



# Causal inference and counterfactual prediction in machine learning for actionable healthcare

Mattia Prosperi<sup>1</sup>✉, Yi Guo<sup>2,3</sup>, Matt Sperrin<sup>4</sup>, James S. Koopman<sup>5</sup>, Jae S. Min<sup>1</sup>, Xing He<sup>2</sup>, Shannan Rich<sup>1</sup>, Mo Wang<sup>6</sup>, Iain E. Buchan<sup>7</sup> and Jiang Bian<sup>2,3</sup>

**Big data, high-performance computing, and (deep) machine learning are increasingly becoming key to precision medicine—from identifying disease risks and taking preventive measures, to making diagnoses and personalizing treatment for individuals. Precision medicine, however, is not only about predicting risks and outcomes, but also about weighing interventions. Interventional clinical predictive models require the correct specification of cause and effect, and the calculation of so-called counterfactuals, that is, alternative scenarios. In biomedical research, observational studies are commonly affected by confounding and selection bias. Without robust assumptions, often requiring a priori domain knowledge, causal inference is not feasible. Data-driven prediction models are often mistakenly used to draw causal effects, but neither their parameters nor their predictions necessarily have a causal interpretation. Therefore, the premise that data-driven prediction models lead to trustable decisions/interventions for precision medicine is questionable. When pursuing intervention modelling, the bio-health informatics community needs to employ causal approaches and learn causal structures. Here we discuss how target trials (algorithmic emulation of randomized studies), transportability (the licence to transfer causal effects from one population to another) and prediction invariance (where a true causal model is contained in the set of all prediction models whose accuracy does not vary across different settings) are linchpins to developing and testing intervention models.**

Advances in computing and machine learning have opened unprecedented paths to processing and inferring knowledge from big data. Deep learning has been game-changing for many analytics challenges, beating humans and other machine learning approaches in decision or action tasks, such as playing games, and aiding or augmenting tasks such as driving or recognizing manipulated images. The potential applications of deep learning in healthcare have been widely speculated<sup>1</sup>, especially in precision medicine—the timely and tailored prevention, identification, diagnosis and treatment of disease. However, the use of data-driven machine learning approaches to model causality—for instance, to uncover new causes of disease or assess treatment effects—carries dangers of unintended consequences. Therefore, the Hippocratic principle of ‘first do no harm’ is being adopted<sup>2</sup> alongside rigorous study design, validation and implementation, with attention to ethics and bias avoidance.

Precision medicine models are not only descriptive or predictive—for example, assessing the mortality risk for a patient undergoing a surgical procedure—but also decision-supporting or interventional—for example, choosing and personalizing the procedure with the highest probability of favourable outcome. Predicting risks and outcomes differs from weighing interventions and intervening. Prediction calculates a future event in the absence of any action or change; intervention presumes an enacted choice that may influence the future, which requires consideration of the underlying causal structure.

Intervention imagines how the world would be if we made different choices—for example, ‘would the patient be cured if we

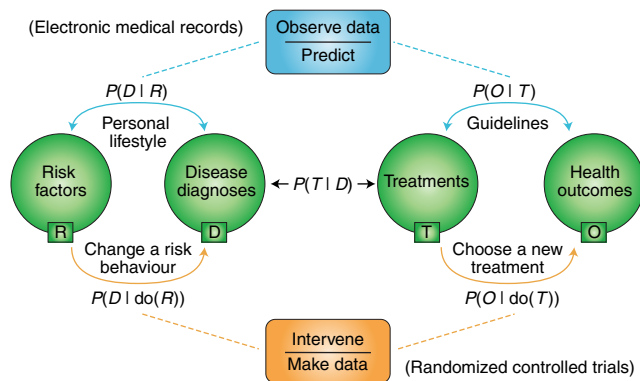
administered amoxicillin instead of a cephalosporin for their upper respiratory tract infection?’ or ‘if a pre-hypertensive patient had accomplished 10–15 min of moderate physical activity per day instead of being prescribed a diuretic, would they have become hypertensive five years later?’ By asking ourselves what would have been the effect of something if we had not taken an action, or vice versa, we are computing so-called counterfactuals. Among different counterfactual options, we choose the ones that minimize harm while maximizing patient benefit. A similar cognitive process happens when a deep learning machine plays a game and must decide on the next move. Such an artificial neural network architecture has been fed the game ruleset and millions of game scenarios, and has learned through trial and error by playing against other human players, networks or itself. In each move of the game, the machine chooses the best counterfactual move based on its domain knowledge<sup>3,4</sup>.

With domain knowledge of the variables involved in a hypothesized cause–effect route and enough data generated at random to cover all possible path configurations, it is possible to deduce causal effects and calculate counterfactuals. Randomization and domain knowledge are key: when either is not met, causal inference may be flawed<sup>5</sup>.

In clinical research, randomized controlled trials (RCTs) permit direct testing of causal hypotheses since randomization is guaranteed a priori by design even with limited domain knowledge. On the other hand, observational data collected retrospectively usually does not fulfil such requirements, thus limiting what secondary data analyses can discover. For instance, databases collating

<sup>1</sup>Department of Epidemiology, College of Public Health and Health Professions, College of Medicine, University of Florida, Gainesville, FL, USA.

<sup>2</sup>Department of Health Outcomes and Biomedical Informatics, College of Medicine, University of Florida, Gainesville, FL, USA. <sup>3</sup>Cancer Informatics and eHealth Core, University of Florida Health Cancer Center, Gainesville, FL, USA. <sup>4</sup>Division of Informatics, Imaging and Data Sciences, University of Manchester, Manchester, UK. <sup>5</sup>Department of Epidemiology, School of Public Health, University of Michigan, Ann Arbor, MI, USA. <sup>6</sup>Department of Management, Warrington College of Business, University of Florida, Gainesville, FL, USA. <sup>7</sup>Institute of Population Health, University of Liverpool, Liverpool, UK. ✉e-mail: [m.prosperi@ufl.edu](mailto:m.prosperi@ufl.edu)



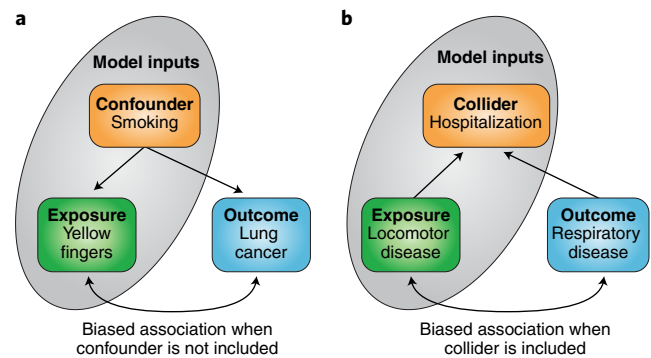
**Fig. 1 | Conditional versus interventional probabilities.** When we observe data (for example, electronic medical records) we can learn a model that predicts the probability of a disease  $D$  given certain risk factors  $R$ , that is,  $P(D | R)$ , or a model that predicts the chance of a health outcome  $O$  for a given treatment  $T$ , that is,  $P(O | T)$ . However, these models cannot be used to support decisions, because they assume that variables of the model remain unchanged, people keep their lifestyles, and the standard of care is followed. When a risk factor is modified or a new treatment is tested, for example, in an RCT, then we ‘make’ new data, and compute different probabilities, which are  $P(D | \text{do}(R))$  and  $P(O | \text{do}(T))$ . Conditional and interventional probabilities are not necessarily the same; for example, treatments are randomized in trials, while they are not in clinical practice.

electronic medical records do not explicitly record domain or contextual knowledge (for example, why one drug was prescribed over another), and are littered with many types of bias, including protopathic bias (when a therapy is given based on symptoms, yet the disease is undiagnosed), indication bias (when a risk factor appears to be associated with a health outcome, but the outcome may be caused by the reason for which the risk factor initially appeared), or selection bias (when a study population does not represent the target population—for example, insured people or hospitalized patients).

Therefore, the development of health intervention models from observational data (no matter how big) is problematic, regardless of the method used (no matter how deep) because of the nature of the data. Fitting a machine learning model to observational data and using it for counterfactual prediction may lead to harmful consequences. One well-known example is that of prediction tools for crime recidivism that convey racial discriminatory bias<sup>6</sup>. Any instrument inferred from existing population data may be biased by gender, sexual orientation, race or ethnicity discrimination, and carry forward such bias when employed to aid decisions<sup>7</sup>.

The health and biomedical informatics community, charged with maximizing the utility of healthcare information, needs to be attuned to the limitations of data-driven inference of intervention models, and needs safeguards for counterfactual prediction modeling. These topics are addressed here as follows: first, we give a brief outline of causality, counterfactuals, and the problem of inferring cause–effect relations from observational data; second, we provide examples where automated learning has failed to infer a trustworthy counterfactual model for precision medicine; third, we offer insights on methodologies for automated causal inference; finally, we describe potential approaches to validate automated causal inference methods, including transportability and prediction invariance.

We aim not to criticize the use of machine learning for the development of clinical prediction models<sup>8,9</sup>, but rather clarify that prediction and intervention models have different developmental paths and intended uses<sup>10</sup>. We recognize the promising applications of machine learning in healthcare for descriptive/predictive tasks rather than interventional tasks—for example, screening



**Fig. 2 | Examples of confounding bias and collider bias.** **a**, Confounding (**a**) can occur when there exists a common cause for both exposure and outcome, while a collider (**b**) is a common effect of both exposure and outcome. Not including a confounder or including a collider in a model results in biased associations.

images for diabetic retinopathy<sup>11</sup>—even when diagnostic models have been shown to be susceptible to errors when applied to different populations<sup>12</sup>. Yet it is important to distinguish prediction works from others that seek to optimize treatment decisions and are clearly interventional<sup>13</sup>, where validating counterfactuals becomes necessary.

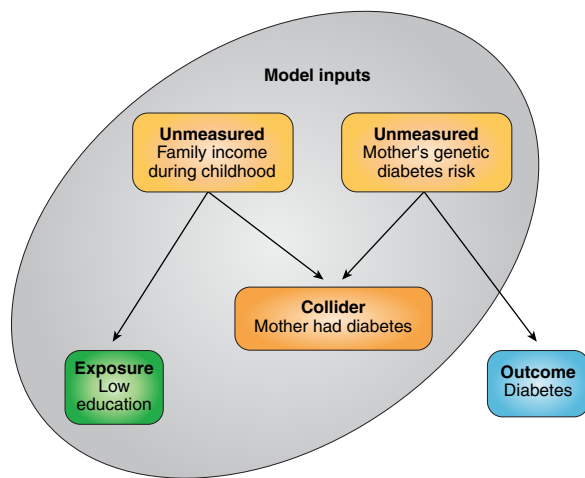
### Causal inference and counterfactuals

Causality has been described in various ways. For the purpose of this work, it is useful to recall the deterministic (yet can be made probabilistic) definitions by means of counterfactual conditionals—for example, the original 1748 proposition by Hume or the 1973 formalization by Lewis<sup>14</sup>—paraphrased as: an event  $E$  causally depends on  $C$  if, and only if,  $E$  always follows after  $C$  and  $E$  does not occur when  $C$  has not occurred (unless something else caused  $E$ ). The counterfactual-based definition contains an implicit time component and works in a chained manner, where effects can become causes of other subsequent effects. Causes can be regarded as necessary, sufficient, contributory or non-redundant<sup>15</sup>.

Causal inference addresses the problem of ascertaining causes and effects from data. Causes can be determined through prospective experiments to observe  $E$  after  $C$  is tried or withheld, by keeping constant all other possible factors that can influence either the choice of  $C$  or the happening of  $E$ , or—what is easier and more often done—randomizing the choice of  $C$ . Formally, we acknowledge that the conditional probability  $P(E | C)$  of observing  $E$  after observing  $C$  can be different from the interventional probability  $P(E | \text{do}(C))$  of observing  $E$  after doing  $C$ . In RCTs—when  $C$  is randomized—the ‘do’ is guaranteed and unconditioned, while with observational data, it is not. Causal calculus helps resolve interventional from conditional probabilities when a causal structure is assumed<sup>16</sup>. Figure 1 illustrates the difference between observing and doing using biomedical target examples.

In a nutshell, the major hurdles to ascertaining causal effects from observational data include: the failure to disambiguate interventional from conditional distributions, to identify all potential sources of bias<sup>17</sup> and to select an appropriate functional form for all variables, that is, model misspecification<sup>18–20</sup>.

Two well-known types of bias are confounding and collider bias (Fig. 2). Given an outcome—that is, the objective of a (counterfactual) prediction—confounding occurs when there exists a variable that causes the outcome and the effect, leading to the conclusion that an exposure is associated with the outcome even though it does not cause it. For instance, cigarette smoking causes both nicotine-stained, yellow fingers and lung cancer. Yellow fingers, as the exposure or independent variable, can be spuriously associated



**Fig. 3 | An example of M-bias.** When estimating the effect of education level on diabetes risk, mother's history of diabetes could be mistaken as a confounder and included in a model, but it is a collider by the effect of history of family income and genetic risk.

with lung cancer if smoking, the underlying confounding variable, is unaccounted. Yellow fingers alone could be used to predict lung cancer but cleaning the skin would not reduce the risk of lung cancer. Therefore, an intervention model that used yellow fingers as the actionable item would be futile, while a model that acted upon smoking (the cause) would be effective in reducing lung cancer risk.

A collider is a variable that is caused by both the exposure (or causes of the exposure) and the outcome (or causes of the outcome). Conditioning on a collider biases the estimate of the causal effect of exposure on outcome. A classic example involves the association between locomotor and respiratory diseases. Originally observed in hospitalized patients and thought biologically plausible, this association could not be established in the general population<sup>21</sup>. In this case, hospitalization status functions as a collider because it introduces selection bias, as people with locomotor disease or respiratory disease have a higher risk of being admitted to a hospital.

Another example of collider bias is the obesity paradox<sup>22</sup>. This paradox refers to the counterintuitive evidence of lower mortality among people who are obese within certain clinical subpopulations—for example, in patients with heart failure. A more careful consideration of the covariate–outcome relationship in this case reveals heart failure is a collider. Had an intervention been developed by means of such model, treating obesity would not have been suggested as an actionable feature to reduce the risk of mortality.

Causal inference can become more complex when a variable may be mistaken for a confounder but actually functions as a collider. This phenomenon is called M-bias since the associated causal diagram is usually drawn in an M-shaped form<sup>23,24</sup>. A classic M-bias example is the effect of education on diabetes, controlled through family history of diabetes and income. In a hypothetical study, it could be reasonable to regard the mother's (or father's) history of diabetes as a confounder, because it is associated with both education level and diabetes status, and it is not caused by either. However, family history of diabetes' associations with the education and diabetes are not causal but are in turn confounded by family income and family genetic risk for diabetes, respectively, that might not be measured as input (Fig. 3). At this point, the mother's diabetes becomes a collider, and including it would induce a biased association between education and diabetes through the links from family income and genetic risk. Specifically, the estimate of the causal effect of education on diabetes would be biased. In general, including the mother's diabetes in the input covariate would lead to bias both if there was a zero or non-zero causal effect;<sup>25</sup> however, if the

unmeasured covariates were included, the bias would be resolved (by a so-called backdoor path blocking)<sup>26</sup>. The M-bias example shows how the causal structure choice (which could be machine learned) can influence the causal effect inference; we will discuss the two in detail later in a specific section.

For brevity, we do not describe moderators, mediators or other important concepts in causality modelling. Nonetheless, it is useful to mention instrumental variables, which determine variation in an explanatory variable (for example, a treatment) but have no independent effect on the outcome of interest. Instrumental variables, therefore, can be used to resolve unmeasured confounding in absence of prospective randomization.

### An old neural network fiasco and a new possible paradox

In 1997, Cooper et al<sup>27</sup> investigated several machine learning models, including rule-based and neural networks, for predicting mortality of hospital patients admitted with pneumonia. The neural network greatly outperformed logistic regression; however, the authors discouraged the use of black box models in clinical practice. They showed how the rule-based method learned that 'IF patient admitted (with pneumonia) has history of asthma THEN patient has lower risk of death from pneumonia'<sup>28</sup>. This counterintuitive association was also later confirmed using generalized additive models<sup>29</sup>. The physicians explained that patients admitted with pneumonia and a known history of asthma are likely to be transferred to intensive care and treated aggressively, thus having higher odds of survival than the general population admitted with pneumonia. The authors recommended the employment of interpretable models instead of black boxes, to identify counterintuitive, surprising patterns and remove them. At this point, the model development is no longer automated and requires domain knowledge. On reflection, those models inferred without modifications, either interpretable or black box, would have worked well at predicting mortality but they could not have been used to test new interventions to reduce mortality, as the recommended actions would have consequentially led to 'less care' for asthmatic patients.

More recently, a possible data-driven improvement in the evaluation of fall risk in hospitals was investigated<sup>30</sup>. Standard practice involves a nurse-led evaluation of patients' history of falls, comorbidities, mental health, gait, ambulatory aids and intravenous therapy summarized with the Morse scale. To assess the predictive ability of the Morse scale (standard practice), its individual components and new expert-based predictors (for example, extended clinical profiles and information on hospital staffing levels), a matched study was performed including patients with and without a fall. Logistic regression and decision trees were used. The additional variables hypothesized by the experts were associated with the outcome and all new models yielded higher discrimination than the Morse scale, but a surprising finding was observed: in all configurations, older patients were at a lower risk of falls. This is contrary to current experts' knowledge, which associates older age with increased frailty and therefore fall risk. If such a model were used for intervention, it would not prioritize the elderly for fall prevention—a potentially devastating consequence of data-driven inference. It is uncertain if this old age paradox is due to a bias. One possible explanation is that older patients are usually monitored and aided more frequently because they are indeed at higher risk, while younger people may be more independent and less prone to accept assistance. Other issues could be survivorship bias, selectively unreported falls and study design. One possible approach to bias reduction is to design the study and extract the data by simulating an RCT, where causal effects on 'randomized' interventions can be estimated directly, as we discuss in the next section.

### The target trial

Target trials refer to RCTs that can be emulated using data from large, observational databases to answer causal questions of comparative

**Table 1 | Emulation of a randomized clinical trial using observational data and algorithmic randomization (the target trial), with the objective to reduce bias and allow more reliable treatment effects estimates**

	RCT	Target trial
Data source	Prospective	Observational
Sample size	Small	Large
Variables	Few	Many
Eligibility and time zero (baseline)	Straightforward	Problematic (for example, multiple baseline points, follow-up requirements)
Treatment assignment	Randomized by design	Randomized algorithmically (for example, via propensity score matching)
Outcome evaluation	Flexible	Flexible (with some caveats for blind outcome studies)
Analysis plan	Relatively straightforward (for example, intention to treat) and flexible (for example, Bayesian adaptive), but can further require bias correction (for example, g-formula)	More complex (need to model treatment assignment) yet can use same techniques as for RCT (for example, g-formula)
Risk of bias	Relatively low	Possible (for example, residual confounding, wrong choice of time zero)
Flexibility to assess extra-protocol causal effects	Limited	High

treatment effect<sup>31</sup>. Although RCTs are the gold standard for discerning causal effects, there exists many scenarios in which they are neither feasible nor ethical to conduct. Alternatively, observational data appropriately adjusted for measured confounding—for instance, via propensity score matching—can be used to emulate randomized treatment assignment; this may be feasible with electronic medical records where many individual-level attributes can be linked to resolve bias. The target trial protocol requires prospective enrolment-like eligibility criteria, a description of treatment strategies and assignment procedures, the identification of time course from a baseline to the outcome, a causal query (for example, treatment effect), and an analysis plan (for example, a regression model), as shown in Table 1.

As an example of the target trial framework, data from public surveillance and clinical claims repositories were used to replicate two RCTs, one investigating treatment effects on colorectal cancer and the other on pancreatic adenocarcinoma<sup>32</sup>. Each study explicitly adhered to the target trial framework, deviating from the RCT design only in the assignment procedures, justifiably due to lack of randomization. The results were consistent with the target trials—all of which reported a null effect. In contrast, when the authors modelled the treatment effects using a non-RCT-like study design with the same variables, the mortality estimates were both inconsistent with the target trials. These examples demonstrate the need to uphold target trial design in the investigation of treatment effects using observational data. Moreover, coupled with machine learning methods equipped to extrapolate more useful information from big data sources, the target trial framework has the potential to serve as the foundation for exploring causal processes currently unknown.

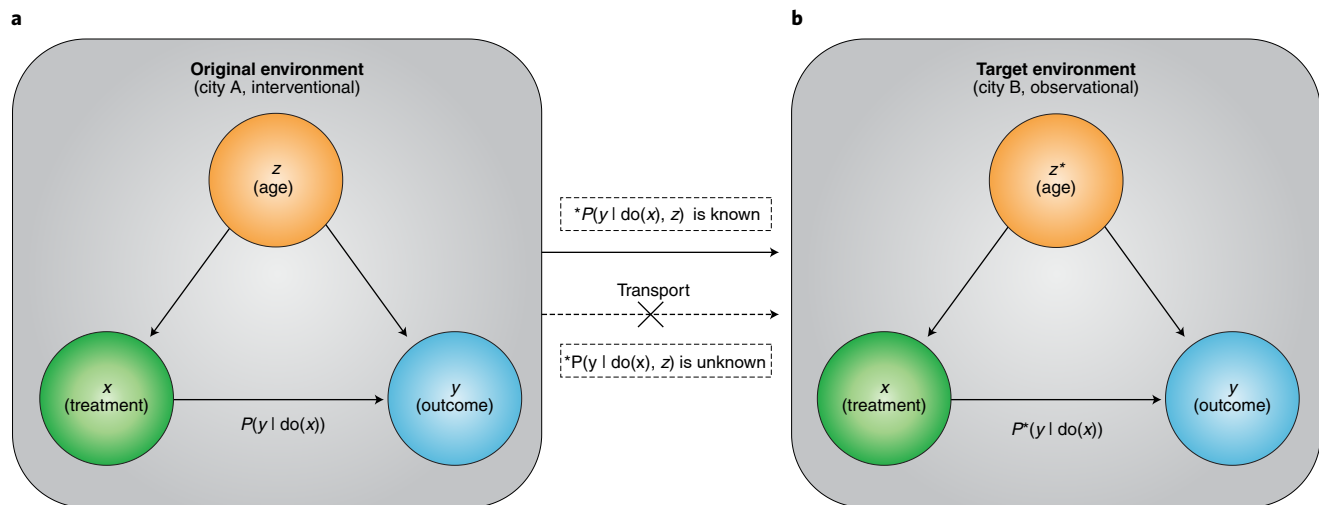
### Causal effect inference and automated learning of causal structures

Prediction models inferred automatically using data without any domain knowledge—from linear regression to deep learning—only approximate a function, and do not necessarily hold a causal meaning. Instead, by estimating interventional in place of conditional probabilities, models can reproduce causal mechanisms. Through counterfactual predictions, models become interventional and actionable, and avoid the pitfalls such as those described in the pneumonia and fall risk examples. In the previous sections we showed that it is possible to directly estimate causal effects by generating data through RCTs or by simulating RCTs with observational

data. Here, we delve further into the approaches to unveiling and disambiguating causality from observational data, including the assumptions to be made. We can categorize two main tasks: (1) estimating causal effects, and (2) learning causal structures. In the first one, a causal structure, or a set of cause–effect relationships and pathways, is defined a priori, input variables are fixed, and causal effects are estimated for a specific input–output pair—for example, the causal effect of diabetes on glaucoma. Directed acyclic graphs (DAGs)<sup>33</sup>—also known as Bayesian networks—and structural equation models<sup>34,35</sup> are often used to model such structures. While with RCT data the estimation of causal effects can be done directly, the estimation of causal effects from observational data requires thoughtful handling of potential biases and confounding. Methods like inverse probability weighting attempt to weigh single observations to mimic the effects of randomization with respect to one variable of interest (for example, an exposure or a treatment)<sup>36</sup>. Other techniques include targeted maximum likelihood estimation<sup>37–39</sup>, g-methods<sup>40,41</sup> and propensity score matching<sup>42</sup>. Often, these estimators can be coupled with machine learning—for example, causal decision trees<sup>43</sup>, Bayesian regression trees<sup>44</sup>, and random forests for estimating individual treatment effects<sup>45</sup>. As previously mentioned, model misspecification—that is, defining the wrong causal structure and the choice of variables to handle confounding and bias—can lead to wrong estimation of causal effects. With big data, especially datasets with numerous features, choosing adjustments, and even over-adjusting using all variables, is problematic. Feature selection algorithms based on conditional independence scoring have been proposed<sup>46</sup>.

Automated causal structure learning uses conditional independence tests and structure search algorithms over given DAGs subject to certain assumptions—for example, ‘causal sufficiency’ that has no unmeasured common causes and no selection bias. In 1990, an important result on independence and conditional independence constraints—the d-separation equivalence theorem<sup>47</sup>—led to the development of automated search and ambiguity resolution of causal structures from data, through so-called patterns and partial ancestral graphs. When assumptions are met (Markov or causal faithfulness<sup>48</sup>), there are asymptotically correct procedures that can predict an effect or raise an ambiguity, and determine graph equivalence<sup>49</sup>. However, the probability of an effect cannot be obtained without deciding on a prior distribution of the graphs and parameters. Also, the number of graphs is super-exponential in the





**Fig. 4 | A selection diagram for illustrating transportability. a,** A causal effect of treatment  $x$  on outcome  $y$ ,  $P(y | do(x))$ , is found through an RCT, and quantified in the original environment of city A. **b,** The  $x$ -to- $y$  causal effect is transportable from city A to city B as  $P^*(y | do(x))$  if both the overall causal effect  $P(y | do(x))$  and the age-specific causal effect  $P(y | do(x), z)$  are known, while it is not transportable if the latter is unknown.

number of observed variables and may even be infinite with hidden variables<sup>50</sup>, making an exhaustive search computationally unfeasible<sup>51</sup>. Today, several heuristic methods for causal structure search are available, from the Peter-Clark (PC) algorithm that assumes causal sufficiency, to others like the fast causal inference (FCI) or really fast causal inference (RFCI) algorithms that extend to latent variables<sup>52,53</sup>.

The enrichment of deep learning with causal methods also provides interesting new insights to address bias. For instance, a theoretical analysis for identification of individual treatment effect under strong ignorability has been derived<sup>54</sup>, and an approach to exploit instrumental variables for counterfactual prediction within deep learning is also available<sup>55</sup>.

### Model transportability and prediction invariance

Validation of causal effects under determined causal structures is especially needed when such effects are estimated in limited settings, for example, RCTs. Transportability is a data fusion framework for external validation of intervention models and counterfactual queries. As defined by Pearl and Bareinboim<sup>56</sup>, transportability is a “license to transfer causal effects learned in experimental studies to a new population, in which only observational studies can be conducted.” By combining datasets generated under heterogeneous conditions, transportability provides formal mathematical tools to (1) evaluate whether results from one study (for example, a causal relationship identified in an RCT) could be used to generate valid estimates of the same causal effect in another study of different setting (for example, an observational study of the same causal effect in a different population); and (2) estimate what the causal effect would have been if the study had been conducted in the new setting<sup>57,58</sup>. The framework utilizes selection diagrams<sup>59</sup>, encoding the causal relationships of variables of interest in a study population, and about the characteristics in which the target and study populations differ. If the structural constraints among variables in the selection diagrams are resolvable through the do-calculus, a valid estimate of the causal effect in the target population can be calculated using the extant causal effect from the original study, which means that the observed causal effect is transportable.

One of Pearl’s transportability examples is shown in Fig. 4. In this example, an RCT is conducted in city A (original environment) and a causal effect of treatment  $x$  on outcome  $y$ ,  $P(y | do(x))$ , is determined.

We wish to generalize if the treatment works also in the population of city B (target environment) where only observational data is available, since it happens that the age distribution in city A,  $P(z)$ , is different than that  $P^*(z)$  in city B. The city B specific  $x$ -to- $y$  causal effect  $P^*(y | do(x))$  is estimated as:

$$P^*(y | do(x)) = \sum_z P(y | do(x), z) P^*(z)$$

In this transport formula, the age-specific causal effect estimated in the RCT,  $P(y | do(x), z)$ , is combined with the observed age distribution in the target population,  $P^*(z)$ , to obtain the causal effect  $P^*(y | do(x))$  in city B.

On the other hand, a causal effect is not always transportable. Following the example above, the  $x$ -to- $y$  causal effect is not transportable from city A to city B if only the overall causal effect  $P(y | do(x))$  is known whereas the age-specific causal effect  $P(y | do(x), z)$  is unknown.

Transportability theory is being extended to a variety of more complex causal relationships—for example, sample selection bias<sup>58</sup>, leaping forward from toy examples to real-world problems<sup>60</sup>. Therefore—linking back with the problematic examples we discussed in the previous sections—one could use transportability to determine how the asthma or old age effects are or are not transportable from one population to another. It is also worth noting how transportability evokes the field of domain adaptation, which aims to learn a model in one source population that can be used in a different target distribution. In fact, domain adaptation has been employed to address sample selection bias<sup>61</sup>.

An interesting next of kin to transportability is prediction invariance<sup>62</sup>. Among all models that show invariance in their predictive accuracy across different experimental settings and interventions, there is a high probability that the causal model will be a member of that set. For example, Schulam and Saria<sup>72</sup> introduced the counterfactual Gaussian process to predict continuous-time trajectories under irregular sampling, handling biases arising from clinical protocols. In another work, aimed at addressing issues of supervised learning when training and target distributions differ (that is, dataset shift), Saria et al<sup>63</sup> proposed the ‘surgery estimator’, defined as an interventional distribution<sup>16</sup> that is invariant to differences across environments. The surgery estimator works by learning a relationship in the training data that is generalizable to

the target population, by incorporating prior knowledge about the data-generating processes that are expected to differ between the original and target populations. It was applied in real-world cases where causal structures were unknown.

## Conclusions

We explored common pitfalls of data-driven developments in machine learning for healthcare, distinguishing between prediction and intervention models that are actionable in support of clinical decision-making. Importantly, the development of intervention models requires careful consideration of causality. Hernan et al.<sup>64</sup> commented that “a recent influx of data analysts, many not formally trained in statistical theory, bring a fresh attitude that does not a priori exclude causal questions” but called—and we strongly endorse such call—for training curricula in data science that properly differentiate descriptive, prediction and intervention modelling.

Undertaking causal machine learning is key to ethical artificial intelligence for healthcare, equivalent to a doctor's oath to ‘first do no harm’<sup>65</sup>. Healthcare intervention models involve actionable inputs and need—implicitly or explicitly—to model causal pathways to compute the correct counterfactuals. There are ongoing discussions in the machine learning community about model explainability for bias avoidance and fairness in decisions<sup>66</sup>. Bias is a core topic in causal theory. Explainability may be a ‘weaker’ model property than causality. Explaining the role of input variables in changing the output of a black box neither assures a correct interpretation of the input–output mechanism nor unveils the cause–effect relationships. For instance, in a deep learning system that predicts the risk of heart attack, a subsequent analysis could be able to explain that the input variables ‘race’ and ‘blood pressure’ affect the risk, but could not say if these findings are causal, since they may be biased by stratification or unmeasured confounders, or mediated by other factors in the causal pathway. Fairness in machine learning aims to develop models that avoid social discrimination due to historically biased data and involves the same conceptual hurdles as learning from observational data. In fact, the usage of causal models has been advocated to identify and mitigate discriminatory relationships in data<sup>67</sup>. Recently, a study in cancer prognostics presented a causal structure coupled to deep learning to eliminate collider bias and provide unbiased individual predictions<sup>68</sup>, although it did not explicitly test for transportability.

For context-specific intervention models, where a causal structure is available or a target trial design can be devised, we then recommend evaluation of model transportability for a given set of action queries—for example, treatment options or risk modifiers. For broader exploratory analyses where causal structures need to be identified or clarified, prediction invariance could be used. Transportability and prediction invariance could become core tools to reporting protocols for intervention models, in line with the current standards for prognostic and diagnostic models<sup>69</sup>. A transportable model can be integrated into clinical guidelines to augment healthcare with action-savvy predictions, in pursuit of better precision medicine.

Received: 25 November 2019; Accepted: 4 June 2020;

Published online: 13 July 2020

## References

- Norgeot, B., Glicksberg, B. S. & Butte, A. J. A call for deep-learning healthcare. *Nat. Med.* **25**, 14–15 (2019).
- Wiens, J. et al. Do no harm: a roadmap for responsible machine learning for health care. *Nat. Med.* **25**, 1337–1340 (2019).
- Silver, D. et al. Mastering the game of Go without human knowledge. *Nature* **550**, 354–359 (2017).
- Jin, P., Keutzer, K. & Levine, S. Regret minimization for partially observable deep reinforcement learning. In *35th Int. Conf. Machine Learning* **80**, 2342–2351 (ICML, 2018).
- Pearl, J. & Mackenzie, D. *The Book of Why: The New Science of Cause and Effect* (Basic Books, 2018).
- Chouldechova, A. Fair prediction with disparate impact: a study of bias in recidivism prediction instruments. *Big Data* **5**, 153–163 (2017).
- Kusner, M., Loftus, J., Russell, C. & Silva, R. Counterfactual fairness. In *Advances in Neural Information Processing Systems* Vol. 31, 4069–4079 (MIT Press, 2017).
- Christodoulou, E. et al. A systematic review shows no performance benefit of machine learning over logistic regression for clinical prediction models. *J. Clin. Epidemiol.* **110**, 12–22 (2019).
- Bian, J., Buchan, I., Guo, Y. & Prosser, M. Statistical thinking, machine learning. *J. Clin. Epidemiol.* **116**, 136–137 (2019).
- Baker, R. E., Peña, J. M., Jayamohan, J. & Jérusalem, A. Mechanistic models versus machine learning, a fight worth fighting for the biological community? *Biol. Lett.* **14**, 20170660 (2018).
- Gulshan, V. et al. Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs. *JAMA* **316**, 2402–2410 (2016).
- Winkler, J. K. et al. Association between surgical skin markings in dermoscopic images and diagnostic performance of a deep learning convolutional neural network for melanoma recognition. *JAMA Dermatol.* **155**, 1135–1141 (2019).
- Komorowski, M., Celi, L. A., Badawi, O., Gordon, A. C. & Faisal, A. A. The Artificial Intelligence Clinician learns optimal treatment strategies for sepsis in intensive care. *Nat. Med.* **24**, 1716–1720 (2018).
- Lewis, D. K. Causation. *J. Philos.* **70**, 556–567 (1973).
- Mackie, J. L. *The Cement of the Universe* (Clarendon, 1974).
- Pearl, J. *Causality: Models, Reasoning and Inference* (Cambridge Univ. Press, 2009).
- Rothman, K. J., Greenland, S. & Lash, T. *Modern Epidemiology* 3rd edn (Lippincott Williams & Wilkins, 2012).
- Lehmann, E. L. Model specification: the views of Fisher and Neyman, and later developments. *Stat. Sci.* **5**, 160–168 (1990).
- Vansteelandt, S., Bekaert, M. & Claeskens, G. On model selection and model misspecification in causal inference. *Stat. Meth. Med. Res.* **21**, 7–30 (2012).
- Asteriou, D., Hall, S. G., Asteriou, D. & Hall, S. G. in *Applied Econometrics* 2nd edn 176–197 (Palgrave Macmillan, 2016).
- Sackett, D. L. Bias in analytic research. *J. Chronic Dis.* **32**, 51–63 (1979).
- Banack, H. R. & Kaufman, J. S. The ‘obesity paradox’ explained. *Epidemiology* **24**, 461–462 (2013).
- Pearl, J. Causal diagrams for empirical research. *Biometrika* **82**, 669–688 (1995).
- Greenland, S., Pearl, J. & Robins, J. M. Causal diagrams for epidemiologic research. *Epidemiology* **10**, 37–48 (1999).
- Westreich, D. & Greenland, S. The table 2 fallacy: Presenting and interpreting confounder and modifier coefficients. *Am. J. Epidemiol.* **177**, 292–298 (2013).
- Wei, L., Brookhart, M. A., Schneeweiss, S., Mi, X. & Setoguchi, S. Implications of m bias in epidemiologic studies: A simulation study. *Am. J. Epidemiol.* **176**, 938–948 (2012).
- Cooper, G. F. et al. An evaluation of machine-learning methods for predicting pneumonia mortality. *Artif. Intell. Med.* **9**, 107–138 (1997).
- Ambrosino, R., Buchanan, B. G., Cooper, G. F. & Fine, M. J. The use of misclassification costs to learn rule-based decision support models for cost-effective hospital admission strategies. In *Proc. Annual Symp. Computer Applications Medical Care* 304–308 (AMIA, 1995).
- Caruana, R. et al. Intelligent models for healthcare: predicting pneumonia risk and hospital 30-day readmission. in *Proc. ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining* 1721–1730 (ACM, 2015).
- Lucero, R. J. et al. A data-driven and practice-based approach to identify risk factors associated with hospital-acquired falls: applying manual and semi- and fully-automated methods. *Int. J. Med. Inform.* **122**, 63–69 (2019).
- Hernán, M. A. & Robins, J. M. Using big data to emulate a target trial when a randomized trial is not available. *Am. J. Epidemiol.* **183**, 758–764 (2016).
- Petito, L. C. et al. Estimates of overall survival in patients with cancer receiving different treatment regimens: emulating hypothetical target trials in the surveillance, epidemiology, and end results (SEER)–Medicare linked database. *JAMA Netw. Open* **3**, e200452–e200452 (2020).
- Pearl, J. Causal diagrams for empirical research. *Biometrika* **82**, 669–688 (1995).
- Westland, J. C. *Structural Equation Models* 1–15 (Springer, 2019).
- Bollen, K. A. & Pearl, J. in *Handbook of Causal Analysis for Social Research* (ed. Morgan, S. L.) 301–328 (Springer, 2013).
- Hernán, M. A. & Robins, J. M. Estimating causal effects from epidemiological data. *J. Epidemiol. Commun. Health* **60**, 553 (2006).
- van der Laan, M. J. & Rubin, D. Targeted maximum likelihood learning. *Int. J. Biostat.* **6**, 2 (2006).
- Schuler, M. S. & Rose, S. Targeted maximum likelihood estimation for causal inference in observational studies. *Am. J. Epidemiol.* **185**, 65–73 (2017).

39. van der Laan, M. J. & Rose, S. *Targeted Learning: Causal Inference For Observational And Experimental Data* (Springer, 2011).
40. Naimi, A. I., Cole, S. R. & Kennedy, E. H. An introduction to g methods. *Int. J. Epidemiol.* **46**, 756–762 (2017).
41. Robins, J. M. & Hernán, M. A. in *Longitudinal Data Analysis* (eds Fitzmaurice, G. et al.) 553–599 (CRC, 2008).
42. Rosenbaum, P. R. & Rubin, D. B. The central role of the propensity score in observational studies for causal effects. *Biometrika* **70**, 41–55 (1983).
43. Li, J., Ma, S., Le, T., Liu, L. & Liu, J. Causal decision trees. *IEEE Trans. Knowl. Data Eng.* **29**, 257–271 (2017).
44. Hahn, P. R., Murray, J. & Carvalho, C. M. Bayesian regression tree models for causal inference: regularization, confounding, and heterogeneous effects. *Bayesian Anal.* <https://doi.org/10.1214/19-BA1195> (2020).
45. Lu, M., Sadiq, S., Feaster, D. J. & Ishwaran, H. Estimating individual treatment effect in observational data using random forest methods. *J. Comput. Graph. Stat.* **27**, 209–219 (2018).
46. Schneeweiss, S. et al. High-dimensional propensity score adjustment in studies of treatment effects using health care claims data. *Epidemiology* **20**, 512–522 (2009).
47. Verma, T. & Pearl, J. in *Machine Intelligence and Pattern Recognition* Vol. 9 (eds Shachter, R. D. et al.) 69–76 (Elsevier, 1990).
48. Jaber, A., Zhang, J. & Bareinboim, E. Causal identification under Markov equivalence. In *34th Conf. Uncertainty in Artificial Intelligence* (UAI, 2018).
49. Richardson, T. in *Compstat* (eds Dutter, R. & Grossmann, W.) 482–487 (Springer, 1994).
50. Heckerman, D., Meek, C. & Cooper, G. In *Innovations in Machine Learning* (eds Holmes, D. E. & Jain, L. C.) 1–28 (Springer, 2006).
51. Peter Spirtes, C. G. and R S. *Causation, Prediction, and Search* 2nd edn (MIT Press, 2003).
52. Glymour, C., Zhang, K. & Spirtes, P. Review of causal discovery methods based on graphical models. *Front. Genet.* **10**, 524 (2019).
53. Colombo, D. & Maathuis, M. H. Order-independent constraint-based causal structure learning. *J. Mach. Learn. Res.* **15**, 3921–3962 (2014).
54. Shalit, U., Johansson, F. D. & Sontag, D. Estimating individual treatment effect: generalization bounds and algorithms. In *Proc. 34th Int. Conf. Machine Learning* Vol. 70 (eds Precup, D. & Teh, Y. W.) 3076–3085 (PMLR, 2017).
55. Hartford, J., Lewis, G., Leyton-Brown, K. & Taddy, M. Deep {IV}: a flexible approach for counterfactual prediction. In *Proc. 34th Int. Conf. Machine Learning* Vol. 70 (eds Precup, D. & Teh, Y. W.) 1414–1423 (PMLR, 2017).
56. Pearl, J. & Bareinboim, E. External validity: from do-calculus to transportability across populations. *Stat. Sci.* **29**, 579–595 (2014).
57. Dahabreh, I. J., Robertson, S. E., Tchetgen, E. J., Stuart, E. A. & Hernán, M. A. Generalizing causal inferences from individuals in randomized trials to all trial-eligible individuals. *Biometrics* **75**, 685–694 (2019).
58. Bareinboim, E. & Pearl, J. Causal inference and the data-fusion problem. *Proc. Natl Acad. Sci. USA* **113**, 7345–7352 (2016).
59. Pearl, J. & Bareinboim, E. Transportability of causal and statistical relations: a formal approach. In *Proc. IEEE Int. Conf. Data Mining* (IEEE, 2011).
60. Lee, S., Cornea, J. D. & Bareinboim, E. General identifiability with arbitrary surrogate experiments. In *Proc. 35th Conf. Uncertainty in Artificial Intelligence* (UAI, 2019).
61. Huang, J., Smola, A. J., Gretton, A., Borgwardt, K. M. & Schölkopf, B. Correcting sample selection bias by unlabeled data. In *Advances in Neural Information Processing Systems* Vol. 19 (eds Schölkopf, B. et al.) 601–609 (MIT Press, 2007).
62. Peters, J., Bühlmann, P. & Meinshausen, N. Causal inference by using invariant prediction: identification and confidence intervals. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **78**, 947–1012 (2016).
63. Subbaswamy, A., Schulam, P. & Saria, S. Preventing failures due to dataset shift: learning predictive models that transport. In *Proc. 22nd Int. Conf. Artificial Intelligence and Statistics* 3118–3127 (AiStats, 2019).
64. Hernán, M. A., Hsu, J. & Healy, B. A second chance to get causal inference right: a classification of data science tasks. *CHANCE* **32**, 42–49 (2019).
65. Wiens, J. et al. Do no harm: a roadmap for responsible machine learning for health care. *Nat. Med.* **25**, 1337–1340 (2019).
66. Rudin, C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nat. Mach. Intell.* **1**, 206–215 (2019).
67. Kusner, M. J. & Loftus, J. R. The long road to fairer algorithms. *Nature* **578**, 34–36 (2020).
68. van Amsterdam, W. A. C., Verhoeff, J. J. C., de Jong, P. A., Leiner, T. & Eijkemans, M. J. C. Eliminating biasing signals in lung cancer images for prognosis predictions with deep learning. *npj Digit. Med.* **2**, 122 (2019).
69. Moons, K. G. M. et al. Transparent reporting of a multivariable prediction model for individual prognosis or diagnosis (TRIPOD): explanation and elaboration. *Ann. Intern. Med.* **162**, W1–W73 (2015).

## Acknowledgements

J.B.'s, Y.G.'s and M.P.'s research for this work was in part supported by the University of Florida (UF)'s Creating the Healthiest Generation—Moonshot initiative, supported by the UF Office of the Provost, UF Office of Research, UF Health, UF College of Medicine and UF Clinical and Translational Science Institute. M.W.'s research for this work was supported in part by the Lanzillotti–McKethan Eminent Scholar Endowment.

## Author contributions

M.P., Y.G., J.B. and M.W. conceived the premise, wrote the paper, designed the figures and tables, and revised the paper. M.S., X.E. and S.R. contributed to specific sections, aided with the figures and tables, and with revision. J.K., I.B. and J.M. contributed to specific sections and helped with revisions.

## Competing interests

The authors declare no competing interests.

## Additional information

Correspondence should be addressed to M.P.

Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

© Springer Nature Limited 2020