

Clinical epidemiology in the era of big data: new opportunities, familiar challenges

Vera Ehrenstein
Henrik Nielsen
Alma B Pedersen
Søren P Johnsen
Lars Pedersen

Department of Clinical Epidemiology,
Aarhus University Hospital,
Aarhus N, Denmark

Abstract: Routinely recorded health data have evolved from mere by-products of health care delivery or billing into a powerful research tool for studying and improving patient care through clinical epidemiologic research. Big data in the context of epidemiologic research means large interlinkable data sets within a single country or networks of multinational databases. Several Nordic, European, and other multinational collaborations are now well established. Advantages of big data for clinical epidemiology include improved precision of estimates, which is especially important for reassuring (“null”) findings; ability to conduct meaningful analyses in subgroup of patients; and rapid detection of safety signals. Big data will also provide new possibilities for research by enabling access to linked information from biobanks, electronic medical records, patient-reported outcome measures, automatic and semiautomatic electronic monitoring devices, and social media. The sheer amount of data, however, does not eliminate and may even amplify systematic error. Therefore, methodologies addressing systematic error, clinical knowledge, and underlying hypotheses are more important than ever to ensure that the signal is discernable behind the noise.

Keywords: electronic health records, healthcare administrative claims, medical record linkage, multicenter studies, validation studies

Introduction

Big data has firmly established itself in the health research,^{1,2} illustrated by publications in high-ranking general-interest biomedical journals, including *The New England Journal of Medicine*,³ *JAMA*,⁴ *Journal of Internal Medicine*,⁵ *Science*,^{6–9} and *Nature*.^{10–13} A basic definition of big data includes “the 3 Vs”: variety (linkage of many data sets from heterogeneous independent sources in a single data set); volume (large number of observations and variables per observation from different sources); and/or velocity (real-time or frequent data updates, often fully or partially automated).¹⁴ Other definitions encompass additional three Vs: value (clinically relevant information); variability (eg, seasonal or secular disease trends); and veracity (data quality).² Routinely recorded health data are large automated data sets stemming from day-to-day activities of health care, such as hospital admissions or claims.^{15–18} These data have evolved from mere by-products of health care delivery or billing into a powerful tool for improving patient care through preventive, etiologic, and prognostic epidemiologic research.⁴ A recent article summarizes 46 most influential studies conducted with big data in health care,¹ while a review from 2015 provides multiple examples of the “variety” V in big data for health.²

The notion of applying lessons from the clinical past to the clinical future is “as old as medicine.”¹⁹ In a simplified form, evidence-based medical care means that a

Correspondence: Vera Ehrenstein
Department of Clinical Epidemiology,
Aarhus University Hospital, Olof Palmes
Allé 43-45, 8200 Aarhus N, Denmark
Tel +45 8716 8063
Fax +45 8716 7215
Email ve@clin.au.dk

clinician can use research results in making treatment decisions in his or her clinical practice, often through explicit literature-based treatment guidelines. For a clinician, this means answers to questions such as: “How likely is my patient with atrial fibrillation on oral anticoagulants to develop a major bleeding? Does the risk vary by type of anticoagulant or patient characteristics?” or “To what extent does comorbidity affect mortality of patients with hip fracture?” To be answered, a clinical question must be first translated into a precise research question and then back-translated and interpreted for clinical decision making. Therefore, it is essential for clinicians and epidemiologists to understand each other’s language. For an epidemiologist, an answer to a research question should be a precise and valid estimate of an underlying population parameter such as mean, risk, incidence rate, or odds ratio. Big data – via the “volume” V – often addresses the precision component, but does little to address validity (the “veracity” V in the big-data vocabulary). Plausible hypotheses, expert knowledge, and accurate measurement tools must be available to ensure validity of research findings, since a highly precise biased result, especially perceived as credible based on precision alone, is more dangerous translated into clinical practice than an imprecise biased result.^{20,21} This paper, using primarily case studies from the Nordic countries, provides a brief overview and examples of use of big data in clinical epidemiology and outlines associated advantages and challenges.

Examples of big data collaborations in epidemiology

Some say that the digitalization of medical records revolutionized the usability of big data in medical research.⁴ Whether or not this claim is accepted, it is important to be aware that the current development follows a long evolution of using register data for medical research. This evolution started with the establishment of the first National Leprosy Register, in Norway, in 1856 (Figure 1),^{22,23} and of the Danish Cancer Registry, in 1943.²⁴ Other Nordic registries followed, most of them established between the 1960s and the early 2000s.^{25,26} Researchers in the Nordic countries have been using the volume component of the big data before the term was invented: for decades, epidemiologists have been conducting epidemiologic studies based on linkage of routinely collected data from multiple administrative, health, and demographic registries, and their potential has been recognized at least since the 1990s,²⁷ if not earlier.²⁸

Estimates of association with narrow confidence intervals often stem from big data analyses of common health



Figure 1 Building that used to house the Norwegian Leprosy Registry, currently home of the Department of Global Public Health and Primary Care, University of Bergen, Norway.

Note: Courtesy: Dr Astrid Lunde.

outcomes in population-based registry data spanning several decades. When the intervention or the outcome of interest is rare, even data from an entire country may be insufficient, requiring that data from different countries are combined. Several formal or ad hoc collaborative networks in observational epidemiology have arisen, often from the need to study benefits and risks of relatively uncommon pharmacological^{16,29–31} or surgical^{32,33} interventions, or vaccines.^{3,30} Examples of pan-Nordic collaborations using combined data from Denmark, Finland, Iceland, Norway, and Sweden^{31,34,35} include studies on prenatal exposure to antidepressants and adverse effects in the offspring^{31,34,35} or the Nordic Arthroplasty Register Association (NARA) database of about 1 million primary hip and knee replacement procedures performed since 1995 in Denmark, Finland, Norway, and Sweden.³⁶ NARA enabled studies of rare risk factors and outcomes, for which single-country data are too sparse.^{32,33} One clinically relevant question is whether a type of fixation used in total hip replacement (THR) is associated with risk of subsequent revision in patients younger than 55 years of age, since these patients may be different from older patients in mobility, post-THR life expectancy, and compliance with treatment. Only 5% of THR procedures are performed in patients younger than 55 years and previous studies, including those based on national hip registries, had insufficient sample size to address the fixation issue in younger patients. Pedersen et al³⁷ used NARA to assemble a study population of ~30,000 patients younger than 55 years undergoing THR, with each fixation technique represented by more than 3,000 observations. The study yielded a clinically relevant message that uncemented implants are associated with a lower long-term risk of aseptic

loosening but a higher short-term risk of revisions. Thus, the purpose of uncemented implants has been achieved in long term, but technical issues causing dislocation, periprosthetic fracture, and infection have been previously overlooked in patients younger than 55 years.

Use of routinely collected data for epidemiologic research has also been possible outside the Nordic countries, including general practice-based data in the UK, or claims-based databases and database networks in the USA. In contrast to the typical European health care databases, which are established to fulfill administrative (health services), clinical quality, or surveillance needs, the US claims databases (eg, Medicare, Medicaid, and commercial insurance records) are by-products of medical accounting. Several European database networks, including those encompassing the Nordic data, have been successfully established and have found ways to overcome challenges of differences in the underlying health care systems, languages, data-sharing laws, record-generating mechanisms, and classifications.^{5,16,30,38,39} Medical data in the Nordic countries are coded using a common basic set of standard classifications (International Classification of Diseases, Nordic Medico-Statistical Committee classification for procedures and causes of injury,^{40,41} or Anatomical Therapeutic Chemical codes for medications), which makes it easier to establish common algorithms. In the USA, Medicare and Medicaid provide financial incentives for “meaningful use” of electronic health records.³ The most prominent big data collaborative models in the USA have been the Mini-Sentinel project and the Observational Medical Outcomes Partnership (OMOP).³ The difference between routine records accumulated in systems like Mini-Sentinel or OMOP and those in Europe is the structure of the health care system, linkage possibilities, and the availability of lifelong complete follow-up. Thus, certain aspects of big data in Nordic countries are more diverse than those in many other databases (the “volume” V and the “variety” V of the big data), thanks to individual-level linkage to both medical and nonmedical data, including education, income, and residence, and because of lifelong follow-up. In 2013, the Mini-Sentinel project covered 360 million person-years of observation representing 150 million lives.³ In 2014, the Danish Civil Registration System, with its linkable network of national registries, covered 400 million person-years of observation from 9.5 million lives.²⁵ Asian countries are building a linkable registry infrastructure with individual-level linkage mimicking those of the Nordic countries.⁴²

The “variety” V of the big data is developing rapidly, whereby previously unused or underused types of data are

incorporated into medical research, including electronic medical records, imaging, biobanks, and patient-reported data (including social media and wearables).^{2,43} Individual linkage may not be always necessary: in a classical ecologic study, hostility of language on Twitter was associated with country-specific mortality from heart diseases.⁴⁴ Pharmacovigilance with social media is already a reality.⁴⁵ Mobile phones can be used to test and subsequently deliver behavioral interventions such as smoking cessation aid⁴⁶ or adherence support.⁴⁷ The type of bias associated with certain types of data may change over time. For example, in the early days of epidemiologic research, random landline phone surveys tended to select the relatively more affluent, the employed, and the young. Today, these groups are more likely to be accessed via social networks and mobile telephony,² while use of landline phones may select for older or disadvantaged population segments.

Assembling database networks carries with it technical, logistical, ethical, and legal challenges.⁴⁸ The last two are often the hardest to overcome because of issues of data access, patient privacy, and potential conflicts of interest. Even in large studies, one has to remain vigilant about patient privacy and the possibility of inadvertently identifying individuals based on a set of rare characteristics. Gini et al¹⁶ provide a practical guide of the different models of data networking, defined on the degree of centralization and harmonization of the different analytic processes. It seems to be practical to designate a single network partner, with adequate resources, to be the coordinating analytic hub. The process starts with raw data from each participating database and ends with the statistical output combining results of individual patients from all databases. Between the starting and the end points, there exist different models for the extent of process automation, autonomy, and control enjoyed by each data partner. A global protocol, with flexibility for local adaptations, is usually followed. Depending on the aims of the study, the analysis may entail as little sharing as contributing country-specific odds ratios for a meta-analysis or as much sharing as harmonization and pooling of individual-level data sets.¹⁶ Harmonization involves transformations, whereby each partner creates standard input data sets according to exact specification – a common data model (CDM) – which dictates the data set types and structure, variable names and attributes, and definitions of derived variables. A single statistical analytic program is then run on the CDM-conforming files either by each network partner locally (“one analyst, many outputs”) or centrally by the hub on the combined data set (“one analyst, one output”). By contrast, the “many

analysts, many outputs” approach is discouraged because it is prone to error and duplicates work. Whether one or many analysts, quality control of programming by another analyst is always necessary.

Health outcomes measured by health care professionals might differ from the outcomes subjectively experienced by patients, and the latter also affects the outcome of treatment. To fill this gap, patient-reported outcome measures (PROMs) are being used increasingly.⁴⁹ An example of incorporation of PROMs in a single-country setting, while capitalizing on unique data linkage capabilities common to the Nordic settings, includes the generic infrastructure for collecting PROM data, AmbuFlex, developed in Denmark by Hjollund et al.⁵⁰ The researchers have successfully implemented a flexible paper-based and electronic data collection on PROMs in more than 20 projects since 2004. Group-level aggregated PROM data, linked with data from routine registries and clinical databases, can be used to monitor national and regional hospital performance in oncology and cardiology care, psychiatry, neurology, and orthopedics. Patient-level PROM data collected on clinic level, in combination with electronic health records, can be used to facilitate screening, clinical decisions, patient–doctor communication, and efficient use of resources in cardiology, rheumatology, and oncology. Response rates exceeded 75% in all and 90% in most cases. A clinical decision support function of PROMs can save clinicians’ time by using an algorithm-based initial identification of patients in need of immediate attention, while presenting data on other patients in a decision-supporting format for clinical judgment.⁵⁰ AmbuFlex is a unique example of implementation in routine care, a generic system integrated with electronic medical records, and is used for longitudinal collection of detailed PROM data on an individual level to personalize the care for the individual patient. This allows the collection of PROM data on large cohorts of chronically ill patients over many years, similar to the systems currently in place for administrative data.

Big data in epidemiology: benefits and challenges

Precision of results is not the only benefit of big data. Observations from large number of individuals allow a rapid detection of potential risk signals associated with newly marketed therapies, for which risks of rare adverse events are rarely known from Phase III preapproval trials (the velocity “V” of the big data).⁵¹ A thought experiment showed that having records of 100 million patients for safety monitoring would have allowed the detection of adverse cardiovascular effects

of rofecoxib (Merck, Kenilworth, NJ, USA) in 3 months instead of 5 years.^{5,52} On the other hand, large data sets help convincingly rule out harmful associations, in the so-called “null studies.” One example is the abovementioned Nordic collaboration on safety of antidepressant use in pregnancy. Less than 2% of pregnant women use selective serotonin reuptake inhibitors (SSRIs) in pregnancy, while birth defects affect about 3% of live births. Therefore it took a pan-Nordic study to assemble a study population of >1.5 million pregnancies with ~73,000 malformation cases, including ~33,000 SSRI-exposed pregnancies with >1,300 cases exposed to SSRIs.³⁴ The study convincingly showed a null association between maternal use of SSRIs and major birth defects, providing reassurance to pregnant women with depression and their physicians. Finally, in analyses based on large data sets, estimates are likely to be “highly statistically significant,” ie, associated with P -values <0.05. This “universal statistical significance” could finally lay to rest reliance on P -values for interpretation of study results, allowing researchers to focus on clinical significance instead.^{53–55}

The perks of big data should not go to our collective heads. Big data does not address the usual epidemiologic challenges related to validity, and may even amplify them.^{15,56} Accurate measurement of study variables remains imperative in big-data settings. An advantage of multinational databases is that estimates originating from different databases to address the same research question amount to reproducibility checks of results under varying assumptions about the record-generating mechanisms and the effects of the underlying health care and social structures. At the same time, in multinational database studies, validity concerns are increased proportional to the number of the databases, with the need of several valid operational definitions for the same clinical characteristic or event, to avoid propagating a systematic error on a large scale.^{53,56} Validation of algorithms in large secondary databases remains imperative for valid inference.^{15,56,57} The NARA collaboration has contributed to improvement of data validity in all four participating countries through regular meetings, where differences in registration practice have been discussed. Also, through different research projects, a number of differences regarding data quality between registries have been pointed out and discussed, and subsequently changes in national registries have been made to achieve uniform data definition, collection, and interpretation.

Large amounts of missing data may cause selection bias and undermine gains in precision afforded by big data, since in multiple regression models, standard statistical software

removes observations with missing values. Reverse causation, immortal time bias,⁵⁸ and healthy user/healthy adherer bias⁵⁹ are likewise not remedied by large amounts of data and need to be addressed in big-data and small-data studies alike. On a pragmatic level, delay of data delivery and changes in coding practice present additional challenges.

Conclusion

Epidemiologic research, including database research, is an “exercise in measurement,”⁶⁰ in an effort to maximize signal-to-noise ratio. The results of big data-based medical research represent a dividend to the public on its investment in the form of contribution to routine databases with data and with tax money. The advantages of big data are precision of results, including precise “null” findings, ability to address clinical questions in patient subgroups, and rapid detection of risk signals. In the Nordic countries, big data is collected and maintained by public institutions and operate in the setting of income-independent access to health care and lifelong follow-up. In other settings, such as US claims databases, demographic or economic disadvantages are better represented, while follow-up is not lifelong and health care access may be interrupted. Combining evidence from different settings and countries creates multiple-informant settings, providing built-in cross-validation and addressing a wide array of clinical questions in a single study. A formal requirement to the big data is that size, complexity, and velocity of the data are too intense for processing and interpretation with existing tools. In the Nordic settings, the volume has been available for some decades, and the variety is increasing rapidly to include data on imaging, behavior, geo-location, ecology, genetics, and patient-reported outcomes. Velocity has not yet reached the real-time update stage, but it is improving, and its value is obvious. Veracity (familiar to epidemiologists as validity) needs to be assured before data can be interpreted. The large amount of data, thus, does not eliminate and may amplify sources of systematic error. To that end, technical expertise, clinical knowledge, and underlying hypotheses are more important than ever to ensure that the signal is not drowned out by noise.

Acknowledgments

We thank Professor Olaf M Dekkers for helpful comments on the early drafts of this manuscript and Dr Astrid Lunde for providing the photo for Figure 1. This paper was funded by the Program for Clinical Research Infrastructure established by the Lundbeck Foundation and the Novo Nordisk Foundation and administered by the Danish Regions.

Disclosure

The authors report no conflicts of interest in this work.

References

- de la Torre Diez I, Cosgaya HM, Garcia-Zapirain B, López-Coronado M. Big Data in health: a literature review from the year 2005. *J Med Syst.* 2016;40(9):209.
- Andreu-Perez J, Poon CC, Merrifield RD, Wong ST, Yang GZ. Big data for health. *IEEE J Biomed Health Inform.* 2015;19(4):1193–1208.
- Psaty BM, Breckenridge AM. Mini-Sentinel and regulatory science—big data rendered fit and functional. *N Engl J Med.* 2014;370(23):2165–2167.
- Murdoch TB, Detsky AS. The inevitable application of big data to health care. *JAMA.* 2013;309(13):1351–1352.
- Trifirò G, Coloma PM, Rijnbeek PR, et al. Combining multiple health-care databases for postmarketing drug and vaccine safety surveillance: why and how? *J Intern Med.* 2014;275(6):551–561.
- Broniatowski DA, Paul MJ, Dredze M. Twitter: big data opportunities. *Science.* 2014;345(6193):148.
- Fung IC, Tse ZT, Fu KW. Converting Big Data into public health. *Science.* 2015;347(6222):620.
- Khoury MJ, Ioannidis JP. Medicine. Big data meets public health. *Science.* 2014;346(6213):1054–1055.
- Lazer D, Kennedy R, King G, Vespignani A. Big data. The parable of Google Flu: traps in big data analysis. *Science.* 2014;343(6176):1203–1205.
- Reardon S. US big-data health network launches aspirin study. *Nature.* 2014;512(7512):18.
- Savage N. Bioinformatics: big data versus the big C. *Nature.* 2014;509(7502):S66–S67.
- Sejdić E. Medicine: adapt current tools for handling big data. *Nature.* 2014;507(7492):306.
- Wilson S. Data protection: big data held to privacy laws, too. *Nature.* 2015;519(7544):414.
- Baro E, Degoul S, Beuscart R, Chazard E. Toward a literature-driven definition of big data in healthcare. *Biomed Res Int.* 2015;2015:639021.
- Gange SJ, Golub ET. From smallpox to big data: the next 100 years of epidemiologic methods. *Am J Epidemiol.* 2016;183(5):423–426.
- Gini R, Schuemie M, Brown J, et al. Data extraction and management in networks of observational health care databases for scientific research: a comparison of EU-ADR, OMOP, Mini-Sentinel and MATRICE strategies. *EGEMS.* 2016;4(1):1189.
- Hernán MA, Savitz DA. From “big epidemiology” to “colossal epidemiology”: when all eggs are in one basket. *Epidemiology.* 2013;24(3):344–345.
- Toh S, Platt R. Is size the next big thing in epidemiology? *Epidemiology.* 2013;24(3):349–351.
- Last JM. What is “clinical epidemiology”? *J Public Health Policy.* 1988;9(2):159–163.
- Hernán MA, Robins JM. Using big data to emulate a target trial when a randomized trial is not available. *Am J Epidemiol.* 2016;183(8):758–764.
- Ioannidis JP. Why most published research findings are false. *PLoS Med.* 2005;2(8):e124.
- Irgens LM. The origin of registry-based medical research and care. *Acta Neurol Scand Suppl.* 2012;(195):4–6.
- Irgens LM, Bjerkedal T. Epidemiology of leprosy in Norway: the history of The National Leprosy Registry of Norway from 1856 until today. *Int J Epidemiol.* 1973;2(1):81–89.
- Gjerstorff ML. The Danish Cancer Registry. *Scand J Public Health.* 2011;39(7 Suppl):42–45.
- Schmidt M, Pedersen L, Sørensen HT. The Danish Civil Registration System as a tool in epidemiology. *Eur J Epidemiol.* 2014;29(8):541–549.
- Furu K, Wettermark B, Andersen M, Martikainen JE, Almarsdottir AB, Sørensen HT. The Nordic countries as a cohort for pharmacoepidemiologic research. *Basic Clin Pharmacol Toxicol.* 2010;106(2):86–94.

27. Sørensen HT. Regional administrative health registries as a resource in clinical epidemiology. A study of options, strengths, limitations and data quality provided with examples of use. *Int J Risk Saf Med.* 1997;10(1):1–22.
28. Baksaas I, Fugelli P, Halvorsen IK, Lunde PKM, Næss K. Prescription of hypotensives in general practice. *Eur J Clin Pharmacol.* 1978;14(5):309–317.
29. FitzHenry F, Resnic FS, Robbins SL, et al. Creating a common data model for comparative effectiveness with the observational medical outcomes partnership. *Appl Clin Inform.* 2015;6(3):536–547.
30. Avillach P, Coloma PM, Gini R, et al. Harmonization process for the identification of medical events in eight European healthcare databases: the experience from the EU-ADR project. *J Am Med Inform Assoc.* 2013;20(1):184–192.
31. Kieler H, Artama M, Engeland A, et al. Selective serotonin reuptake inhibitors during pregnancy and risk of persistent pulmonary hypertension in the newborn: population based cohort study from the five Nordic countries. *BMJ.* 2012;344:d8012.
32. Havelin LI, Fenstad AM, Salomonsson R, et al. The Nordic Arthroplasty Register Association: a unique collaboration between 3 national hip arthroplasty registries with 280,201 THRs. *Acta Orthop.* 2009;80(4):393–401.
33. Robertsson O, Bizjajeva S, Fenstad AM, et al. Knee arthroplasty in Denmark, Norway and Sweden. A pilot study from the Nordic Arthroplasty Register Association. *Acta Orthop.* 2010;81(1):82–89.
34. Furu K, Kieler H, Haglund B, et al. Selective serotonin reuptake inhibitors and venlafaxine in early pregnancy and risk of birth defects: population based cohort study and sibling design. *BMJ.* 2015;350:h1798.
35. Stephansson O, Kieler H, Haglund B, et al. Selective serotonin reuptake inhibitors during pregnancy and risk of stillbirth and infant mortality. *JAMA.* 2013;309(1):48–54.
36. Havelin LI, Robertsson O, Fenstad AM, Overgaard S, Garellick G, Furnes O. A Scandinavian experience of register collaboration: the Nordic Arthroplasty Register Association (NARA). *J Bone Joint Surg Am.* 2011;93 (Suppl 3):13–19.
37. Pedersen AB, Mehnert F, Havelin LI, et al. Association between fixation technique and revision risk in total hip arthroplasty patients younger than 55 years of age. Results from the Nordic Arthroplasty Register Association. *Osteoarthritis Cartilage.* 2014;22(5):659–667.
38. Coloma PM, Schuemie MJ, Trifiro G, et al. Combining electronic healthcare databases in Europe to allow for large-scale drug safety monitoring: the EU-ADR Project. *Pharmacoepidemiol Drug Saf.* 2011;20(1):1–11.
39. Patadia VK, Coloma P, Schuemie MJ, et al. Using real-world healthcare data for pharmacovigilance signal detection – the experience of the EU-ADR project. *Expert Rev Clin Pharmacol.* 2015;8(1):95–102.
40. NOMESCO. Nordic Medico-Statistical Committee (NOMESCO) Classification of Surgical Procedures. Available from: <http://nowbase.org/Publikationer/-/media/Projekt%20sites/Nowbase/Publikationer/NCSP/NCSP%2014.ashx>. Accessed May 17, 2015.
41. Nordic Medico-Statistical Committee's (NOMESCO) Classification of External Causes of Injuries (NCECI). Nordic Medico-Statistical Committee. Copenhagen, 1990. Available from: http://www.nordclass.se/ncsp_e.htm. Accessed July 27, 2011.
42. Hsing AW, Ioannidis JP. Nationwide population science: lessons from the Taiwan National Health Insurance Research Database. *JAMA Intern Med.* 2015;175(9):1527–1529.
43. Lo BPL, Ip H, Yang GZ. Transforming Health Care: body sensor networks, wearables, and the Internet of things. *Pulse EMBS.* 2016;7(1):4–8.
44. Eichstaedt JC, Schwartz HA, Kern ML, et al. Psychological language on Twitter predicts county-level heart disease mortality. *Psychol Sci.* 2015;26(2):159–169.
45. Sarker A, Ginn R, Nikfarjam A, et al. Utilizing social media data for pharmacovigilance: a review. *J Biomed Inform.* 2015;54:202–212.
46. Vodopivec-Jamsek V, de Jongh T, Gurol-Urganci I, Atun R, Car J. Mobile phone messaging for preventive health care. *Cochrane Database Syst Rev.* 2012;12:CD007457.
47. Sarfo FS, Treiber F, Jenkins C, et al. Phone-based intervention under nurse guidance after stroke (PINGS): study protocol for a randomized controlled trial. *Trials.* 2016;17(1):436.
48. Ludvigsson JF, Håberg SE, Knudsen GP, et al. Ethical aspects of registry-based research in the Nordic countries. *Clin Epidemiol.* 2015;7:491–508.
49. Nelson EC, Eftimovska E, Lind C, Hager A, Wasson JH, Lindblad S. Patient reported outcome measures in practice. *BMJ.* 2015;350:g7818.
50. Hjollund NH, Larsen LP, Biering K, Johnsen SP, Riiskjær E, Schougaard LM. Use of patient-reported outcome (PRO) measures at group and patient levels: experiences from the generic integrated PRO system, WestChronic. *Interact J Med Res.* 2014;3(1):e5.
51. Sørensen HT, Lash TL, Rothman KJ. Beyond randomized controlled trials: a critical comparison of trials with nonrandomized studies. *Hepatology.* 2006;44(5):1075–1082.
52. McClellan M. Drug safety reform at the FDA – pendulum swing or systematic improvement? *N Engl J Med.* 2007;356(17):1700–1702.
53. Chiolerio A. Big data in epidemiology: too big to fail? *Epidemiology.* 2013;24(6):938–939.
54. Rothman KJ. Significance questing. *Ann Intern Med.* 1986;105(3):445–447.
55. Lang JM, Rothman KJ, Cann CI. That confounded P-value. *Epidemiology.* 1998;9(1):7–8.
56. Toh S, Platt R. Big data in epidemiology: too big to fail? *Epidemiology.* 2013;24(6):939.
57. Ehrenstein V, Petersen I, Smeeth L, et al. Helping everyone do better: a call for validation studies of routinely recorded health data. *Clin Epidemiol.* 2016;8:49–51.
58. Suissa S. Immortal time bias in pharmacoepidemiology. *Am J Epidemiol.* 2008;167(4):492–499.
59. Shrank WH, Patrick AR, Brookhart MA. Healthy user and related biases in observational studies of preventive interventions: a primer for physicians. *J Gen Intern Med.* 2011;26(5):546–550.
60. Rothman KJ, Greenland S. Causation and causal inference in epidemiology. *Am J Public Health.* 2005;95 (Suppl 1):S144–S150.

Clinical Epidemiology

Publish your work in this journal

Clinical Epidemiology is an international, peer-reviewed, open access, online journal focusing on disease and drug epidemiology, identification of risk factors and screening procedures to develop optimal preventive initiatives and programs. Specific topics include: diagnosis, prognosis, treatment, screening, prevention, risk factor modification,

Submit your manuscript here: <https://www.dovepress.com/clinical-epidemiology-journal>

systematic reviews, risk and safety of medical interventions, epidemiology and biostatistical methods, and evaluation of guidelines, translational medicine, health policies and economic evaluations. The manuscript management system is completely online and includes a very quick and fair peer-review system, which is all easy to use.

Dovepress