ALGORITHM/PHENOTYPE: SERIOUS INFECTION

	A	B	C	D	E	F	G	H	I	J	K
1	InfectionType	ICD10	ICD10Text	Exclusion	ExclusionText	Notes	DiagnosesTypes	PatientType			
2	Intra-abdominal infe	A00	Cholera				C_ADIAG	0 (Inpatient 1+ overnight stay)			
3	Intra-abdominal infe	A01	Typhoid and paratyphoid fevers				C_ADIAG	0 (Inpatient 1+ overnight stay)			
4	Intra-abdominal infe	A02	Other salmonella infections	A021, A022C	Excluding A021 Salmonella sepsis, A022C S	Included in other types	C_ADIAG	0 (Inpatient 1+ overnight stay)			
5	Sepsis	A021	Salmonella sepsis				C_ADIAG	0 (Inpatient 1+ overnight stay)			
6	Infections of CNS	A022C	Salmonella meningitis				C_ADIAG	0 (Inpatient 1+ overnight stay)			
7	Intra-abdominal infe	A03	Shigellosis				C_ADIAG	0 (Inpatient 1+ overnight stay)			
8	Intra-abdominal infe	A04	Other bacterial intestinal infections				C_ADIAG	0 (Inpatient 1+ overnight stay)			
9	Intra-abdominal infe	A05	Other bacterial foodborne intoxications, not elsewhere classified				C_ADIAG	0 (Inpatient 1+ overnight stay)			
10	Intra-abdominal infe	A06	Amoebiasis				C_ADIAG	0 (Inpatient 1+ overnight stay)			
11	Intra-abdominal infe	A07	Other protozoal intestinal diseases				C_ADIAG	0 (Inpatient 1+ overnight stay)			
12	Intra-abdominal infe	A08	Viral and other specified intestinal infections				C_ADIAG	0 (Inpatient 1+ overnight stay)			
13	Intra-abdominal infe	A09	Other gastroenteritis and colitis of infectious and unspecified origin				C_ADIAG	0 (Inpatient 1+ overnight stay)			
14	Infections of CNS	A170	Tuberculous meningitis				C_ADIAG	0 (Inpatient 1+ overnight stay)			
15	Other	A20	Plague	A203	Exclude A203 Plague meningitis	Included in other types	C_ADIAG	0 (Inpatient 1+ overnight stay)			
16	Infections of CNS	A203	Plague meningitis				C_ADIAG	0 (Inpatient 1+ overnight stay)			
17	Other	A21	Tularaemia				C_ADIAG	0 (Inpatient 1+ overnight stay)			
18	Other	A22	Anthrax	A227	Exclude A227 Anthrax sepsis	Included in other types	C_ADIAG	0 (Inpatient 1+ overnight stay)			
19	Sepsis	A227	Anthrax sepsis				C_ADIAG	0 (Inpatient 1+ overnight stay)			
20	Other	A23	Brucellosis				C_ADIAG	0 (Inpatient 1+ overnight stay)			
21	Other	A24	Glanders and melioidosis				C_ADIAG	0 (Inpatient 1+ overnight stay)			
22	Other	A25	Rat-bite fevers				C_ADIAG	0 (Inpatient 1+ overnight stay)			
23	Other	A26	Erysipeloid				C_ADIAG	0 (Inpatient 1+ overnight stay)			
24	Other	A27	Leptospirosis				C_ADIAG	0 (Inpatient 1+ overnight stay)			
25	Other	A28	Other zoonotic bacterial diseases, not elsewhere	A282B		Included in other types	C_ADIAG	0 (Inpatient 1+ overnight stay)			
26	Sepsis	A282B	Yersinia sepsis				C_ADIAG	0 (Inpatient 1+ overnight stay)			
27	Other	A30	Leprosy [Hansen disease]				C_ADIAG	0 (Inpatient 1+ overnight stay)			
28	Other	A31	Infection due to other mycobacteria				C_ADIAG	0 (Inpatient 1+ overnight stay)			

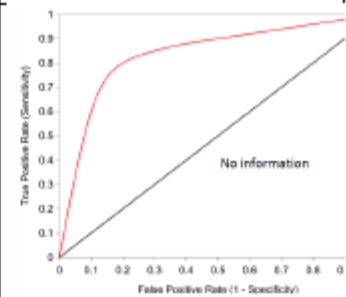
Table 1. Glossary of important epidemiologic terms used in this paper.

Term	Definition/elaboration
Algorithm	An algorithm is defined as a combination of people with a specific health condition, health state, used to classify persons with algorithm may be as simple as a single meeting algorithm criteria are classified
Alloyed gold standard	See <i>reference standard</i> .
Area under the curve	Area under the ROC curve, used as a single values of a continuous variable. See <i>Table S2</i> for details.
Completeness	The completeness of a data source is the population captured in that data source
Confounder	A confounder can be conceptualized as a development. Thus, a confounder cannot study, a characteristic that fulfills the exposure categories and (2) predicts the
Covariate	A characteristic of members of the study population is often used for predictors of the associated with the exposure of interest
Diagnostic odds ratio	The ratio of the odds of algorithm positivity those without the health state of interest
Exposure	A health state that, in a given study, explicit treatment allocation in randomized trial status (medicinal, surgical, device), but socioeconomic status). Conceptually, exposure
External validation	Validation of events of a health state identified in an ideal sense, a method that classifies standards rarely exist (see <i>reference standard</i>) gold standards may be time-dependent methods previously considered gold standard
Health state	A generic term used in this paper for any health state may play the role of exposure
Information bias	Discrepancy between the true value of a given value. Errors may be inherent in a measurement (eg, data entry errors). <i>Information bias</i> is
Internal validation	Validation of events of a health state identified in a given study. See <i>information bias</i> .
Measurement error	Incorrect classification by an algorithm of a body mass index in the normal range or weight).
Misclassification	Proportion of patients who are truly negative algorithm. See <i>Table S2</i> for details.
Negative predictive value	A health state that, in a given study, plays a role in the outcome of treatment (benefit or risk), (eg, infarction), or death. Conceptually, exposure
Outcome	Proportion of patients who are truly positive algorithm. See <i>Table S2</i> for details.
Positive predictive value	Data that originate in the course of routine clinical practice, usually by contrast with data originating in phase I-III clinical trials.
Real-world data	Evidence generated using RWD.
Real-world evidence	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). See <i>Table S2</i> for details.
ROC curve	A method or data source that is expected to classify individuals with respect to a health state of interest better than the algorithm being validated. Alternative term: <i>alloyed gold standard</i> .
Reference standard	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 nonhealth characteristics such as income, education, and employment, if they are important in studying health states. Routinely collected health data are a type of RWD.
Routinely collected health data	Proportion of persons with a given health state according to a gold/reference standard who are classified as positive by an (RWD) algorithm. See <i>Table S2</i> for details.
Sensitivity	Proportion of persons without a given health state according to a gold/reference standard who are classified as negative by an (RWD) algorithm. See <i>Table S2</i> for details.
Specificity	Study population used to evaluate measures of algorithm validity.
Validation cohort	Correspondence between an estimate and the true value of a parameter. In epidemiologic studies, a valid (unbiased) estimate may estimate the occurrence of a health event, or an association between exposure and outcome.
Validity (internal)	

Table S1. Examples of RWD-based algorithms and validation studies of RWD-based algorithms

Health state or event being validated	Country	RWD source(s)	RWD-based algorithm	Gold/reference standard
In-hospital chemotherapy treatment	Denmark	Danish National Patient Registry	Agent-specific treatment codes	In-hospital pharmacy production system and chart review.
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.
Cardiac interventions	Denmark	Danish National Patient Registry	Surgical and procedure codes according to the NSCP classification	Chart review.
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.
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.
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.
Hepatic encephalopathy	United States	Administrative claims	ICD-10 codes and treatment proxies after switch to ICD-10 in 2015 rendered algorithms	Independent prospectively acquired sample containing both true-positive and

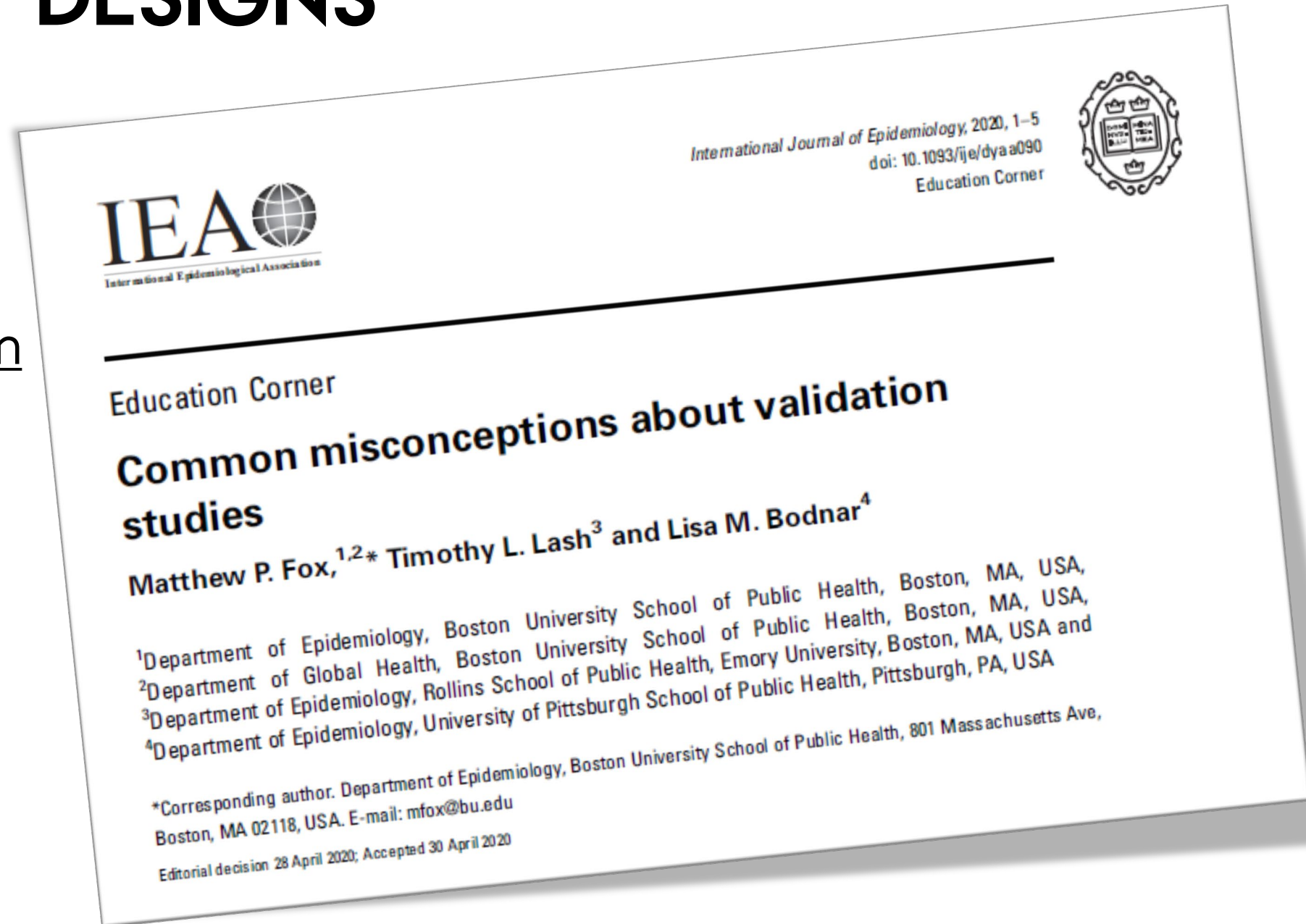
Table S2. Measures of algorithm validity

Measure	Range	Description	Formulae
Sensitivity	[0, 1]	Proportion of persons with a given health state according to gold/reference standard who are classified as such by an algorithm.	Health state status based on gold/reference standard: Positive 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.	Health state status based on algorithm: Positive Negative
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.	$Sensitivity = \frac{A}{(A+B)}$ $Specificity = \frac{D}{(C+D)}$ $PPV = \frac{A}{(A+C)}$ $NPV = \frac{D}{(B+D)}$
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. Chaback et al and Benschimol et al. provide formulae connecting sensitivity and specificity, PPV and NPV with prevalence of a health state. ¹¹	
Receiver operating characteristic (ROC) curve		Definition of an RWD-based algorithm may be based on a cut-off value of a measured continuous variable (eg, hypotension based on serum potassium levels in the study by Holland-Hill et al ⁹). 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) an ROC is a plot of sensitivity (or true-positive proportion) on the y-axis against (1-specificity) or false-positive proportion on the x-axis. For a (hypothetical) perfect algorithm, sensitivity=specificity=1.	
Area under the curve (AUC)		AUC is derived from the 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 no better than a coin toss would have an AUC=0.5.	Example of a receiver operating characteristic (ROC) curve and area under the curve (AUC). (Sorenson and Vandembroucke (in press). ⁷)
Diagnostic odds ratio (DOR)	[0, ∞)	An odds ratio of dichotomous tests in diagnostic applications 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. ⁸ 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 corresponds to higher probability of an algorithm to be positive among true cases than in non-cases of a given health state.	$DOR = \frac{A \text{ (True Positive)}}{B \text{ (False Negative)}} \div \frac{C \text{ (False Positive)}}{D \text{ (True Negative)}}$
Kappa statistic	[-1, 1]	Kappa statistic is used to quantify inter-rater variability. ¹² In evaluating an RWD-based algorithm, it can be used to quantify agreement between two algorithms or data sources, none of which can be considered better (a reference standard) relative to the other. Based on the value of kappa statistic, agreement may be qualified on a scale 'less than chance' to 'almost perfect' (Viers and Garrett 2005 ¹²). Kappa Agreement	Health state status according to data source 1: Positive Negative Health state status according to data source 2: Positive Negative Expected agreement between two data sources/algorithms, none of which is superior to the other in classifying a given health state is to present/absent

< 0	Less than chance	Observed agreement $P_{obs} = \frac{(a+d)}{n}$ Expected agreement $P_{exp} = \left[\left(\frac{a}{n} \right) \times \left(\frac{a}{n} \right) \right] + \left[\left(\frac{d}{n} \right) \times \left(\frac{d}{n} \right) \right]$ Kappa statistic $\kappa = \frac{(P_{obs} - P_{exp})}{(1 - P_{exp})}$
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	

VALIDATION 'DESIGNS'

- 1: select on algorithm
- 2: select on truth
- 3: take them all



SCENARIO 1: SELECT ON ALGORITHM

- Select **algorithm-positive** (and negative) persons
- Estimate **PPV** (and NPV)
 - Any component
 - Each component

	Event truly present	Event truly absent	
Algorithm+	a True Positive	b False Positive	a + b
Algorithm-	c False Negative	d True Negative	c + d
	a + c	b + d	

Positive predictive value = $a/(a+c)$

SCENARIO 2: SELECT ON GOLD STANDARD

- Identify an independent sample of **gold-standard-positive events** (and nonevents)
- Estimate **sensitivity** (and specificity)
 - Any component
 - Each component

	Event truly present	Event truly absent	
Algorithm+	a True Positive	b False Positive	a + b
Algorithm-	c False Negative	d True Negative	c + d
	a + c	b + d	

Sensitivity = $a/(a+c)$

SCENARIO 3: SELECT INDEPENDENTLY OF BOTH

- Select population regardless of algorithm/gold standard status
- Estimate **sensitivity, specificity, PPV, NPV**
 - Any component
 - Each component

	Event truly present	Event truly absent	
Algorithm+	a True Positive	b False Positive	a + b
Algorithm-	c False Negative	d True Negative	c + d
	a + c	b + d	

GOLD STANDARD

- Perfect (better) measure of event of interest
- **Case vs. Non-case definition**
- Examples
 - Medical charts
 - Patient histories e-review
 - Another database



MEDICAL CHART REVIEW

- Ascertainment time windows
- Key record types required for identifying the gold/standard (eg, laboratory or diagnostic results)
- Key specialists/departments who should provide those records
- Combine or decouple chart abstraction and case adjudication
- Obtaining additional information, resolving disagreements, and quantifying interrater variability
- Staff skills and need for training in standardized chart abstraction (a pilot medical record review ahead of a large-scale effort will inform the process and help use resources efficiently)
- Whether/to what extent to adopt the event definitions from RCTs

Routinely collected data: the importance of high-quality diagnostic coding to research

Stuart G. Nicholls PhD, Sinéad M. Langan MSc PhD, Eric I. Benchimol MD PhD

■ Cite as: *CMAJ* 2017 August 21;189:E1054-5. doi: 10.1503/cmaj.170807

Follow this preprint

Trends of use of drugs with suggested shortages and their alternatives across 52 real world data sources and 18 countries in Europe and North America

[Marta Pineda-Moncusí](#), [Alexandros Rekkas](#), [Álvaro Martínez Pérez](#), [Angela Leis](#),
[Carlos Lopez Gomez](#), [Eric Fey](#), [Erwin Bruninx](#), [Filip Maljković](#), [Francisco Sánchez-Sáez](#),
[Jordi Rodeiro](#), [Loretta Zsuzsa Kiss](#), [Michael Franz](#), [Miguel-Angel Mayer](#), [Neva Eleangovan](#),
[Pericàs Pulido Pau](#), [Pantelis Natsiavas](#), [Selçuk Şen](#), [Steven Cooper](#), [Sulev Reisberg](#), [Katrin Manlik](#),
[Beatriz del Pino](#), [Albert Prats Uribe](#), [Ali Yağız Üresin](#), [Ana Danilović Bastić](#), [Ana Maria Rodrigues](#),
[Ângela Afonso](#), [Anna Palomar-Cros](#), [Annelies Verbiest](#), [Antonella Delmestri](#), [Barış Erdoğan](#),
[Carina Dinkel-Keuthage](#), [Carmen Olga Torre](#), [Caroline de Beukelaar](#), [Caroline Eteve-Pitsaer](#),
[Cátia F. Gonçalves](#), [Costantino de Palma](#), [Cristina Gavina](#), [Daniel Dedman](#), [David Brendan Price](#),
[Denisa Gabriela Balan](#), [Dirk Enders](#), [Edward Burn](#), [Elisa Henke](#), [Elyne Scheurwegs](#), [Emma Callewaert](#),
[Encarnación Pérez Martínez](#), [Eng Hooi Tan](#), [Fabian Prasser](#), [Francois Antonini](#), [Frank Staelens](#),
[Fredrik Nyberg](#), [Geoffray Agard](#), [Gianluigi Galli](#), [Gianmario Candore](#), [Gianny Mestdach](#), [Hadas Shachaf](#),
[Harri Rantala](#), [Huiqi Li](#), [Ines Reinecke](#), [Irene López-Sánchez](#), [Jaime E Poquet-Jornet](#),
[Javier de la Cruz Bertolo](#), [Jelle Evers](#), [João Firmino-Machado](#), [Jonas Wastesson](#),
[Juan Luis Cruz Bermúdez](#), [Juan Manuel Ramírez-Anguita](#), [Kimmo Porkka](#), [Kristina Johnell](#),
[Laurent Boyer](#), [Lieselot Cool](#), [Luca Moschetti](#), [Manon Merkelbach](#), [Mariana Canelas-Pais](#),
[Massimo Dominici](#), [Máté Szilcz](#), [Matteo Puntoni](#), [Mees Mosseveld](#), [Mina Tadrous](#), [Miquel Oltra-Sastre](#),
[Mona Bové](#), [Nadav Rappoport](#), [Noelia García Barrio](#), [Otto Ettala](#), [Paolo Baili](#),
[Paula Rubio Mayo](#), [Peter Prinsen](#), [Raeleesha Norris](#), [Ravinder Claire](#), [Reut Sherman Yackob](#),
[Roberto Lillini](#), [Salvador Garcia-Torrens](#), [Sampo Kukkurainen](#), [Silvia Lazzarelli](#), [Talita Duarte-Salles](#),
[Tiago Taveira-Gomes](#), [Tim Jansen](#), [Ulrich Keilholz](#), [Wai Yi Man](#), [Xintong Li](#), [Zsolt Bagyura](#),
[Daniel Prieto-Alhambra](#), [Peter R. Rijnbeek](#), [Theresa Burkard](#)

doi: <https://doi.org/10.1101/2024.08.28.24312695>

This article is a preprint and has not been peer-reviewed [what does this mean?]. It reports new medical research that has yet to be evaluated and so should not be used to guide clinical practice.

Increasing study size

- Reduces random error
- Does not reduce systematic error



Helping everyone do better: a call for validation studies of routinely recorded health data

Vera Ehrenstein¹
Irene Petersen^{1,2}
Liam Smeeth³
Susan S Jick⁴
Eric I Benchimol^{5,6}
Jonas F Ludvigsson^{7,8}
Henrik Toft Sørensen¹

This article was
Clinical Epidemiol
12 April 2016
Number of tim

There has
electronic
and diseas
and Drug
of validate
focusing o
papers of t

EPIDEMIOLOGY

Articles & Issues ▾ Collections Multimedia ▾ For Authors ▾ Journal Info ▾

[Next Article >](#)

EPIDEMIOLOGY Announces the “Validation Study” Submission Category

Lash, Timothy L.; Olshan, Andrew F.

Epidemiology: September 2016 - Volume 27 - Issue 5 - p 613–614
doi: 10.1097/EDE.0000000000000532
Commentary


Received: 7 January 2018 | Revised: 5 October 2018 | Accepted: 11 October 2018

DOI: 10.1002/pds.4694

COMMENTARY

WILEY

Pharmacoepidemiology and Drug Safety's special issue on validation studies

Danielle S. Chun  | Jennifer L. Lund  | Til Stürmer 

Epidemiology, University of North Carolina at Chapel Hill Gillings School of Global Public Health, Chapel Hill, North Carolina, USA

Correspondence

D. S. Chun, Epidemiology, University of North Carolina at Chapel Hill Gillings School of Global Public Health, Chapel Hill, NC, USA
Email: danielle.chun@unc.edu



..research accomplishes more good and brings greater health to humanity than more intensive doctoring would ever accomplish...

NOTES FROM BEYOND

My Interview With John Snow

Kenneth J. Rothman

Epidemiology. 2004;15



A huge responsibility!



CONSIDERATIONS FOR RESEARCHER

Are there validated algorithms for your variables/database?

- If yes, is their validity satisfactory given
 - Study aims
 - Variable roles

If no, is a bespoke validation study possible (resources)?

If resources are limited, what is the best solution?





AARHUS
UNIVERSITY