

particular ecosystems. The reinvasion risk we discussed is not the instant one but the one that may happen months or years after control campaigns. What we argued was about the sustainable success of control across decades, because this is the right time-window for sustained development of affected countries. However, we acknowledge that our data cope with uncontrolled populations and may not predict accurately what will happen after control. Additionally, what needs to be clarified is that what we measured was not dispersal between favorable and unfavorable (e.g., treated) sites, but colonization and sustained establishment of alleles from one favorable site to another, in each particular situation defined by the tsetse species and the zone of study.

To conclude, the aim of our opinion article was to point out a drastic negative relationship between effective population density and average parent–offspring distance in tsetse fly populations across countries and species. Because this may predict unsustainability of controlled campaigns, especially so in high-density areas, we stressed the need to develop appropriate studies to clarify this issue and then find applicable tools to fix it in case of confirmation. We indeed believe that our ability to implement strategies with long-term sustainability will rely on such approaches.

¹Intertryp, IRD, Cirad, Univ Montpellier, Montpellier, France

²Astre, Cirad, Inra, Montpellier, France

³Insect Pest Control Laboratory, Joint Food and Agriculture Organization of the United Nations/International Atomic Energy Agency Program of Nuclear Techniques in Food and Agriculture, A-1400 Vienna, Austria

*Correspondence:
thierry.demeus@ird.fr
<https://doi.org/10.1016/j.pt.2019.07.009>

© 2019 Elsevier Ltd. All rights reserved.

References

1. Lord, J.S. (2019) Comments on T. De Meeûs et al.'s article. *Trends Parasitol.* 35, 741
2. De Meeûs, T. et al. (2019) Negative density dependent dispersal in tsetse flies: a risk for control campaigns? *Trends Parasitol.* 35, 615–621
3. Berté, D. et al. (2019) Population genetics of *Glossina palpalis palpalis* in sleeping sickness foci of Côte d'Ivoire before and after vector control. *Infect. Genet. Evol.* 75, 103963
4. Coe, M.J. et al. (1976) Biomass and production of large african herbivores in relation to rainfall and primary production. *Oecologia* 22, 341–354
5. Cecilia, H. et al. (2019) Environmental heterogeneity drives tsetse fly population dynamics. *bioRxiv*. Published online 13 December 2018. <https://doi.org/10.1101/493650>
6. Dicko, A.H. et al. (2014) Using species distribution models to optimize vector control in the framework of the tsetse eradication campaign in Senegal. *Proc. Natl. Acad. Sci. U. S. A.* 111, 10149–10154
7. Hargrove, J.W. (2001) The effect of temperature and saturation deficit on mortality in populations of male *Glossina m. morsitans* (Diptera: Glossinidae) in Zimbabwe and Tanzania. *Bull. Entomol. Res.* 91, 79–86
8. Vale, G.A. et al. (1984) The use of small plots to study populations of tsetse (Diptera, Glossinidae): difficulties associated with population dispersal. *Insect Sci. Appl.* 5, 403–410

Forum

Causal Inference in Spatial Mapping

Moritz U.G. Kraemer,^{1,2,3,*}
Robert C. Reiner Jr.,⁴
and Samir Bhatt^{5,6,*}

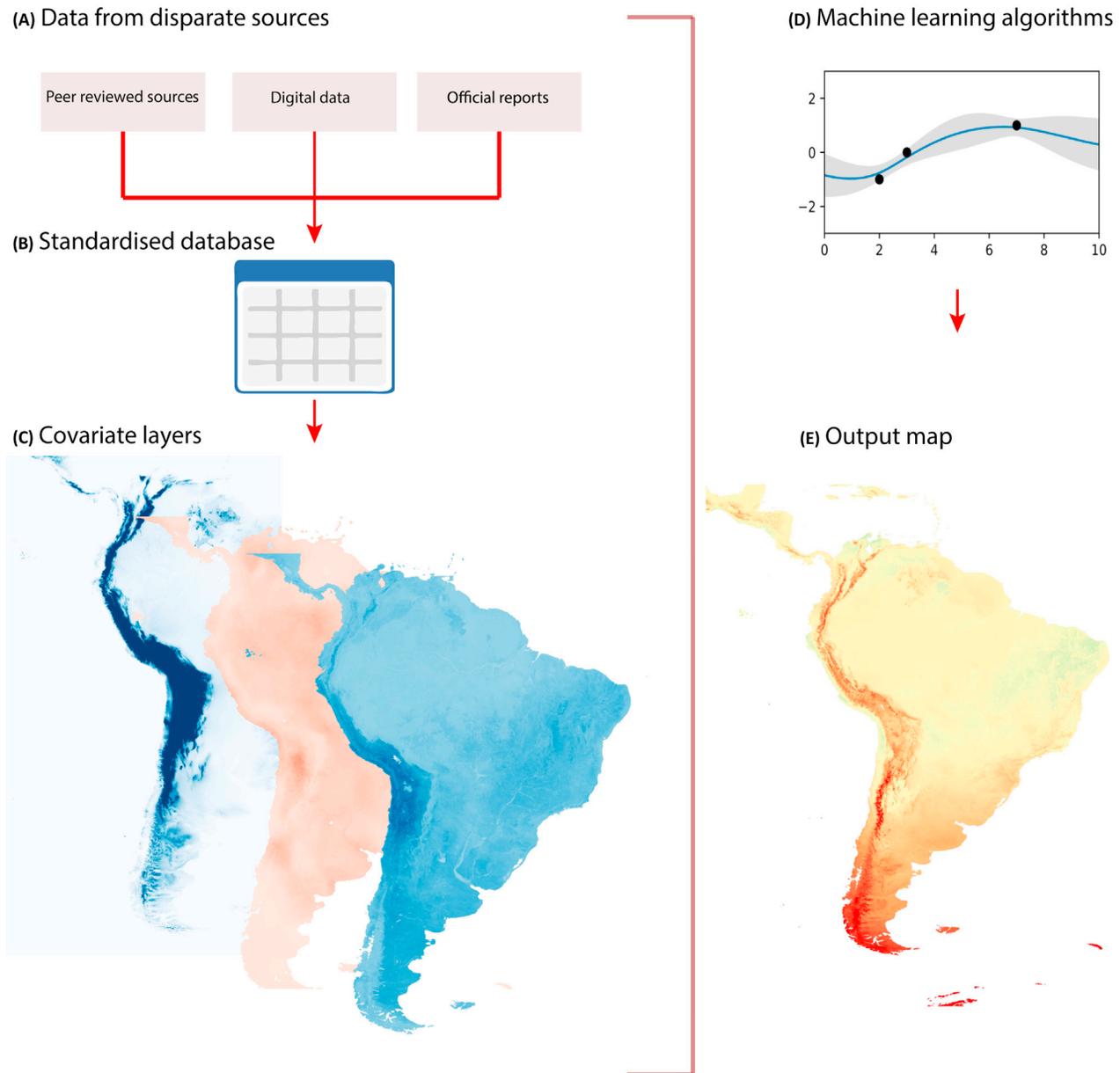
Disease mapping has evolved to a powerful field in epidemiology and public health to focus interventions. Increased precision has come at the expense of interpretability. We propose that future efforts should focus on causal inference to evaluate and predict the effectiveness of intervention strategies to guide decisions more effectively.

With increasing availability of data, and efficient computational methods, tremendous progress has been made to map human health and related health indicators around the world [1–3]. Maps provide visually interpretable and comparable estimates across geographic areas and time. They have been used to expose large variations in disease risk or incidence rates, thereby providing tools for policy makers to target their efforts, resources, or interventions to particular areas. The temporal and spatial resolution of these precision maps is often chosen opportunistically, depending on computational feasibility and the resolution of the underlying observational data. Research and applications in mapping geographic distributions of disease risk have primarily focused on making the most accurate predictions; this has come at the expense of interpretability, making it difficult to base decisions on how to best curb disease burden directly on these high dimensional maps. In most instances spatial mapping approaches use a suite of covariates that are available to match the observed data but apply machine-learning approaches where the contribution of each covariate cannot be interpreted in isolation. In tandem, novel research in causal inference has shown that there is great potential to extract causal structures (variable A has a causal effect on the outcome of interest) from observational data [4]. In this forum article we argue that causal inference should become a priority area of research in disease mapping, an area of research which we show has all the ingredients for successful implementations of causal inference.

Spatial Mapping and Machine Learning

The central question in disease mapping concerns predicting the





Trends in Parasitology

Figure 1. Spatial Mapping and Machine Learning with No Causal Structure.

Typically, data are collated using multiple data sources that include published reports, digital data from social media or smartphones, and official reports (e.g., World Health Organization) (A). Data are then standardized and geographically matched with explanatory variables (B,C). These data are then fed into a learning algorithm of variable complexity (D) that produces the output map at high spatial resolution (E). Maps are generated using the statistical software R using simulated data.

geographic distribution of disease (including mortality) or related indicators such as access to healthcare, socio-demographic factors, and education

among many others. The strength of maps is that they are produced at a unified spatial resolution (e.g., 1×1 km grid cells) which makes them

comparable across space and time and easy to aggregate to any higher spatial resolution (districts, countries). The process of generating maps usually

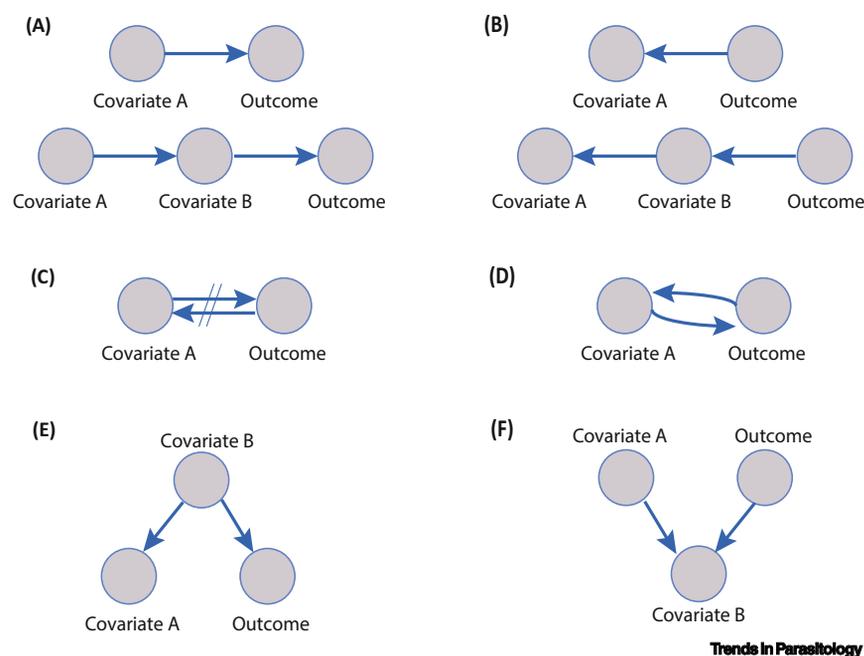


Figure 2. Causal Pathways That Can Exist between an Explanatory Variable and Outcome (e.g., Disease, Development Indicator).

(A) Covariate A causes outcome (altitude causes temperature) directly or indirectly. (B) Outcome causes covariate A (malaria causes anemia, but not vice versa) directly or indirectly. (C) Covariate and outcome are not causally linked. (D) Feedback relationship between two variables: HIV causes poverty, and poverty causes HIV. (E) Outcome and covariate A are linked through a hidden confounder (here, covariate B). An example is: poverty causes malaria and HIV, but HIV and malaria are independent. (F) The observed relationship between Covariate A and outcome is linked by a hidden variable (here, covariate B). An example is: malaria and influenza both cause childhood mortality but are independent from each other. Figure adapted from [4].

starts with collecting data of the outcome of interest (e.g., disease incidence or prevalence) that are abstracted from a variety of sources, including user-generated big data (Figure 1A). These observed data-points are then geo-coded, time-stamped, and paired with relevant covariates (Figure 1B). Covariates considered relevant may be environmental variables such as temperature or precipitation but could also include sociodemographic variables such as urbanization (Figure 1C). The central modeling aim is to make accurate predictions to new locations where no observed data are available. Over the years, models have matured from simple linear models to complex nonlinear

machine-learning approaches, making intuitive interpretations of the model structure difficult (Figure 1D). In the process of generating these maps, however, massive amounts of data have become available across multiple spatial and temporal scales that have already been used to assess the causal relationship between, for example, air quality and infant mortality [5].

Increasingly, there is an important need to ensure that model outputs (disease maps and many more not discussed in this forum) can directly be translated and interpreted in such a way that they can aid public health interventions. For example, a very intuitive question to ask is: what should we do to reduce

the burden of disease in a high-risk area? More specifically, policy makers may ask counterfactual type questions such as: if we reduced poverty by 10%, how many HIV infections could be averted? Or, if we reduced HIV incidence by 10%, how many people could be lifted out of poverty? Answering such questions is very rewarding. In order to compute these so-called counterfactuals, it is not sufficient to rely on correlated associations: we need to understand the causal relationships between the outcome variable and factors that influence them. Understanding the causal direction of complex processes, such as malaria transmission across space and time, is extremely challenging, however, and requires careful attention. Prominent examples of disease maps today do not allow us to make claims about the underlying causal structure.

The Causal Question in Spatial Mapping

Understanding the causal relations from observational data is challenging [6]. It is a widely observed fact that correlation does not (even) imply causation, but having information about the causal relationship between two variables would enable predictions of the effects of actions that change the observed system (e.g., effectiveness of interventions). The main aim of applying causal models to observational data is to recover the underlying causal mechanisms between the outcome variable and the explanatory variables. The gold standard for assessing causal relationships is based on randomized controlled experiments where the effectiveness of a treatment can be directly inferred from the data. These are, however, expensive and impractical to implement in many settings, and cannot be used retrospectively. Matching algorithms such as those developed by Imai *et al.* [7] can be

used to balance datasets to resemble those collected through randomized trials, but careful knowledge of the problem is required to use these approaches effectively. Figure 2 describes common relationships between explanatory covariates and observed data, with intuitive examples.

A few approaches have been proposed using causal inference on large observational data, for example, Pearls Do-calculus [8]; here, Bayesian reasoning is used to infer probabilistic causality. In short, theory is used to measure how the outcome variable may change as a result of an intervention as we would in a controlled experiment (for example, a randomized controlled trial). The trick is that the intervention is not carried out but rather inferred from observational data. Given the large amount of observational data that are collected to carry out spatial mapping (i.e., multiple locations, and over long periods of time), such a method would be a promising extension of current applications. This method, however, may be insufficient to understand complex geo-statistical problems, as pointed out by Mooij *et al.* [9] (problem of conditional independence). For pairwise causal relationships machine-learning additive noise models are proposed by Hoyer *et al.* [10] and Mooij *et al.* [11] to infer causal relations by analyzing differences in the residual errors. This method is very flexible, uses straightforward and efficient regression algorithms, and has already been shown to work for some real-world datasets, including models with more than one variable and which have nonlinear rela-

tionships. The vast amount of literature on the underlying biological mechanisms may help to stratify which covariates should be used to infer causal pathways and which should be left to explain the residual unexplained variation when doing large-scale spatial mapping. Readers can find a summary of other causal models applied to observational data in [11].

Implications for Public Health Interventions

Given the recent advances in theory and application of causal inference on large observational data, we propose that it should become an integral part of future research in mapping disease and development indicators. If done carefully and correctly [12], this field of research can deliver high-resolution maps and, in addition, provide recommendations for targeted intervention that can be evaluated in real-time as new data become available. Similarly, causal discovery from observational data is an active area of research where significant advances can be expected over the next decade.

¹Department of Zoology, University of Oxford, Oxford, UK

²Harvard Medical School, Boston, MA, USA

³Computational Epidemiology Group, Boston Children's Hospital, Boston, MA, USA

⁴Institute for Health Metrics and Evaluation, University of Washington, Seattle, WA, USA

⁵Department of Infectious Disease Epidemiology, Imperial College, London, UK

⁶Big Data Institute, Nuffield Department of Medicine, University of Oxford, Oxford, UK

*Correspondence:
moritz.kraemer@zoo.ox.ac.uk,
s.bhatt@imperial.ac.uk

<https://doi.org/10.1016/j.pt.2019.06.005>

© 2019 The Author(s). Published by Elsevier Ltd. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

References

1. Tusting, L.S. *et al.* (2019) Mapping changes in housing in sub-Saharan Africa from 2000 to 2015. *Nature* 568, 391–394
2. Kraemer, M.U.G. *et al.* (2019) Past and future spread of the arbovirus vectors *Aedes aegypti* and *Aedes albopictus*. *Nat. Microbiol.* 4, 854–863
3. Osgood-Zimmerman, A. *et al.* (2018) Mapping child growth failure in Africa between 2000 and 2015. *Nature* 555, 41–47
4. Peters, J. *et al.* (2019) Elements of Causal Inference (MIT Press)
5. Heft-Neal, S. *et al.* (2018) Robust relationship between air quality and infant mortality in Africa. *Nature* 559, 254–258
6. Hernán, M.A. (2018) The C-word: scientific euphemisms do not improve causal inference from observational data. *Am. J. Public Health* 108, 616–619
7. Imai, K. *et al.* (2010) Identification, inference and sensitivity analysis for causal mediation effects. *Stat. Sci.* 25, 51–71
8. Pearl, J. (1995) Casual diagrams for empirical research. *Biometrika* 82, 669–710
9. Mooij, J. *et al.* (2009) Regression by dependence minimization and its application to causal inference in additive noise models. In *Proceedings of the 26th Annual International Conference on Machine Learning – ICML '09*, pp. 1–8, (ACM)
10. Hoyer, P.O. *et al.* (2009) Nonlinear causal discovery with additive noise models. *Adv. Neural Inf. Process. Syst.* 21, 689–696
11. Mooij, J.M. *et al.* (2016) Distinguishing cause from effect using observational data: methods and benchmarks. *J. Mach. Learn. Res.* 17, 1–102
12. Hernán, M.A. (2019) Comment: spherical cows in a vacuum: data analysis competitions for causal inference. *Stat. Sci.* 34, 69–71