# A Survey on Causal Inference

LIUYI YAO, University at Buffalo, USA

ZHIXUAN CHU, University of Georgia, USA

SHENG LI, University of Georgia, USA

YALIANG LI, Alibaba Group, USA

JING GAO, University at Buffalo, USA

AIDONG ZHANG, University of Virginia, USA

Causal inference is a critical research topic across many domains, such as statistics, computer science, education, public policy and economics, for decades. Nowadays, estimating causal effect from observational data has become an appealing research direction owing to the large amount of available data and low budget requirement, compared with randomized controlled trials. Embraced with the rapidly developed machine learning area, various causal effect estimation methods for observational data have sprung up. In this survey, we provide a comprehensive review of causal inference methods under the potential outcome framework, one of the well known causal inference framework. The methods are divided into two categories depending on whether they require all three assumptions of the potential outcome framework or not. For each category, both the traditional statistical methods and the recent machine learning enhanced methods are discussed and compared. The plausible applications of these methods are also presented, including the applications in advertising, recommendation, medicine and so on. Moreover, the commonly used benchmark datasets as well as the open-source codes are also summarized, which facilitate researchers and practitioners to explore, evaluate and apply the causal inference methods.

## 1 INTRODUCTION

In everyday language, correlation and causality are commonly used interchangeably, although they have quite different interpretations. Correlation indicates a general relationship: two variables are correlated when they display an increasing or decreasing trend [6]. Causality is also referred to as cause and effect where the cause is partly responsible for the effect, and the effect is partly dependent on the cause. Causal inference is the process of drawing a conclusion about a causal connection based on the conditions of the occurrence of an effect. The main difference between causal inference and inference of correlation is that the former analyzes the response of the effect variable when the cause is changed [91, 130].

It is well known that "*correlation does not imply causation.*" For example, a study showed that girls have breakfast normally have lightweight than the girls who don't, and thus concluded that having breakfast can help to lose weight. But in fact, these two events may just have correlation instead of causality. Maybe the girls who have breakfast everyday have a better lifestyle, such as exercise frequently, sleep regularly, and have a healthy diet, which finally makes them have lightweight. In this case, having a better lifestyle is the common cause of both having breakfast and lightweight, and thus we also can treat it as a confounder of the causality between having breakfast and lightweight.

In many cases, it seems obvious that one action can cause another; however, there exists also many cases that we cannot easily tease out and make sure the relationship. Therefore, learning causality is one dauntingly challenging problem. The most effective way of inferring causality is to conduct a randomized controlled trial, which randomly assigns participants into a treatment group or a control group. As the randomized study is conducted, the only expected difference between the control and treatment groups is the outcome variable being studied. However, in reality, randomized controlled trials are always time-consuming and expensive, and

thus the study cannot involve many subjects, which may be not representative of the real-world population a treatment/intervention would eventually target. Another issue is that the randomized controlled trials only focus on the average of samples, and it doesn't explain the mechanism or pertain for individual subjects. In addition, ethical issues also need to be considered in most of the randomized controlled trials, which largely limits its applications. Therefore, instead of the randomized controlled trials, the observational data is a tempting shortcut. Observational data is obtained by the researcher simply observing the subjects without any interfering. That means, the researchers have no control over treatments and subjects, and they just observe the subjects and record data based on their observations. From the observational data, we can find their actions, outcomes, and information about what has occurred, but cannot figure out the mechanism why they took a specific action. For the observational data, the core question is how to get the counterfactual outcome. For example, we want to answer this question "would this patient have different results if he received a different medication?" Answering such counterfactual questions is challenging due to two reasons [119]: the first one is that we only observe the factual outcome and never the counterfactual outcomes that would potentially have happened if they have chosen a different treatment option. The second one is that treatments are typically not assigned at random in observational data, which may lead the treated population differs significantly from the general population.

To solve these problems in causal inference from observational data, researchers develop various frameworks, including the potential outcome framework [111, 129] and the structural causal model [89, 92, 94]. The potential outcome framework is also known as the Neyman-Rubin Potential Outcomes or the Rubin Causal Model. In the example we mentioned above, a girl would have a particular weight if she had breakfast normally everyday, whereas she would have a different weight if she didn't have breakfast normally. To measure the causal effect of having breakfast normally for a girl, we need to compare the outcomes for the same person under both situations. Obviously, it is impossible to see both potential outcomes at the same time, and one of the potential outcomes is always missing. The potential outcome framework aims to estimate such potential outcomes and then calculate the treatment effect. Therefore, the treatment effect estimation is one of the central problems in causal inference under the potential outcome framework. Another influential framework in causal inference is the structural causal model (SCM), which includes the causal graph and the structural equations. The structural causal model describes the causal mechanisms of a system where a set of variables and the causal relationship among them are modeled by a set of simultaneous structural equations.

Causal inference has a close relationship with the machine learning area. In recent years, the magnificent bloom of the machine learning area enhances the development of the causal inference area. Powerful machine learning methods such as decision tree, ensemble methods, deep neural network, are applied to estimate the potential outcome more accurately. In addition to the amelioration on the outcome estimation model, machine learning methods also provide a new aspect to handle the confounders. Benefit from the recently deep representation learning methods, such as generative adversarial neural network, the confounder variables are adjusted by learning the balanced representation for all covariates, so that conditioning on the learned representation, the treatment assignment is independent of the confounder variables. In machine learning, the more data the better. However, in causal inference, the more data alone is not yet enough. Having more data only helps to get more precise estimates, but it cannot make sure these estimates are correct and unbiased. Machine learning methods enhance the development of causal inference, meanwhile, causal inference also helps machine learning methods. The simple pursuit of predictive accuracy is insufficient for modern machine learning research, and correctness and interpretability are also the targets of machine learning methods. Causal inference is starting to help to improve machine learning, such as recommender systems or reinforcement learning.

In this article, we provide a comprehensive review of the causal inference methods under the potential outcome framework. We first introduce the basic concepts of the potential outcome framework as well as its three critical assumptions to identify the causal effect. After that, various causal inference methods with these three assumptions are discussed in detail, including re-weighting methods, stratification methods, matching

based methods, tree-based methods, representation-based methods, multi-task learning based methods, and meta-learning methods. Additionally, causal effect estimation methods that relax the three assumptions are also described to fulfill the needs in different settings. After introducing various causal effect estimation methods, the real-world applications that the discussed methods have great potential to benefit are discussed, including the advertisement area, recommendation area, medicine area, and reinforcement learning area as the representative examples.

To the best of our knowledge, this is the first paper that provides a comprehensive survey for causal inference methods under the potential outcome framework. There also exist several surveys that discuss one category of the causal effect estimation methods, such as the survey of matching based methods [131], survey of tree-based and ensemble-based method [10], and the review of dynamic treatment regimes [25]. For the structural causal model, it is suggested to refer the survey [91] or the book [90]. We will also briefly discuss the relation and difference between the two causality frameworks at the end of our survey. There is also a survey about learning causality from observational data [44] whose content ranges from inferring the causal graph from observational data, structural causal model, potential outcome framework and their connection to machine learning. Compared with the surveys mentioned above, this survey paper mainly focuses on the theoretical background of the potential outcome framework, the representative methods across the statistic domain and machine learning domain, and how this framework and the machine learning area enhance each other.

To summarize, our contributions of this survey are as follows:

- *New taxonomy:* We separate various causal inference methods into two major categories based on whether they require the three assumptions of the potential outcome framework. The category requiring three assumptions are further divided into seven sub-categories based on the way to handle the confounder variables.
- *Comprehensive review:* We provide a comprehensive survey of the causal inference methods under the potential outcome framework. In each category, the detailed descriptions of the representative methods, the connection and comparison between the mentioned methods, and the general summation are provided.
- *Abundant resources:* In this survey, we list the state-of-art methods, the benchmark data sets, open-source codes, and representative applications.

The rest of the paper is organized as follows. In section 2, the background of the potential outcome framework is introduced, including the basic definitions, the assumptions, and the fundamental problems with their general solutions. In Section 3, the methods under three assumptions are presented. Then, in Section 4, we discuss the problem when some assumptions are not satisfied, and describe the methods that relax those assumptions. Next, we provide experimental guidelines in Section 5. Afterward, in Section 6, the typical applications of causal inference are illustrated. Finally, Section 7 summarizes the paper.

## 2  BASIC OF CAUSAL INFERENCE

In this section, we present the background knowledge of causal inference, including task description, mathematical notions, assumptions, challenges and general solutions. We also give an illustrative example that will be used throughout this survey.

Generally speaking, the task of causal inference is to estimate the outcome changes if another treatment had been applied. For example, suppose there are two treatments that can be applied to patients: Medicine A and Medicine B. When applying Medicine A to the interested patient cohort, the recovery rate is 70%, while applying Medicine B to the same cohort, the recovery rate is 90%. The change of recovery rate is the effect that treatment (i.e., medicine in this example) asserts on the recovery rate.

The above example describes an ideal situation to measure the treatment effect: applying different treatments to the same cohort. In real-world scenarios, this ideal situation can only be approximated by a randomized

experiment, in which the treatment assignment is controlled, such as a completely random assignment. In this way, the group receives a specific treatment can be viewed as an approximation to the cohort we are interested in.

However, performing randomized experiments are expensive, time-consuming, and sometimes even unethical. Therefore, estimating the treatment effect from observational data has attracted growing attention due to the wide availability of observational data. Observational data usually contains a group of individuals taken different treatments, their corresponding outcomes, and possibly more information, but *without direct access to the reason/mechanism why they took the specific treatment*. Such observational data enable researchers to investigate the fundamental problem of learning the causal effect of a certain treatment without performing randomized experiments. To better introduce various treatment effect estimation methods, the following section introduces several definitions including unit, treatment, outcome, treatment effect, and other information (pre- and post-treatment variables) provided by observational data.

## 2.1 Definitions

Here we define the notations under the potential outcome framework [111, 129], which is logically equivalent to another framework, the structural causal model framework [62]. The foundation of potential outcome framework is that the causality is tied to treatment (or action, manipulation, intervention), applied to a unit [59]. The treatment effect is obtained by comparing units' potential outcomes of treatments. In the following, we first introduce three essential concepts in causal inference: unit, treatment, and outcome.

DEFINITION 1. *Unit. A unit is the atomic research object in the treatment effect study.*

A unit can be a physical object, a firm, a patient, an individual person, or a collection of objects or persons, such as a classroom or a market, at a particular time point [59]. Under the potential outcome framework, the atomic research objects at different time points are different units. One unit in the dataset is a sample of the whole population, so in this survey, the term "sample" and "unit" are used interchangeably.

DEFINITION 2. *Treatment. Treatment refers to the action that applies (exposes, or subjects) to a unit.*

Let $W$ ($W \in \{0, 1, 2, \ldots, N_W\}$) denote the treatment, where $N_W + 1$ is the total number of possible treatments. In the aforementioned medicine example, Medicine A is a treatment. Most of the literatures consider the binary treatment, and in this case, the group of units applied with treatment $W = 1$ is the *treated group*, and the group of units with $W = 0$ is the *control group*.

DEFINITION 3. *Potential outcome. For each unit-treatment pair, the outcome of that treatment when applied on that unit is the potential outcome [59].*

The potential outcome of treatment with value $w$ is denoted as $Y(W = w)$.

DEFINITION 4. *Observed outcome. The observed outcome is the outcome of the treatment that is actually applied.*

The observed outcome is also called factual outcome, and we use $Y^F$ to denote it where F stands for "factual". The relation between the potential outcome and the observed outcome is: $Y^F = Y(W = w)$ where $w$ is the treatment actually applied.

DEFINITION 5. *Counterfactual outcome: Counterfactual outcome is the outcome if the unit had taken another treatment.*

The counterfactual outcomes are the potential outcomes of the treatments except the one actually taken by the unit. Since a unit can only take one treatment, only one potential outcome can be observed, and the remaining unobserved potential outcomes are the counterfactual outcome. In the multiple treatment case, let $Y^{CF}(W = w')$ denote the counterfactual outcome of treatment with value $w'$. In the binary treatment case, for

notation simplicity, we use $Y^{CF}$ to denote the counterfactual outcome, and $Y^{CF} = Y(W = 1 - w)$, where $w$ is the treatment actually taken by the unit.

In the observational data, besides the chosen treatments and the observed outcome, the units' other information is also recorded, and they can be separated as pre-treatment variables and the post-treatment variables.

DEFINITION 6. *Pre-treatment variables: Pre-treatment variables are the variables that will not be affected by the treatment.*

Pre-treatment variables are also named as *background variables*, and they can be patients' demographics, medical history, and etc. Let $X$ denote the pre-treatment variables.

DEFINITION 7. *post-treatment variables: The post-treatment variables are the variables that are affected by the treatment.*

One example of post-treatment variables is the intermediate outcome, such as the lab test after taking the medicine in the aforementioned medicine example.

In the following sections, the terminology *variable* refers to the pre-treatment variable unless otherwise specified.

**Treatment Effect**. After introducing the observational data and the key terminologies, the treatment effect can be quantitatively defined using the above definitions. The treatment effect can be measured at the population, treated group, subgroup, and individual levels. To make these definitions clear, here we define the treatment effect under binary treatment, and it can be extended to multiple treatments by comparing their the potential outcomes.

At the population level, the treatment effect is named as the Average Treatment Effect (ATE), which is defined as:

$$\text{ATE} = \mathbb{E}[\mathbf{Y}(W = 1) - \mathbf{Y}(W = 0)], \tag{1}$$

where $\mathbf{Y}(W = 1)$ and $\mathbf{Y}(W = 0)$ are the potential treated and control outcome of the whole population respectively.

For the treated group, the treatment effect is named as Average Treatment effect on the Treated group (ATT), and it is defined as:

$$\text{ATT} = \mathbb{E}[\mathbf{Y}(W = 1)|W = 1] - \mathbb{E}[\mathbf{Y}(W = 0)|W = 1], \tag{2}$$

where $\mathbf{Y}(W = 1)|W = 1$ and $\mathbf{Y}(W = 0)|W = 1$ are the potential treated and control outcome of the treated group respectively.

At the subgroup level, the treatment effect is called Conditional Average Treatment Effect (CATE), which is defined as:

$$\text{CATE} = \mathbb{E}[\mathbf{Y}(W = 1)|X = x] - \mathbb{E}[\mathbf{Y}(W = 0)|X = x], \tag{3}$$

where $\mathbf{Y}(W = 1)|X = x$ and $\mathbf{Y}(W = 0)|X = x$ are the potential treated and control outcome of the subgroup with $X = x$, respectively. CATE is a common treatment effect measurement under the case where the treatment effect varies across different subgroups, which is also known as the heterogeneous treatment effect.

At the individual level, the treatment effect is called Individual Treatment Effect (ITE), and the ITE of unit $i$ is defined as:

$$\text{ITE}_i = Y_i(W = 1) - Y_i(W = 0), \tag{4}$$

where $Y_i(W = 1)$ and $Y_i(W = 0)$ are the potential treated and control outcome of unit $i$ respectively. In some literatures [60, 122], the ITE is viewed equivalent to the CATE.

**Objective**. For causal inference, our objective is to estimate the treatment effects from the observational data. Formally speaking, given the observational dataset, $\left\{X_i, W_i, Y_i^F\right\}_{i=1}^N$, where $N$ is the total number of units in the datasets, the goal of the causal inference task is to estimate the treatment effects defined above.

## 2.2 An Illustrative Example

To better illustrate causal inference, we use the following example combined with the notations defined above to give an overview. In this example, we want to evaluate the treatment effects of several different medications for one disease, by exploiting the observational data (i.e., the electronic health records) that include demographic information of patients, the specific medication with the specific dosage taken by patients, and the outcome of medical tests. Obviously, we can only get one factual outcome for a specific patient from electronic health records, and thus the core task is to predict what would have happened if a patient took another treatment (i.e., a different medication or the same medication with a different dosage). Answering such counterfactual questions is very challenging. Therefore, we want to use causal inference to predict all of the potential outcomes for each patient over all of the medications with different dosages. Then, we can reasonably and accurately evaluate and compare the treatment effect of different medications for this disease.

One particular point to keep in mind is that for each medication, they may have different dosages. For example, for medication A, the dosage range can be a continuous variable in the range $[a, b]$ while for medication B, the dosage can be a categorical variable that has several specific dosage regimens.

In the aforementioned example, the units are the patients with the studied disease. The treatments refer to the different medications with specific dosages for this disease, and we use $W$ ($W \in \{0, 1, 2, \ldots, N_W\}$) to denote these treatments. For example, $W_i = 1$ can represent the medication $A$ with a specific dosage is taken by the unit $i$, and $W_i = 2$ represents the medication $B$ with a specific dosage is taken by the unit $i$. $Y$ is the outcome, such as one type of blood test that can measure the medication's ability to destroy the disease and lead to the recovery of the patients. Let $Y_i(W = 1)$ denote the potential outcome of medication $A$ with a specific dosage on patient $i$. The features of patients may include age, gender, clinical presentation, and some other medical tests, etc. Among these features, age, gender and other demographic information are pre-treatment variables that cannot be affected by taking a treatment. Some clinical presentations and medical tests are affected by taking medications, and they are post-treatment variables. In this example, our goal to estimate the treatment effects of different medications for this disease based on the provided observational data.

In the following sections, we will continuously use this example to explain more concepts and illustrate intuitions behind various causal inference methods.

## 2.3 Assumptions

In order to estimate the treatment effect, the following assumptions are commonly used in the causal inference literature.

ASSUMPTION 2.1. **Stable Unit Treatment Value Assumption (SUTVA).** *The potential outcomes for any unit do not vary with the treatment assigned to other units, and, for each unit, there are no different forms or versions of each treatment level, which lead to different potential outcomes.*

This assumption emphasizes two points: The first point is the independence of each unit, that is, there are no interactions between units. In the context of the above illustrative example, one patient's outcome will not affect other patients' outcomes.

The second point is the single version for each treatment. In the above example, Medicine A with different dosages are different treatments under the SUTVA assumption.

ASSUMPTION 2.2. **Ignorability.** *Given the background variable, $X$, treatment assignment $W$ is independent to the potential outcomes, i.e., $W \perp\!\!\!\perp Y(W = 0), Y(W = 1)|X$.*

In the context of the illustrative example, this ignorability assumption indicates two folds: First, if two patients have the same background variable $X$, their potential outcomes should be the same whatever the treatment assignment is, i.e., $p(Y_i(0), Y_i(1)|X = x, W = W_i) = p(Y_j(0), Y_j(1)|X = x, W = W_j)$. Analogously, if two patients

have the same background variable value, their treatment assignment mechanism should be same whatever the value of potential outcomes they have, i.e., $p(W|X = x, Y_i(0), Y_i(1)) = p(W|X = x, Y_j(0), Y_j(1))$. The ignorability assumption is also named as unconfoundedness assumption. With this unconfoundedness assumption, for the units with the same background variable $X$, their treatment assignment can be viewed as random.

ASSUMPTION 2.3. **Positivity**. *For any value of $X$, treatment assignment is not deterministic:*

$$P(W = w|X = x) > 0, \quad \forall w \text{ and } x. \tag{5}$$

If for some values of $X$, the treatment assignment is deterministic; then for these values, the outcomes of at least one treatment could never be observed. In this case, it would be unable and meaningless to estimate the treatment effect. To be more specific, suppose there are two treatments: Medicine A and Medicine B. Let's assume that patients with age greater than 60 are always assigned with medicine A, then it will be unable and meaningless study the outcome of medicine B on those patients. In other words, the positivity assumption indicates the variability, which is important for treatment effect estimation.

In [59], the ignorability and the positivity assumptions together are called *Strong Ignorability* or *Strongly Ignorable Treatment Assignment*.

With these assumptions, the relationship between the observed outcome and the potential outcome can be rewritten as:

$$
\begin{aligned}
\mathbb{E}[Y(W = w)|X = x] &= \mathbb{E}[Y(W = w)|W = w, X = x] \text{ (Ignorability)} \\
&= \mathbb{E}[Y^F|W = w, X = x],
\end{aligned}
\tag{6}
$$

where $Y^F$ is the random variable of the observed outcome, and $Y(W = w)$ is the random variable of the potential outcome of treatment $w$. If we are interested in the potential outcome of one specific group (either the subgroup, the treated group, or the whole population), the potential outcome can be obtained by taking expectation of the observed outcome over that group.

With the above equation, we can rewrite the treatment effect defined in Section 2.1 as follows:

$$
\begin{aligned}
\text{ITE}_i &= W_i Y_i^F - W_i Y_i^{CF} + (1 - W_i)Y_i^{CF} - (1 - W_i)Y_i^F \\
\text{ATE} &= \mathbb{E}_X \left[ \mathbb{E}[Y^F|W = 1, X = x] - \mathbb{E}[Y^F|W = 0, X = x] \right] \\
&= \frac{1}{N} \sum_i (Y_i(W = 1) - Y_i(W = 0)) = \frac{1}{N} \sum_i \text{ITE}_i \\
\text{ATT} &= \mathbb{E}_{X_T} \left[ \mathbb{E}[Y^F|W = 1, X = x] - \mathbb{E}[Y^F|W = 0, X = x] \right] \\
&= \frac{1}{N_T} \sum_{\{i:W_i=1\}} (Y_i(W = 1) - Y_i(W = 0)) = \frac{1}{N_T} \sum_{\{i:W_i=1\}} \text{ITE}_i \\
\text{CATE} &= \mathbb{E}[Y^F|W = 1, X = x] - \mathbb{E}[Y^F|W = 0, X = x] \\
&= \frac{1}{N_x} \sum_{\{i:X_i=x\}} (Y_i(W = 1) - Y_i(W = 0)) = \frac{1}{N_x} \sum_{\{i:X_i=x\}} \text{ITE}_i
\end{aligned}
\tag{7}
$$

where $Y_i(W = 1)$ and $Y_i(W = 0)$ are the potential treated/control outcomes of unit $i$, $N$ is the total number of units in the whole population, $N_T$ is the number of units in the treated group, and $N_x$ is the number of units in the group with $X = x$. The second line in the ATE, ATT and CATE equations are their empirical estimations. Empirically, the ATE can be estimated as the average of ITE on the entire population. Similarly, ATT and CATE can be estimated as the average of ITE on the treated group and specific subgroup separately.

However, due to the fact that the potential treated/control outcomes can never be observed simultaneously, the key point in the treatment effect estimation is how to estimate the counterfactual outcome in ITE estimation

or how to estimate the $\frac{1}{N_*} \sum_i Y_i(W = 1)$ and $\frac{1}{N_*} \sum_i Y_i(W = 0)$, where $N_*$ denotes $N$, $N_T$ or $N_x$. In the following section, we will discuss the challenges in estimation these terms and briefly introduce the general solutions.

## 2.4 Confounders and General Solutions

As mentioned above, how to estimate the average potential treated/control outcome over a specific group is the core of causal inference. Let's take ATE as a case study: When estimating the ATE, a natural idea is to directly use the average of observed treated/control outcomes, i.e., $\hat{\text{ATE}} = \frac{1}{N_T} \sum_{i=1}^{N_T} Y_i^F - \frac{1}{N_C} \sum_{i=1}^{N_C} Y_j^F$, where $N_T$ and $N_C$ is the number of units in the treated and control group, respectively. However, due to the existence of *confounders*, there is a serious problem in this estimation: this calculated ATE includes a spurious effect brought by the confounders.

DEFINITION 8. *Confounders. Confounders are the variables that affect both the treatment assignment and the outcome.*

Confounders are some special pre-treatment variables, such as age in the medicine example. When directly using the average of observed treated/control outcome, the calculated ATE not only includes the effect of treatment on the outcome, but also includes the effect of confounders on the outcome, which leads to the **spurious effect**. For example, in the medicine example, age is a confounder. Age affects the recovery rate: in general, young patients have better chance to recover compared to older patients. Age also affects the treatment choice: young patients may prefer to take medicine A while older patients prefer medicine B, or for the same medicine, young patients have a different dosage with elder patients. The observational data is shown in Table 1, and let's estimate ATE according to the above equation: $\hat{\text{ATE}} = \frac{1}{N_A} \sum_{i=1}^{N_A} Y_i^F - \frac{1}{N_B} \sum_{i=1}^{N_B} Y_j^F = 295/350 - 273/350 = 5\%$, where $N_A$ and $N_B$ is the number of patients taking Medicine A and B, respectively. However, we cannot conclude that Treatment A is more effective than Treatment B, because the high average recovery rate of the group taking Treatment A may be caused by the fact that most patients of this group (270 out of 350) are young patients. Thus the effect of age on the recovery rate is the spurious effect, as it is mistakenly counted into the effect of treatment on the outcome.

| Recovery Rate ⟍ Treatment  Age | Treatment A | Treatment B |
|---|---|---|
| Young | 234/270 = 87% | 81/87 = **92**% |
| Older | 55/80 = 69% | 192/263 = **73**% |
| Overall | 289/350 = **83**% | 273/350 = 78% |

Table 1. An example to show the spurious effect of confounder variable *Age*.

From Table 1, we can observe another interesting phenomenon, *Simpson's paradox* (or Simpson's reversal, Yule-Simpson effect, amalgamation paradox, reversal paradox) [18, 42], brought by the confounder. It can be observed that: in both Young and Older patient groups, Medicine B has a higher recovery rate than Medicine A; but when combining these two groups, Medicine A is the one with a higher recovery rate. This paradox is caused by the confounder variable: When compare the recovery rate in the whole group, most of the people taking medicine A are young, and the comparison shown in the table fails to eliminate the effect of age on the recovery rate.

In addition to the spurious effect in treatment effect estimation, confounders also cause problems in counterfactual outcome estimation. As shown in Eqn. (7), counterfactual outcome estimation is an alternative way to
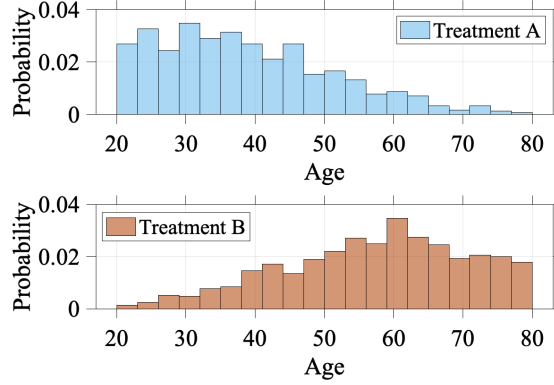
Fig. 1. An example to show the selection bias caused by confounder variable *Age*.

estimate the ATE. Confounder variables cause selection bias, which makes the counterfactual outcome estimation more difficult.

**Selection bias** is the phenomenon that the distribution of the observed group is not representative to the group we are interested in, i.e., $p(X_{obs}) \neq p(X_*)$, where $p(X_{obs})$ and $p(X_*)$ are the distributions of the variables in the observed group and the interested group, respectively. Confounder variables affect units' treatment choices, which leads to the selection bias. In the medicine example, age is a confounder variable, so that people of different ages have different treatment preferences. Fig 1 shows the age distribution of the observed treated/control group. Apparently, the age distribution of the observed treated group is different from the age distribution of the observed control group. This phenomenon exacerbates the difficulty of counterfactual outcome estimation as we need to estimate the control outcome of units in the treated group based on the observed control group, and similarly, estimate the treated outcome of units in the control group based on the observed treated group. If we directly train the potential outcome estimation model $\hat{Y}(x, w) = f_w(x)$ on the data with $W = w$ without handling the selection bias, the trained model would work poorly in estimating the potential outcome of $W = w$ for the units in the other group. This problem brought by the selection is also named as covariate shift in the Machine Learning community.

Handing the problems caused by confounder variables is the essential part of causal inference, and the procedure of handing confounder variables is called *adjust confounders*. The following part of this section briefly discusses the general solutions to tackle the above two problems caused by confounders under the ignorability assumption. The problem when there exists unobserved confounders will be discussed in Section 4.2.

To solve the spurious effect problem, we should take the effect of confounder variables on outcomes into consideration. A general approach along this direction first estimates the treatment effect conditioning on the confounder variables and then conducts weighted averaging over the confounder according to its distribution. To be more specific,

$$
\begin{aligned}
\hat{\text{ATE}} &= \sum_x p(x)\mathbb{E}[Y^F|X = x, W = 1] - \sum_x p(x)\mathbb{E}[Y^F|X = x, W = 0] \\
&= \sum_{\mathcal{X}^*} p(X \in \mathcal{X}^*)\left(\frac{1}{N_{\{i:X_i \in \mathcal{X}^*, W_i=1\}}} \sum_{\{i:X_i \in \mathcal{X}^*, W_i=1\}} Y_i^F\right) - \sum_{\mathcal{X}^*} p(X \in \mathcal{X}^*)\left(\frac{1}{N_{\{j:X_j \in \mathcal{X}^*, W_j=1\}}} \sum_{\{j:X_j \in \mathcal{X}^*, W_j=0\}} Y_j^F\right),
\end{aligned}
\tag{8}
$$

where $\mathcal{X}^*$ is a set of $X$ values, $p(X \in \mathcal{X}^*)$ is the probability of the background variables in $\mathcal{X}^*$ over the whole population, $\{i : x_i \in \mathcal{X}^*, W_i = w\}$ is the subgroup of units whose background variable values belong to $\mathcal{X}^*$ and

treatment is equal to $w$. Stratification, which will be discussed in details later, is a representative method of this category.

For the selection bias problem, there are two general approaches to solve it. The first general approach handles selection bias by creating a pseudo group which is approximately close to the interested group. Possible methods include sample re-weighting, matching, tree-based methods, confounder balancing, balanced representation learning methods, multi-task based methods. The created pseudo-group alleviates the negative influence of the selection bias, and better counterfactual outcome estimations can be obtained. The other general approach first trains the base potential outcome estimation models solely on the observed data, and then correct the estimation bias caused by the selection bias. Meta-learning based methods belong to this category.

## 3 CAUSAL INFERENCE METHODS RELYING ON THREE ASSUMPTIONS

In this section, we introduce existing causal inference methods that rely on the three assumptions introduced in Section 2. According to the way to control confounders, we divide these methods into the following categories: (1) Re-weighting methods; (2) Stratification methods; (3) Matching methods; (4) Tree-based methods; (5) Representation based methods; (6) Multi-task methods; and (7) Meta-learning methods.

### 3.1 Re-weighting Methods

Due to the existence of confounders, the covariate distributions of the treated group and control group are different, which leads to the *selection bias* problem as described in Section 2.4. In other words, the treatment assignment is correlated with covariates in the observational data. Sample re-weighting is an effective approach to overcome the selection bias. By assigning appropriate weight to each unit in the observational data, a pseudo-population can be created on which the distributions of the treated group and control group are similar.

In sample re-weighting methods, a key concept is *balancing score*. Balancing score $b(x)$ is a general weighting score, which is the function of $x$ satisfying: $W \perp\!\!\!\perp x | b(x)$ [59], where $W$ is the treatment assignment and $x$ is the background variables. There are various designs of the balancing score, and apparently, the most trivial balancing score is $b(x) = x$. Besides, propensity score is also a special case of balancing score.

DEFINITION 9. *Propensity score: The propensity score is defined as the conditional probability of treatment given background variables [106]:*

$$e(x) = Pr(W = 1 | X = x) \tag{9}$$

In detail, a propensity score indicates the probability of a unit being assigned to a particular treatment given a set of observed covariates. Balancing scores that incorporate propensity score are the most common approach.

A summarization of the algorithms mentioned in this section is shown in Fig 2. The propensity score based sample reweighting will be introduced in the next section, followed by methods that weigh both samples and the covariates.

*3.1.1 Propensity score based sample re-weighting.* Propensity scores can be used to reduce selection bias by equating groups based on these covariates. Inverse propensity weighting (IPW) [105, 106], also named as inverse probability of treatment weighting (IPTW), assigns a weight $r$ to each sample:

$$r = \frac{W}{e(x)} + \frac{1-W}{1-e(x)}, \tag{10}$$

where $W$ is the treatment assignment ($W = 1$ denotes being treated group; $W = 0$ denotes the control group), and $e(x)$ is the propensity score defined in Eqn. (9).

After re-weighting, the IPW estimator of average treatment effect (ATE) is:

$$\hat{\text{ATE}}_{IPW} = \frac{1}{n} \sum_{i=1}^{n} \frac{W_i Y_i^F}{\hat{e}(x_i)} - \frac{1}{n} \sum_{i=1}^{n} \frac{(1-W_i) Y_i^F}{1 - \hat{e}(x_i)}, \tag{11}$$
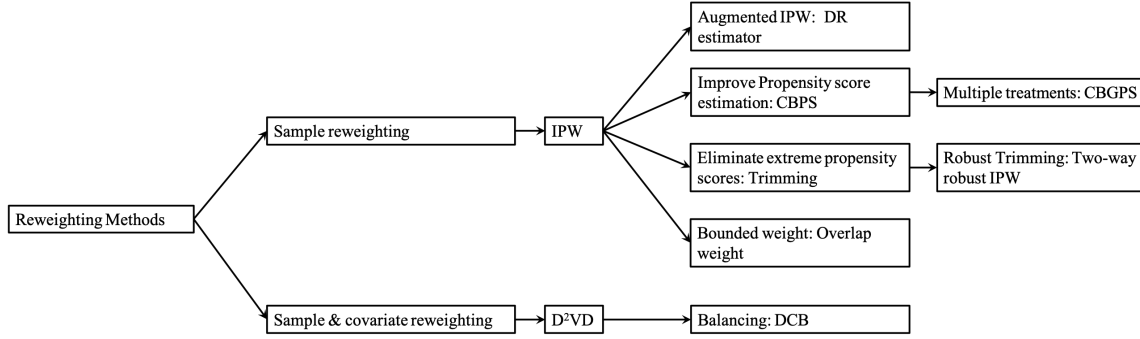
Fig. 2. Categorization of Reweighting Methods

and its normalized version, which is preferred especially when the propensity scores are obtained by estimation [58]:

$$\hat{\text{ATE}}_{IPW} = \sum_{i=1}^{n} \frac{W_i Y_i^F}{\hat{e}(x_i)} \bigg/ \sum_{i=1}^{n} \frac{W_i}{\hat{e}(x_i)} - \sum_{i=1}^{n} \frac{(1-W_i)Y_i^F}{1-\hat{e}(x_i)} \bigg/ \frac{(1-W_i)}{1-\hat{e}(x_i)}. \tag{12}$$

Both large and small sample theory show that adjustment for the scalar propensity score is enough to remove bias due to all observed covariates [106]. The propensity score can be used to balance the covariates in the treatment and control groups and therefore reduce the bias through matching, stratification (subclassification), regression adjustment, or some combination of all three. [31] discusses the use of propensity score to reduce the bias, which also provides examples and detailed discussions.

However, in practice, the correctness of the IPW estimator highly relies on the correctness of the propensity score estimation, and slightly misspecification of propensity scores would cause ATE estimation error dramatically [57]. To handle this dilemma, Doubly Robust estimator (DR) [103], also named as Augmented IPW (AIPW), is proposed. DR estimator combines the propensity score weighting with the outcome regression, so that the estimator is robust even when one of the propensity score or outcome regression is incorrect (but not both). In detail, the DR estimator is formalized as:

$$\begin{aligned}
\hat{\text{ATE}}_{DR} &= \frac{1}{n} \sum_{i=1}^{n} \left\{ \left[ \frac{W_i Y_i^F}{\hat{e}(x_i)} - \frac{W_i - \hat{e}(x_i)}{\hat{e}(x_i)} \hat{m}(1, x_i) \right] - \left[ \frac{(1-W_i)Y_i^F}{1-\hat{e}(x_i)} - \frac{W_i - \hat{e}(x_i)}{1-\hat{e}(x_i)} \hat{m}(0, x_i) \right] \right\} \\
&= \frac{1}{n} \sum_{i=1}^{n} \left\{ \hat{m}(1, x_i) + \frac{W_i(Y_i^F - \hat{m}(1, x_i))}{\hat{e}(x_i)} - \hat{m}(0, x_i) - \frac{(1-W_i)(Y_i^F - \hat{m}(0, x_i))}{1-\hat{e}(x_i)} \right\},
\end{aligned} \tag{13}$$

where $\hat{m}(1, x_i)$ and $\hat{m}(0, x_i)$ are the regression model estimations of treated and control outcomes. The DR estimator is consistent and therefore asymptotically unbiased, if either the propensity score is correct or the model correctly reflects the true relationship among exposure and confounders with the outcome [38]. In reality, one definitely cannot guarantee whether one model can accurately explain the relationship among variables. The combination of outcome regression with weighting by propensity score ensures that the estimators are robust to misspecification of one of these models [12, 101, 103, 117].

The DR estimator consults outcomes to make the IPW estimator robust when propensity score estimation is not correct. An alternative way is to improve the estimation of propensity scores. In the IPW estimator, propensity score serves as both the probability of being treated and the covariate balancing score, covariate balancing propensity score (CBPS) [57] is proposed to exploit such dual characteristics. In particular, CBPS estimates

propensity scores by solving the following problem:

$$\mathbb{E}\left[\frac{W_i \tilde{x}_i}{e(x_i; \beta)} - \frac{(1 - W_i)\tilde{x}_i}{1 - e(x_i; \beta)}\right] = 0, \tag{14}$$

where $\tilde{x}_i = f(x_i)$ is a predefined vector-valued measurable function of $x_i$. By solving the above problem, CBPS directly constructs the covariate balancing score from the estimated parametric propensity score, which increase the robustness to the misspecification of the propensity score model. An extension of CBPS is the covariate balancing generalized propensity score (CBGPS) [39], which enables to handle the treatment with continuous value. Due to the continuous valued treatment, it's hard to directly minimized the covariates distribution distance between the control and treated group. CBGPS solves this problem by mitigating the definition of the balancing score. Based on the definition, the treatment assignment is conditionally independent of the background variables, CBGPS directly minimize the correlation between the treatment assignment and the covariates after weighting. In specific, the objective of CBGPS is to learn a propensity score based weight so that the weighted correlation between the treatment assignment and the covariates are minimized:

$$\mathbb{E}\left(\frac{p(t^*)}{p(t^*|x^*)}t^*x^*\right) = \int \left\{\int \frac{p(t^*)}{p(t^*|x^*)}t^*dP(t^*|x^*)\right\} x^*dP(x^*)$$
$$= \mathbb{E}(t^*)\mathbb{E}(x^*) = 0, \tag{15}$$

where $p(t^*|x^*)$ is the propensity score, and $\frac{p(t^*)}{p(t^*|x^*)}$ is the balancing weight, $t^*$ and $x^*$ is the treatment assignment and the background variables after centering and orthogonalizing (i.e., normalization). In a nutshell, both CBPS and CBGPS learns the propensity score based sample weight directly towards the covariate balancing goal, which can alleviates negative effect brought by model misspecification of propensity score.

Another drawback of the IPW estimator is that it might be unstable if the estimated propensity scores are small. If the probability of either treatment assignment is small, the logistic regression model can become unstable around the tails, causing the IPW to also be less stable. To overcome this issue, trimming is routinely employed as a regularization strategy, which eliminates the samples whose propensity scores are less than a pre-defined threshold [73]. However, this approach is highly sensitive to the amount of trimming [82]. Also, theoretical results in [82] show that the small probability of propensity scores and the trimming procedure may result in different non-Gaussian asymptotic distribution of IPW estimator. Based on this observation, a two-way robustness IPW estimation algorithm is proposed in [82]. This method combines subsampling with a local polynomial regression based trimming bias corrector, so that it is robust to both small propensity score and the large scale of trimming threshold. An alternative approach to overcome the instability of IPW under small propensity scores is to redesign the sample weight so that the weight is bounded. In [75], the overlap weight is proposed, in which each unit's weight is proportional to the probability of that unit being assigned to the opposite group. In detail, the overlap weight $h(x)$ is defined as $h(x) \propto 1 - e(x)$, where $e(x)$ is the propensity score. The overlap weight is bounded within the interval $[0, 0.5]$, and thus it is less sensitive to extreme vale of propensity score. Recent theoretical results show that the overlap weight has the minimum asymptotic variance among all balancing weights [75].

*3.1.2 Confounder balancing.* The aforementioned sample re-weighting methods could achieve balance in the sense that the observed variables are considered equally as confounders. However, in real cases, not all the observed variables are confounders. Some of the variables, named as adjustment variables, are only predictive to the outcome, and some others might be irrelevant variables. Adjusting on the adjustment variables by Lasso, although it cannot reduce the bias, helps decrease the variance [17, 116]. While including the irrelevant variables would cause overfitting.

Based on the separateness assumption that the observed variables can be decomposed into confounders, adjusted variables and the irrelevant variables, in [69], the Data-Driven Variable Decomposition (D$^2$VD) algorithm is

proposed to distinguish the confounders and adjustment variables, and meanwhile, eliminate the irrelevant variables. In detail, the adjusted outcome is written as:

$$Y^*_{\text{D}^2\text{VD}} = \left(Y^F - \phi(\mathbf{z})\right) \frac{W - p(x)}{p(x)(1 - p(x))}, \tag{16}$$

where $\mathbf{z}$ is the adjustment variables. Therefore, the ATE estimator of $\text{D}^2\text{VD}$ is:

$$\text{ATE}_{\text{D}^2\text{VD}} = \mathbb{E}\left[\left(Y^F - \phi(\mathbf{z})\right) \frac{W - p(x)}{p(x)(1 - p(x))}\right]. \tag{17}$$

To get $\text{ATE}_{\text{D}^2\text{VD}}$, the $Y^*_{\text{D}^2\text{VD}}$ is regressed on all observed variables. The objective function is $l_2$ loss between $Y^*_{\text{D}^2\text{VD}}$ and the linear regression function on all observed variables, along with sparse regularization to distinguish the confounder, adjusted variables, and irrelevant variables. However, little prior knowledge about the interactions among observed variables is provided in practice, and the data are usually high-dimensional and noisy. To solve this, Differentiated Confounder Balancing (DCB) algorithm [68] is proposed to select and differentiate confounders to balance the distributions. Overall, DCB balances the distributions by re-weighting both the samples and confounders.

## 3.2 Stratification Methods

Stratification, also named as *subclassification* or *blocking* [59], is a representative method to adjust the confounders. The idea of stratification is to adjust the bias that stems from the difference between the treated group and the control group by splitting the entire group into homogeneous subgroups (blocks). Ideally, in each subgroup, the treated group and the control group are similar under certain measurements over the covariates, therefore, the units in the same subgroup can be viewed as sampled from the data under randomized controlled trials. Based on the homogeneity of each subgroup, the treatment effect within each subgroup (i.e., CATE) can be calculated through the method developed on RCTs data. After having the CATE of each subgroup, the treatment effect over the interested group can be obtained by combining the CATEs of subgroups belonging to that group, as shown in (8). In the following, we adopt the calculation of ATE as an example. In detail, if we separate the whole dataset into $J$ blocks, the ATE is estimated as:

$$\text{ATE}_{\text{strat}} = \hat{\tau}^{\text{strat}} = \sum_{j=1}^{J} q(j) \left[\bar{Y}_t(j) - \bar{Y}_c(j)\right], \tag{18}$$

where $\bar{Y}_t(j)$ and $\bar{Y}_c(j)$ are the average of the treated outcome and control outcome in the $j$-th block, respectively. $q(j) = \frac{N(j)}{N}$ is the portion of the units in the $j$-th block to the whole units.

Stratification effectively decreases the bias of ATE estimation compared with the difference-estimator where ATE is estimated as: $\text{ATE}_{\text{diff}} = \hat{\tau}^{diff} = \frac{1}{N_t} \sum_{i:W_i=1} Y^{CF}_i - \sum_{i:W_i=0} Y^{CF}_i$. In particular, if we assume the outcome is linear with the covariates, i.e., $\mathbb{E}[Y_i(w)|X_i = x] = \alpha + \tau * w + \beta * x$. The bias of the difference-estimator is:

$$\mathbb{E}[\hat{\tau}^{\text{diff}} - \tau | X, W] = (\bar{X}_t - \bar{X}_c)\beta. \tag{19}$$

While, the bias of the stratification estimator is the weighted average of the within-block bias:

$$\mathbb{E}[\hat{\tau}^{\text{strat}} - \tau | X, W] = \left(\sum_{j=1}^{J} q(j) \left(\bar{X}_t(j) - \bar{X}_c(j)\right)\right)\beta. \tag{20}$$

Compared with the difference estimator, the stratification estimator reduces the bias per covariate by the factor:

$$\gamma_k = \frac{\sum_j q(j) \left(\bar{X}_{t,k}(j) - \bar{X}_{c,k}(j)\right)}{\bar{X}_{t,k} - \bar{X}_{c,k}}, \tag{21}$$

where $\bar{X}_{t,k}(j)$ $(\bar{X}_{c,k}(j))$ is the average of $k$-th covariate of treated (control) group in $j$-th block, and $\bar{X}_{t,k}$ $(\bar{X}_{c,k})$ is the average of $k$-th covariate in the whole treated (control) group.

The key component of stratification methods is how to create the blocks and how to combine the created blocks. The equal frequency [106] is a common strategy to create blocks. Equal-frequency approach split the block by the appearance probability, such as the propensity score, so that the covariates have the same appearance probability (i.e., the propensity score) in each subgroup (block). The ATE is estimated by weighted average of each block's CATE, with the weight as the fraction of the units in this block. However, this approach suffers from high variance due to the insufficient overlap between treated and control groups in the blocks whose propensity score is very high or low. To reduce the variance, in [55], the blocks, which divided according to the propensity score, are re-weighted by the inverse variance of the block-specific treatment effect. Although this method reduces the variance of equal-frequency method, it unavoidably increases the estimation bias.

The stratification methods described above are all splitting the blocks according to the pre-treatment variables. However, in some real-world applications, it is required to compare the outcome conditioned on some post-treatment variables, denoted as $S$. For example, the "surrogate" markers of disease progression (i.e., intermediate outcome) like CD4 count and measures of viral load in AIDS are the post-treatment variables [40]. In the studies comparing drugs for AIDS patients, the researchers are interested in the effect of AIDS drugs on group with CD4 count lower than 200 cell/mm$^3$. However, directly comparing the observed outcomes on the group with $S^{obs} < 200$ is not the true effect because the compared two subgroups: $\{i : W_i = 1, S^{obs} < 200\}$ and $\{j : W_j = 0, S^{obs} < 20\}$, where $S^{obs}$ is the observed post-treatment values, have great discrepency if the treatment has effect on the intermediate results. To solve this, principle stratification [40] constructs the subgroup based on the potential values of the pre-treatment variables. Analogous to the potential outcome defined in 2.1, potential pre-treatment variables value, denoted as $S(W = w)$, is the potential value of $S$ under treatment with value $w$. With the nature assumption that potential value of $S$ is independent of the treatment assignment, the treatment effect of subgroup can be obtained by comparing the outcomes of two sets: $\{Y_i^{obs} : W_i = 1, S_i(W_i = 1) = v_1, S_i(W_i = 0) = v_2\}$ and $\{Y_j^{obs} : W_j = 0, S_j(W_j = 1) = v_1, S_j(W_j = 0) = v_2\}$, where $v_1$ and $v_2$ are two post-treatment values. The comparison based on the potential values of post-treatment variables ensures that the compared two set are similar, so that the obtained treatment effect is the true effect.

## 3.3 Matching Methods

As mentioned previously, missing counterfactuals and confounder bias are two major challenges in treatment effect estimation. Matching based approaches provide a way to estimate the counterfactual and, at the same time, reduce the estimation bias brought by the confounders. In general, the potential outcomes of the $i$-th unit estimated by matching are [1]:

$$\hat{Y}_i(0) = \begin{cases} Y_i & \text{if } W_i = 0, \\ \frac{1}{\#\mathcal{J}(i)} \sum_{l \in \mathcal{J}(i)} Y_l & \text{if } W_i = 1; \end{cases} \qquad \hat{Y}_i(1) = \begin{cases} \frac{1}{\#\mathcal{J}(i)} \sum_{l \in \mathcal{J}(i)} Y_l & \text{if } W_i = 0, \\ Y_i & \text{if } W_i = 1; \end{cases} \tag{22}$$

where $\hat{Y}_i(0)$ and $\hat{Y}_i(1)$ are the estimated control and treated outcome, $\mathcal{J}(i)$ is the matched neighbors of unit $i$ in the opposite treatment group [11].

The analysis of the matched sample can mimic that of an RCT: one can directly compare outcomes between the treated and control group within the matched sample. In the context of an RCT, one expects that, on average, the distribution of covariates will be similar between treated and control groups. Therefore, matching can be used to reduce or eliminate the effects of confounding when using observational data to estimate treatment effects [11].

*3.3.1 Distance Metric.* Various distances have been adopted to compare the closeness between units [43], such as the widely used Euclidean distance [110] and Mahalanobis distance [113]. Meanwhile, many matching methods

develop their own distance metrics, which can be abstracted as: $D(\mathbf{x}_i, \mathbf{x}_j) = ||f(\mathbf{x}_i) - f(\mathbf{x}_j)||_2$. The existing distance metrics mainly vary in how to design the transformation function $f(\cdot)$.

*Propensity score based transformation.* Original covariates of units can be represented by propensity scores. As a result, the similarity between two units can be directly calculated as: $D(\mathbf{x}_i, \mathbf{x}_j) = |e_i - e_j|$, where $e_i$, and $e_j$ are the propensity scores of $\mathbf{x}_i$ and $\mathbf{x}_j$, respectively. Later, the linear propensity score based distance metric is also proposed, which is defined as $D(\mathbf{x}_i, \mathbf{x}_j) = |\text{logit}(e_i) - \text{logit}(e_j)|$. This improved version is recommended since it can effectively reduce the bias [131]. Furthermore, the propensity score based distance metric can be combined with other existing distance metrics, which provides a fine-grained comparison. In [113], when the difference of two unit's propensity scores is within a certain range, they are further compared with other distances on some key covariates. Under this metric, the closeness of two units contains two criteria: they are relatively close under propensity score measure, and they particularly similar under the comparison of the key covariates [131].

*Other transformations.* Propensity score only adopts the covariate information, while some other distance metrics are learned by utilizing both the covariates and the outcome information so that the transformed space can preserve more information. One representative metric is the prognosis score [49], which is the estimated control outcome. The transformation function is represented as: $f(x) = \hat{Y}_c$. However, the performance of the prognosis score relies on modeling the relationship between the covariates and control outcomes. Moreover, the prognosis score only takes the control outcome into consideration and ignores the treated outcome. The Hilbert-Schmidt Independence Criterion based nearest neighbor matching (HSIC-NNM) proposed in [26] could overcome the drawbacks of prognosis score. HSIC-NNM learns two linear projections for control outcome estimation task and treated outcome estimation task separately. To fully explore the observed control/treated outcome information, the parameters of linear projection is learned by maximizing the nonlinear dependency between the projected subspace and the outcome: $M_w = \arg\max_{M_w} \text{HSIC}(\mathbf{X}_w M_w, Y_w^F) - \mathcal{R}(M_w)$, where $w = 0, 1$ represent the control group and treated group, respectively. $\mathbf{X}_w M_w$ is the transformed subspace with the transformation function as: $f(x) = xM_w$. $Y_w^F$ is the observed control/treated outcome, and $\mathcal{R}$ is the regularization to avoid overfitting. The objective function ensures the learned transformation functions project the original covariates to an information subspace where similar units will have similar outcomes.

Compared with propensity score based distance metric that focuses on balancing, prognosis score and HSIC-NNM focus on embedding the relationship between the transformed space and the observed outcome. These two lines of methods have different advantages, and some recent work tries to integrate these advantages together. In [77], the balanced and nonlinear representation (BNR) is proposed to project the covariates into a balanced low-dimensional space. In detail, the parameters in the nonlinear transformation function is learned by jointly optimizing the following two objectives: (1) Maximizing the differences of noncontiguous-class scatter and within-class scatter so that the units with the same outcome prediction shall have similar representations after transformation; (2) Minimizing the maximum mean discrepancy between the transformed control and outcome group in order to get the balanced space after transformation. A series of works that have similar objectives but vary in balancing regularization have been proposed, such as using the conditional generative adversarial network to ensure the transformation function blocks the treatment assignment information [74, 151].

The methods mentioned above adopt either one or two transformations for treated and control groups separately. Different from the existing method, Randomized Nearest Neighbor Matching (RNNM) [78] adopts a number of random linear projections as the transformation function, and the treatment effects are obtained as the median treatment effect by nearest-neighbor matching in each transformed subspace. The theoretical motivation of this approach is the Johnson-Lindenstrauss (JL) lemma, which guarantees that the pairwise similarity information of the points in the high-dimensional space can be preserved through random linear projection. Powered by the JL lemma, RNNM ensembles the treatment effect estimation results of several linear random transformations.
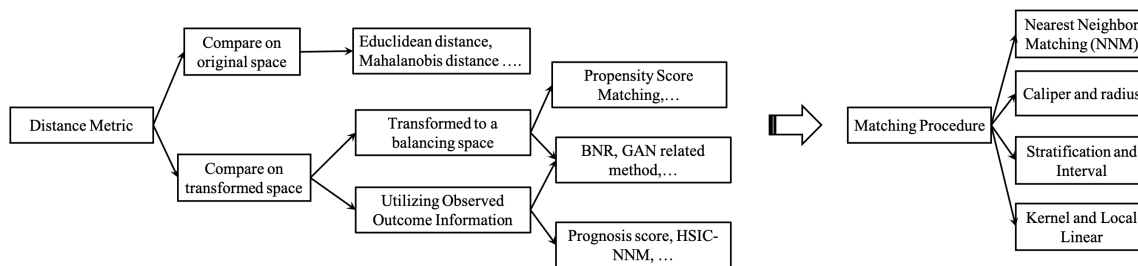
Fig. 3. Categorization of Matching Methods.

*3.3.2 Choosing a Matching Algorithm.* After defining the similarity metric, the next step is to find the neighbors. In [23], existing matching algorithms are divided into four essential approaches, including the nearest neighbor matching, caliper, stratification and kernel, as shown in Fig. 3. The most straightforward matching estimator is nearest neighbour matching (NNM). In particular, a unit in the control group is chosen as the matching partner for a treated unit, so that they are closest based on a similarity score (e.g., propensity score). The NNM has several variants like NNM with replacement and NNM without replacement. Treated units are matched to one control, called pair matching or 1-1 matching, or treated units are matched to two controls, called 1-2 matching, and so on. It's a trade-off to determine the number of neighbors, since a large number of neighbors may result in the treatment effect estimator with high bias but low variance, while small number results in low bias but high variance. It is known, however, that the optimal structure is a full matching in which a treated unit may have one or several controls or a control may have one or several treated units [43].

NNM may have bad matches if the closest partner is far away. One can set a tolerance level on the maximum propensity score distance (caliper) to avoid this problem. Hence, caliper matching is one form of imposing a common support condition.

The stratification matching is to partition the common support of the propensity score into a set of intervals and then to take the mean difference in outcomes between treated and control observations in order to calculate the impact within each interval. This method is also known as interval matching, blocking and subclassification [108].

The matching algorithms discussed above have in common that only a few observations in the control group are used to create the counterfactual outcome of a treatment observation. Kernel matching (KM) and local linear matching (LLM) are nonparametric matching that use weighted averages of observations in the control group to create the counterfactual outcome. Thus, one major advantage of these approaches is the lower variance, because we use more information to create counterfactual outcome.

Here, we also want to introduce another matching method called Coarsened Exact Matching (CEM) proposed in [56]. Because either the 1-k matching or the full matching fails to consider the extrapolation region, where few or no reasonable matches exist in the other treatment group, CEM was proposed to handle this problem. CEM first coarsen the selected important covariate,i.e., discretization, and then perform exact matching on the coarsened covariates. For example, if the selected covariates are age (age > 50 is 1, and others are 0) and gender (female as 1, and male as 0). A female patient with age 50 in the treated group is represented by the coarsen covariates as (1, 1). She will only match the patients in the treated group with exactly the same coarsened covariates value. After exact matching, the whole data is separated into two subsets. In one subset, every unit has its exact matched neighbors and it is the opposite in the other subset which contains the units in the extrapolation region. The outcomes of units in the extrapolation region are estimated by the outcome prediction model trained on the matched subset. So far, the treatment effect on the two subsets can be estimated separately, and the final step is to combine treatment effect on the two subsets by weighted average.

We have provided several different matching algorithms, but the most important question is how we should select a perfect matching method. Asymptotically all matching methods should yield the same results as the sample size grows and they will become closer to comparing only exact matches [128]. When we only have small samples size, this choice will be important [52]. There is one trade-off between bias and variance.

*3.3.3 Variables to include.* The above two subsections illustrate the key steps in matching procedure, and in this subsection, we briefly discuss what kinds of variables should be included in the matching, a.k.a feature selection, to improve the matching performance. Many literatures [41, 52, 112] suggest to include as many variables that are related to the treatment assignment and the outcome as possible, in order to satisfy the strong ignorablity assumption. However, post-treatment variables, which are the variables affected by the treatment assignment, should be excluded in the matching procedure [107]. Moreover, besides the post-treatment variables, researchers also suggest excluding the instrumental variables [93, 148], because they tend to amplify the bias of treatment effect estimator.

## 3.4 Tree-based Methods

Another popular method in causal inference is based on decision tree learning, which is one of the predictive modeling approaches. Decision tree is a non-parametric supervised learning method used for classification and regression. The goal is to create a model that predicts the value of a target variable by learning simple decision rules inferred from data.

Tree models where the target variable is discrete are called classification trees with prediction error measured based on misclassification cost. In these tree structures, leaves represent class labels and branches represent conjunctions of features that lead to those class labels. Decision trees where the target variable is continuous are called regression trees with prediction error measured by the squared difference between the observed and predicted values. The term Classification And Regression Tree (CART) analysis is an umbrella term used to refer to both of the above procedures [21]. In CART model, the data space is partitioned and a simple prediction model for each partition space is fitted, and therefore every partitioning can be represented graphically as a decision tree [80].

For estimating heterogeneity in causal effects, a data-driven approach [9] based on CART is provided to partition the data into subpopulations that differ in the magnitude of their treatment effects. The valid confidence intervals can be created for treatment effects, even with many covariates relative to the sample size, and without "sparsity" assumptions. This approach is different from conventional CART in two aspects. First, it focuses on estimating conditional average treatment effects instead of directly predicting outcomes as in the conventional CART. Second, different samples are used for constructing the partition and estimating effects each subpopulation, which is referred to as the honest estimation. However, in conventional CART, the same samples are used for these two tasks.

In CART, a tree is built up until a splitting tolerance is reached. There is only one tree, and it is grown and pruned as needed. However, BART is an ensemble of trees, so it is more comparable to random forests. A Bayesian "sum-of-trees" model called Bayesian Additive Regression Trees (BART) is developed in [28] [29]. Every tree in BART model is a weak learner, and it is constrained by a regularization prior. Information can be extracted from the posterior by a Bayesian backfitting MCMC algorithm. BART is a nonparametric Bayesian regression model, which uses dimensionally adaptive random basis elements. Let $W$ be a binary tree which has a set of interior node decision rules and terminal nodes, and let $M = \{\mu_1, \mu_2, ..., \mu_B\}$ be parameters associated with each of the $B$ terminal nodes for $W$. We use $g(x; W, M)$ to assign a $\mu_b \in M$ to input vector $x$. The sum-of-trees model can be expressed as:

$$Y = g(x; W_1, M_1) + g(x; W_2, M_2) + \cdots + g(x; W_m, M_m) + \varepsilon, \tag{23}$$

$$\varepsilon \sim N(0, \sigma^2), \tag{24}$$

BART has a couple of advantages. It is very easy to implement and only needs to plug in the outcome, treatment assignment, and confounding covariates. In addition, it doesn't require any information about how these variables are parametrically related, so that it requires less guess when fitting the model. Moreover, it can deal with a mass of predictors, yield coherent uncertainty intervals, and handle continuous treatment variables and missing data [53].

BART is proposed to estimate average causal effects. In fact, it can also be used to estimate individual-level causal effects. BART not only can easily identify the heterogeneous treatment effects, but also get more accurate estimates of average treatment effects compared to other methods like propensity score matching, propensity score weighting, and regression adjustment in the nonlinear simulation situations examined [53].

In most previous methods, the prior distribution over treatment effects is always induced indirectly, which is difficult to be attained. A flexible sum of regression trees (i.e., a forest) can address this issue by modeling a response variable as a function of a binary treatment indicator and a vector of control variables [48]. This approach interpolates between two extremes: entirely and separately modeling the conditional means of treatment and control, or only the treating treatment assignment as another covariate.

Random forest is a classifier consisting of a combination of tree predictors, in which each tree depends on a random vector that is independently sampled and has the identical distribution for all trees [20]. This model can also be extended to estimate heterogeneous treatment effects based on the Breiman's random forest algorithm [141]. Trees and forests can be considered as nearest neighbor methods with an adaptive neighborhood metric. Tree-based methods seek to find training examples that are close to a point $x$, but now closeness is defined with respect to a decision tree. And the closest points to $x$ are those that fall in the same leaf as it. The advantage of using trees is that their leaves can be narrower along the directions where the signal is changing fast and wider along the other directions, potentially leading to a substantial increase in power when the dimension of the feature space is even moderately large.

The tree-based framework also can be extended to uni- or multi-dimensional treatments [142]. Each dimension can be discrete or continuous. A tree structure is used to specify the relationship between user characteristics and the corresponding treatment. This tree-based framework is robust to model misspecification and highly flexible with minimal manual tuning.

## 3.5  Representation Learning Methods

### 3.5.1  *Balanced representation learning.* The most basic assumption used in statistical learning theory is that training data and test data are drawn from the same distribution. However, in most practical cases, the test data are drawn from a distribution that is only related, but not identical, to the distribution of the training data. In causal inference, this is also a big challenge. Unlike the randomized control trials, the mechanism of treatment assignment is not explicit in observational data. Therefore, interventions of interest are not independent of the property of the subjects. For example, in an observational study of the treatment effect of a medicine, the medicine is assigned to individuals based on several factors, including the known confounders and some unknown confounders. As a result, the counterfactual distribution will generally be different from the factual distribution. Thus, it is necessary to predict counterfactual outcomes by learning from the factual data, which converts the causal inference problem to a domain adaptation problem.

Extracting effective feature representations is critical for domain adaptation. A model [14] with a generalization bound is proposed to formalize this intuition theoretically, which can not only explicitly minimize the difference between the source and target domains, but also maximize the margin of the training set. Building on this work [14], the discrepancy distance between distributions is tailored to adaptation problems with arbitrary loss functions [83]. In the following discussions, the discrepancy distance plays an important role in addressing the domain adaptation problem in causal inference.

So far, we can see a clear connection between counterfactual inference and domain adaptation. An intuitive idea is to enforce the similarity between the distributions of different treatment groups in the representation space. The learned representations trade-off three objectives: (1) low-error prediction over the factual representation, (2) low-error prediction over counterfactual outcomes by taking into account relevant factual outcomes, and (3) the distance between the distribution of treatment population and that of control population [60]. Following this motivation, [122] gives a simple and intuitive generalization-error bound. It shows that the expected ITE estimation error of representation is bounded by a sum of the standard generalization-error of that representation and the distance between the treated and control distributions based on representation. Integral probability metric (IPM) is used to measure the distances between distributions, and explicit bounds are derived for the Wasserstein distance and Maximum Mean Discrepancy (MMD) distance. The goal is to find a representation $\Phi : X \rightarrow R$ and hypothesis $h : X \times \{0, 1\} \rightarrow Y$ that minimizes the following objective function:

$$\min_{h,\Phi} \frac{1}{n} \sum_{i=1}^{n} r_i \cdot L(h(\Phi(x_i), W_i), y_i) + \lambda \cdot R(h) + \alpha \cdot IPM_G(\{\Phi(x_i)\})_{i:W_i=0}, \{\Phi(x_i)\})_{i:W_i=1}), \qquad (25)$$

where $w_i = \frac{W_i}{2u} + \frac{1-W_i}{2(1-u)}$, $u = \frac{1}{n} \sum_{i=1}^{n} W_i$, and the weights $r_i$ compensate for the difference in treatment group size. $R$ is a model complexity term. Given two probability density functions $p$, $q$ defined over $S \subseteq R^d$, and a function family $G$ of functions $g : S \rightarrow R$, the IPM is defined as:

$$IPM_G(p, q) := \sup_{g \in G} |\int_S g(s)(p(s) - q(s))ds|. \qquad (26)$$

This model allows for learning complex nonlinear representations and hypotheses with large flexibility. When the dimension of $\Phi$ is high, it risks losing the influence of $t$ on $h$ if the concatenation of $\Phi$ and $W$ is treated as input. To address this problem, one approach is to parameterize $h_1(\Phi)$ and $h_0(\Phi)$ as two separate "heads" of the joint network. $h_1(\Phi)$ is used to estimate the outcome under treatment and $h_0(\Phi)$ is for the control group. Each sample is used to update only the head corresponding to the observed treatment. The advantage is that statistical power is shared in the common representation layers and the influence of treatment is retained in the separate heads [122]. This model can also be extended to any number of treatments, as described in the perfect match (PM) approach [120]. Following this idea, a few improved models have been proposed and discussed. For example, [61] brings together shift-invariant representation learning and re-weighting methods. [51] presents a new context-aware weighting scheme based on the importance sampling technique, on top of representation learning, to alleviate the selection bias problem in ITE estimation.

Existing ITE estimation methods mainly focus on balancing the distributions of control and treated groups, but ignore the local similarity information that provides meaningful constraints on the ITE estimation. In [149, 150], a local similarity preserved individual treatment effect (SITE) estimation method is proposed based on deep representation learning. SITE preserves local similarity and balances data distributions simultaneously. The framework of SITE contains five major components: representation network, triplet pairs selection, position-dependent deep metric (PDDM), middle point distance minimization (MPDM), and the outcome prediction network. To improve the model efficiency, SITE takes input units in a mini-batch fashion, and triplet pairs could be selected from every mini-batch. The representation network learns latent embeddings for the input units. With the selected triplet pairs, PDDM and MPDM can preserve the local similarity information and meanwhile achieve the balanced distributions in the latent space. Finally, the embeddings of mini-batch are fed forward to a dichotomous outcome prediction network to get the potential outcomes. The loss function of SITE is as follows:

$$L = L_{FL} + \beta L_{PDDM} + \gamma L_{MPDM} + \lambda ||M||_2 \qquad (27)$$

where $L_{FL}$ is the factual loss between the estimated and observed factual outcomes. $L_{PDDM}$ and $L_{MPDM}$ are the loss functions for PDDM and MPDM, respectively. The last term is $L_2$ regularization on model parameters $M$ (except the bias term).

Most models focus on covariates with numerical values, while how to handle covariates with textual information for treatment effect estimation is still an open question. One major challenge is how to filter out the nearly instrumental variables which are the variables more predictive to the treatment than the outcome. Conditioning on those variables to estimate the treatment effect would amplify the estimation bias. To address this challenge, a conditional treatment-adversarial learning based matching (CTAM) method is proposed in [151]. CTAM incorporates the treatment-adversarial learning to filter out the information related to nearly instrumental variables when learning the representations, and then it performs matching among the learned representations to estimate the treatment effects. The CTAM contains three major components: text processing, representation learning, and conditional treatment discriminator. Through the text processing component, the original text is transformed into vectorized representation $S$. After that, $S$ is concatenated with the non-textual covariates $X$ to construct a unified feature vector, which is then fed into the representation neural network to get the latent representation $Z$. After learning the representation, $Z$, together with potential outcomes $Y$, are fed into the conditional treatment discriminator. During the training procedures, the representation learner plays a minimax game with the conditional treatment discriminator: By preventing the discriminator from assigning correct treatment, the representation learner can filter out the information related to nearly instrumental variables. The final matching procedure is performed in the representation space $Z$. The conditional treatment-adversarial learning helps reduce the bias of treatment effect estimation.

Compared to the above regression-based methods after representation learning, matching method is more interpretable, because any sample's counterfactual outcome is directly set to be the factual outcome of its nearest neighbor in the group receiving the opposite treatment. Nearest neighbor matching (NNM) sets the counterfactual outcome of any treatment (control) sample to be equal to the factual outcome of its nearest neighbor in the control (treatment) group. Although being simple, flexible and interpretable, most NNM approaches could be easily misled by variables that do not affect the outcome. To address this challenge, matching can be performed on subspaces that are predictive of the outcome variable for both the treatment group and the control group. Applying NNM in the learned subspaces leads to a more accurate estimation of the counterfactual outcomes and therefore the accurate estimation of treatment effects. [26] estimates the counterfactual outcomes of treatment samples by learning a projection matrix that maximizes the nonlinear dependence between the subspace and outcome variable for control samples. Then it directly applies the learned projection matrix to all the samples and finds every treatment sample's matched control sample in the subspace.

## 3.6 Multitask Learning Methods

Treatment group and control group always share some common features except for their idiosyncratic characteristics. Naturally, causal inference can be conceptualized as a multitask learning problem with a set of shared layers for treated group and control group together, and a set of specific layers for treated group and control group separately. The impact of selection bias in multi-task learning problem can be alleviated via a propensity-dropout regularization scheme [4], in which the network is thinned for every training example via a dropout probability that depends on the associated propensity score. The dropout probability is higher for subjects with features that belong in a region of poor overlap in the feature space between treatment and control group.

The Bayesian method also can be extended under multi-task model. A nonparametric Bayesian method [3] uses a multi-task Gaussian process with a linear coregionalization kernel as a prior over the vector-valued reproducing kernel Hilbert space. The Bayesian approach allows computing individualized measures of confidence in our estimates via pointwise credible intervals, which are crucial for realizing the full potential of precision medicine.

The impact of selection bias is alleviated via a risk-based empirical Bayes method for adapting the multi-task GP prior, which jointly minimizes the empirical error in factual outcomes and the uncertainty in counterfactual outcomes.

The multi-task model can be extended to multiple treatments even with continuous parameters in each treatment. The dose response network (DRNet) architecture [119] with shared base layers, $N_W$ intermediary treatment layers, and $N_W \times E$ heads for the multiple treatment setting with an associated dosage parameter $s$. The shared base layers are trained on all samples, and the treatment layers are only trained on samples from their respective treatment category. Each treatment layer is further subdivided into $E$ head layers. Each head layer is assigned a dosage stratum that subdivides the range of potential dosages $[a_t, b_t]$ into $E$ partitions of equal width $\frac{b-a}{E}$.

### 3.7 Meta-Learning Methods

When designing the heterogeneous treatment effect estimation algorithms, two key factors should be considered: 1) Control the confounders, i.e., eliminate the spurious correlation between the confounder and the outcome; 2) Give an accurate expression of the CATE estimation [87]. The methods mentioned in the previous sections seek to satisfy the two requirements simultaneously, while meta-learning based algorithms separate them into two steps. In general, the meta-learning based algorithms have the following procedures: (1) Estimate the conditional mean outcome $\mathbb{E}[Y|X = x]$, and the prediction model learned in this step is the base learner. (2) Derive the CATE estimator based on the difference of results obtained from step (1). Existing meta-learning methods include T-learner [70], S-learner [70], X-learner [70], U-learner [87] and R-learner [87], which are introduced in the following.

In detail, the T-learner [70] adopts two trees to estimate the conditional treated/control outcomes, which are denoted as $\mu_0(x) = \mathbb{E}[Y(0)|X = x]$ and $\mu_1(x) = \mathbb{E}[Y(1)|X = x]$, respectively. Let $\hat{\mu}_0(x)$ and $\hat{\mu}_0(x)$ denote the trained tree model on the control/treated group. Then the CATE of T-learner estimation is obtained as: $\hat{\tau}_T(x) = \hat{\mu}_1(x) - \hat{\mu}_0(x)$. T-learner trains two base models for control and treated groups (the name "T" comes from two base model), while S-learner[70] views the treatment assignment as one feature and estimate the combined outcome as: $\mu(x, w) = \mathbb{E}[Y^F|X = x, W = w]$ (The name "S" denotes single). $\mu(x, w)$ can be any base model, and we denote the trained model as $\hat{\mu}(x, w)$. The CATE estimator provided by S-learner is then given as: $\hat{\tau}_S(x) = \hat{\mu}(x, 1) - \hat{\mu}(x, 0)$.

However, T-learner and S-learner highly rely on the performance of the trained base models. When the number of units in two groups are extremely unbalanced (i.e., the number of one group is much larger than the other), the performance of the base model trained on the small group would be poor. To overcome this problem, X-learner [70] is proposed, which adopts information from the control group to give a better estimator on the treated group and vice versa. The cross-group information usage is where X-learner comes from, and the X denotes "cross group". In detail, X-learner contains three key steps. The first step of X-learner is the same as T-learner, and the trained base learners are denoted as $\hat{\mu}_0(x)$ and $\hat{\mu}_1(x)$. In the second step, X-learner calculates the difference between the observed outcome and the estimated outcome as the imputed treatment effect: In the control group, the difference is the estimated treated outcome subtracts the observed control outcome, denoted as: $\hat{D}_i^C = \hat{\mu}_1(x) - Y^F$; Similarly, in the treated group, the difference is formulated as: $\hat{D}_i^T = Y^F - \hat{\mu}_0(x)$. After the difference calculation, the dataset is transformed into two groups with imputed treatment effect: control group: $(X_C, \hat{D}^C)$ and treated group: $(X_T, \hat{D}^T)$. On two imputed datasets, the two base learners of treatment effect $\tau_1(x)$($\tau_0(x)$) are trained with $X_C(X_T)$ as the input and $\hat{D}^C(\hat{D}^T)$ as the output. The last step is to combine the two CATE estimators by weighted average: $\tau_X(x) = g(x)\hat{\tau}_0(x) + (1 - g(x))\hat{\tau}_1(x)$, where $g(x)$ is the weighting function ranging from 0 to 1. Overall, with the cross information usage and the weighted combination of two CATE base estimator, X-learners can handle the case where the number of units in two groups are unbalanced [70].

Different from the regular loss function adopted in X-learner, R-learner, proposed in [87] designs the loss function for CATE estimator based on the Robinson transformation [104]. The character "R" in R-learner denotes the Robinson transformation. The Robinson transformation can be derived by rewriting the observed outcome and the conditional outcome: Rewrite the observed outcome as:

$$Y_i(W = w_i) = \hat{\mu}_0(x_i) + w_i * \tau(x_i) + \epsilon_i(w_i), \tag{28}$$

where $\hat{\mu}_0$ is the already-trained control outcome estimator(base learner), $\tau(x_i)$ is the CATE estimator, and $E[\epsilon_i(w_i)|x_i, w_i] = 0$ (under ignorability). The conditional mean outcome can be also rewritten as:

$$\hat{m}(x_i) = E[Y|X] = \hat{\mu}_0(x_i) + \hat{e}(x_i) * \tau(x_i), \tag{29}$$

where $\hat{e}(x)$ is the already-trained propensity score estimator(base learner). Robinson transformation is obtained by subtracting Eqn. (28) and Eqn. (29):

$$Y_i^F - \hat{m}(x_i) = (w_i - \hat{e}(x_i))\tau(x_i) + \epsilon(w_i) \tag{30}$$

Based on the Robinson transformation, a good CATE estimator should minimize the difference between $Y_i^F - \hat{m}(x_i)$ and $(w_i - \hat{e}(x_i))\tau(x_i)$. Therefore, the objective function of R-learner is as follows:

$$\tau(\cdot) = \arg\min_{\tau} \left\{ \frac{1}{n} \sum_{i=1}^{n} \left( (Y_i^F - \hat{m}(x_i)) - (w_i - \hat{e}(x_i)) \tau(x_i) \right)^2 + \Lambda(\tau(\cdot)) \right\}, \tag{31}$$

where $\hat{m}(x_i)$ and $\hat{e}(x_i)$ are pre-trained outcome estimator and propensity score estimator, and $\Lambda(\tau(\cdot))$ is the regularization on $\tau(\cdot)$.

## 4 METHODS RELAXING THREE ASSUMPTIONS

In Section 3, the causal inference methods based on three assumptions have been introduced in detail, which are the stable unit treatment value assumption (SUTVA), ignorability assumption, and positivity assumption. However, in practice, for some specific applications like social media analysis, which involves dependent network information, special data types (e.g., time series data) or particular conditions (e.g., the existence of unobserved confounders), these three assumptions cannot always hold. In this section, the methods that try to relax certain assumptions will be discussed.

### 4.1 Stable unit treatment value assumption (SUTVA) assumption

Stable Unit Treatment Value Assumption (SUTVA) states that the potential outcomes for any unit do not vary with the treatment assigned to other units, and, for each unit, there are no different forms or versions of each treatment level, which lead to different potential outcomes. This assumption mainly focuses on two aspects: (1) Units are independent and identically distributed (i.i.d.); (2) there only exists a single level for each treatment. An extensive literature exists on making causal inferences under SUTVA, but when considering many real-world situations, it may not always be the case. In the following, SUTVA will be discussed from these two aspects.

The assumption of independent and identically distributed samples is ubiquitous in most causal inference methods, but this assumption cannot hold in many research areas, such as social media analytics [46] [123], herd immunity, and signal processing [140] [133]. Causal inference in non-i.i.d. contexts is challenging due to the presence of both unobserved confounding and data dependence. For example, in social networks, subjects are connected and influenced by each other.

For such network data, SUTVA cannot hold anymore. Under this situation, instances are inherently interconnected with each other through the network structure and hence their features are not independent identically distributed samples drawn from a certain distribution. Applying Graph Convolutional Networks into causal inference model is an approach to handle the network data [46]. In particular, the original features of subjects

and the network structure are mapped to a representation space, in order to get the representation of confounders. Furthermore, the potential outcomes could be inferred using treatment assignments and confounder representations.

The dependence in data often leads to interference because some subjects' treatments can affect others' outcomes [54, 88]. This difficulty can impede the identification of causal parameters of interest. Extensive work has been developed on identification and estimation of causal parameters under interference [54, 88, 95, 136]. For this problem, a strategy proposed by [125] is to use segregated graphs [127], a generalization of latent projection mixed graphs [139], to represent causal models.

Modeling time series data is another important problem in causal inference, which does not satisfy the independent and identically distributed assumption. Most of the existing methods use regression models for this problem but the accuracy of inference depends greatly on whether the model fits the data. Therefore, selecting a right and appropriate regression model is of crucial importance, but in practice, it is not easy to find the perfect one. [27] proposes a supervised learning framework that uses a classifier to replace regression models. It presents a feature representation that employs the distance between the conditional distributions given past variable values and shows experimentally that the feature representation provides sufficiently different feature vectors for time series with different causal relationships. For the time series data, another issue that needs to be considered is hidden confounders. A time series deconfounder [15] was developed, which leverages the assignment of multiple treatments over time to enable the estimation of treatment effects even in the presence of hidden confounders. This time series deconfounder uses a recurrent neural network architecture with multitask output to build a factor model over time and infer substitute confounders, which render the assigned treatments conditionally independent. Then it performs causal inference using the substitute confounders.

For the second direction in SUTVA assumption, it assumes that there only exists one version for each treatment. However, if adding one continuous parameter into the treatment, this assumption cannot hold anymore. For example, estimating individual dose response curves for a couple of treatments requires adding an associated dosage parameter (categorical or continuous) for each treatment. Under this situation, for each treatment, it will have multiple versions for categorical dosage parameters or infinite versions for continuous dosage parameters. One way to solve this problem is to convert the continuous dosage into a categorical variable and then treat every medication with specific dosage as one new treatment, so that it will satisfy the SUTVA assumption again [119].

Another example that breaks the SUTVA is the dynamic treatment regime, which consists of a sequence of decision rules, one per stage of intervention [24]. One useful application of dynamic treatment is precision medication. It includes more individualization to adjust which type of treatment should be used, or how many the dosage is best in response to the patient's background characteristics, the illness severity and other heterogeneity, aiming to get the optimal treatment strategy. These heterogeneities are called tailoring variables. To get a useful dynamic treatment regime, [72] introduces one 'biased coin adaptive within-subject' (BCAWS) design. Then, [85] presents one general framework of this type of design, which uses sequential multiple assignment randomized trials (SMART) for developing decision rules in that each individual may be randomized multiple times and the multiple randomizations occur sequentially over time.

For estimating optimal dynamic decision rules from observational data, Q [144, 145] and A [84, 102] learning are two main approaches for estimating the optimal dynamic treatment regime. Q in Q-learning denotes "quality". Q-learning is a model-free reinforcement learning algorithm, which employs posited regression models for estimating outcome at each decision point given units' information. In advantage learning (A-learning), models are posited only for the part of the regression including contrasts among treatments and for the probability of observed treatment assignment at each decision point, given units' information. Both methods are implemented through a backward recursive fitting procedure that is related to dynamic programming [13].

## 4.2 Unconfoundedness assumption

The ignorability assumption is also named as unconfoundedness assumption. Given the background variable, $X$, the treatment assignment $W$ is independent to the potential outcomes, i.e., $W \perp Y(W = 0), Y(W = 1)|X$. With this unconfoundedness assumption, for the units with the same background variable $X$, their treatment assignment can be viewed as random. Obviously, identifying and collecting all of background variables is impossible, and this assumption is very difficult to satisfy. For example, in an observational study that tries to estimate the individual treatment effect of a medicine, instead of randomized experiments, the medicine is assigned to individuals based on a series of factors. Some factors (e.g., socioeconomic status) are challenging to measure and therefore become hidden confounders. Existing work overwhelmingly relies on the unconfoundedness assumption that all confounders can be measured. However, this assumption might be untenable in practice. In the above example, units' demographic attributes, such as their home address, consumption ability or employment status, may be the proxies for socioeconomic status. Leveraging big data, it is possible to find a proxy for the latent and unobserved confounders.

Variational autoencoder has been used to infer the complex non-linear relationships between the observed confounders and joint distribution of the latent confounders, treatment assignment and outcomes [81]. The joint distribution of the latent confounders and the observed confounders can be approximately recovered from the observations. An alternative way is to capture their patterns and control their influence by incorporating the underlying network information. Network information is also a reasonable proxy for the unobserved confounding. [46] applies GCN on network information to get the representation of hidden confounders. Moreover, in [45], graph attention layers are used to map the observed features in networked observational data to the D-dimensional space of partial latent confounders, by capturing the unknown edge weights in the real-world networked observational data.

An interesting insight mentioned in [138] is that, even if the confounders are observed, it doesn't mean all the information they contain is useful to infer the causal effect. Instead, requiring the part of confounders actually used by the estimator is sufficient. Therefore, if a good predictive model for the treatment can be built, one may only need to plug the outputs into a causal effect estimate directly, without any need to learn the whole true confounders. In [138], the main idea is to reduce the causal estimation problem to a semi-supervised prediction of both the treatments and outcomes. Networks admit high-quality embedding models that can be used for this semi-supervised prediction. In addition, embedding methods can also offer an alternative to fully specified generative models.

Only using observational data to solve the confoundings problem is always difficult. Another way is to combine the experimental data and observational data together. In [63], limited experimental data is used to correct the hidden confounding in causal effect models trained on larger observational data, even if the observational data do not fully overlap with the experimental data. This method makes strictly weaker assumptions than existing approaches.

For estimating treatment effects from longitudinal observational data, existing methods usually assume that there are no hidden confounders. This assumption is not testable in practice and, if it does not hold, leads to biased estimates. [15] infers substitute confounders that render the assigned treatments conditionally independent. Then it performs causal inference using the substitute confounders. This method can help estimate treatment effects for time series data in the presence of hidden confounders.

Above methods all aim to solve the problems about the observed and unobserved confounders. Are there any other ways to get around the unconfoundedness assumption and conduct causal inference? One way is to use instrumental variables that only affect treatment assignment but not the outcome variable. Changes in the instrumental variables would lead to the different assignment of treatment, which is independent of the latent variables, and this assignment is as good as randomization for the purposes of causal inference. [50] broke

instrumental variables analysis into two supervised stages that can each be targeted with deep networks. It models the conditional distribution of the treatment variable given the instruments and covariates, and then employs a loss function involving integration over the conditional treatment distribution. The deep instrumental variables framework also takes advantage of existing supervised learning techniques to estimate causal effects.

### 4.3  Positivity assumption

The positivity assumption, also known as covariate overlap or common support, is a necessary assumption for the identification of treatment effect in the observational study. However, little literature discusses the satisfaction of this assumption in the high dimensional datasets. [32] argues that the positivity assumption is a strong assumption and is more difficult to be satisfied in the high-dimensional datasets. To support the claim, the implication of the strict overlap assumption is explored, and it shows that strict overlap restricts the general discrepancies between the control and treated covariates. Therefore, the positivity assumption is stronger than the investigator expected. Based on the above implication, methods that eliminate the information about the treatment assignment while still hold the unconfounderness assumption are recommended, such as trimming [30, 97, 106] that drops the records in the region without overlap, and instrumental variable adjustment methods [35, 86, 93] that eliminate the instrumental variables from covariates.

## 5   GUIDELINE ABOUT EXPERIMENT

In this section, we provide the related experimental information, including the available datasets that are commonly adopted in the experiments, and the open-source codes of the methods mentioned in the previous two sections.

### 5.1  Available Datasets

*5.1.1  Datasets for Section 3.* Because the counterfactual outcome can never be observed, it's hard to find the dataset that perfectly satisfies the requirements of the experiment that it is an observational dataset with the ground truth ATE (or ITE) available. The datasets used in the literature are often semi-synthetic datasets. Some datasets, such as IHDP dataset, are obtained from the randomized dataset by generating their observed outcome according to a certain generation process and removing a biased subset to mimic the selection bias in the observational dataset. Some datasets, such as Jobs dataset, combine the randomized dataset and the observational control dataset together to create the selection bias. The details of the available benchmark datasets are in the following.

**IHDP**. This dataset is a commonly adopted benchmark dataset. This dataset is generated based on the randomized controlled experiment conducted by Infant Health and Development Program [22], whose targets are low-birthweight, premature infants. The pre-treatment covariates are 25 variables measuring the aspects about the children and their mothers, such as birth weight, head circumference, neonatal health index, prenatal care, mother's age, education, drugs, alcohol, etc. In the treated group, the infants are provided with both intensive high-quality childcare and specialist home visits [53]. The outcome is the infants' cognitive test score and can be simulated through the NPCI package[1]. Besides, a biased subset of the treated group is required to be removed to simulate the selection bias. An example of IHDP dataset whose outcome is simulated by the setting "A" of NPCI package can be downloaded from http://www.mit.edu/~fredrikj/files/ihdp_100.tar.gz.

**Jobs**. The jobs datasets used in the observational study [33, 34, 122] is the combination of Lalonde experiment data and the PSID comparison group. Both Lalonde and PSID datasets can be downloaded from NBER website[2]. The pre-treatment covariates are 8 variables such as age, education, ethnicity, as well as earnings in 1974 and 1975.

---

[1]https://github.com/vdorie/npci
[2]http://users.nber.org/~rdehejia/data/nswdata2.html

The people in the treated group take part in the job training while in the control group are not. The outcome is employment status.

**Twins**. Twins dataset is first introduced in [81] and is adopt by various observational studies [81, 149, 152]. Twins dataset is constructed on the data of twins birth in the USA between 1989-1991 [5][3]. In [81], the twins whose gender is the same and weight is less than 2000$g$ are selected into records. For each twin pair, there are in total 40 pre-treatment covariates measuring the pregnancy, twins birth and the parents, such as the number of gestation weeks before birth, the quality of care during pregnancy, pregnancy risk factors (Anemia, alcohol use, tobacco use, etc.), adequacy of care, residence and so on. The outcome is the one-year mortality. The treatment is being the heavier one in the twins, and the outcome is the one-year mortality. In twins dataset, both treated (the heavier one in the twin) and control (the lighter one in the twin) outcomes are observed. The treatment assignment usually defined by the users to simulate the selection bias. For example, in [81, 152], the selection bias is created by the following procedure: $W_i|\mathbf{X}_i \sim Bern(Sigmoid(\mathbf{w}'\mathbf{X}_i) + n)$, where $\mathbf{w} \sim U(-0.1, 0.1)^{40 \times 1}$ and $n \sim \mathbf{N}(0, 0.1)$.

**ACIC datasets**. Starting from 2016, every year, the Atlantic Causal Inference Conference holds the causal inference data analysis challenge, which provides some datasets targeting different causal inference problems.

2016 Challenge: The goal of ACIC 2016 Challenge is to better understand which approaches to causal inference perform well in particular observational study settings[4]. The datasets contain 77 datasets with varying degrees of non-linearity, sparsity, correlation between treatment assignment and outcome, non-linearity of treatment effect, overlapping. The covariates are real-world data from the Infant Health and Development Program dataset [22], which consists of 58 variables. The treatment, factual outcome, and counterfactual outcome are all generated by simulation, and the selection bias is created by removing treated children with non-white moms. The whole datasets can be downloaded from https://drive.google.com/file/d/0B7pG5PPgj6A3N09ibmFwNWE1djA/view, and the summarization of this year's challenge is in [36].

2017 Challenge: ACIC 2017 challenge focused on the estimation and inference for conditional average treatment effects (CATEs) in the presence of targeted selection. Targeted selection means the likelihood that an individual receives treatment is a function of the expected response of that individual if left untreated, which leads to strong confounding [47]. The same as the previous year's challenge, the covariates are from Infant Health and Development Program dataset [22], but only 8 variables are used. The outcomes and the treatment assignments are generated according to 32 distinct, fixed, data generating processes representing four different types of errors. For every data generating process, 250 independent replicate data sets were produced, and overall, there are a total of 8,000 data sets.

2018 Challenge: The ACIC 2018 challenge has two different tasks focusing on two sub-challenges: censoring and scaling. Censoring means some of the samples may not have observed outcomes. Therefore, the dataset used by censoring challenges contains missing outcome values for some of the samples. The dataset for scaling challenge contains 48 datasets whose data sizes, and they are not censored. The details of the above datasets are available at https://www.synapse.org/#!Synapse:syn11294478/wiki/486304

2019 Challenge: This challenge focuses on estimating the ATE on the quasi real-world dataset with low dimensional data and high dimensional data[5]. The datasets for this challenge contain several datasets with different variables size and record size, and the R code for data generation is available at https://drive.google.com/file/d/1Qqgmb3R9Vt9KTx6t8i_5IbFenylsPfrK/view.

**IBM causal inference benchmark**. This dataset is created in [126] and is available at https://github.com/IBM-HRL-MLHLS/IBM-Causal-Inference-Benchmarking-Framework. This dataset uses the cohort of 100K samples in

---

[3]www.nber.org/data/linked-birth-infant-death-data-vital-statistics-data.html
[4]https://jenniferhill7.wixsite.com/acic-2016/competition
[5]https://sites.google.com/view/acic2019datachallenge/data-challenge?authuser=0

Linked Births and Infant Deaths Database (LBIDD)[6] as the fundamental set of covariates. The treatment, factual outcome, and the counterfactual outcome are generated by simulation.

**BlogCatalog**. This dataset is used for causal inference with networked observational data [46]. It is a social blogger network. A blogger is one observation. The bloggers are connected by some social relationships in this dataset. The features are bag-of-words representations of keywords in bloggers' descriptions. The outcomes are the opinions of readers on each blogger. A blogger belongs to the treated group (control group) if her blogs are read more on mobile devices (desktops).

**Flickr**. This dataset includes networked observational data in [46]. Flickr is a photo-sharing platform and social network where users upload photos for others to see. In this dataset, the users with Flickr account are observations, and the users are connected by some social relationships. The features of each user are the tags of interest. The outcomes and treatment assignment are the same as BlogCatalog.

**News**. The News benchmark includes 5000 randomly sampled news articles from the NY Times corpus. It contains the data on the opinion of media consumers on news items and was originally introduced as a benchmark for counterfactual inference in the setting with two treatment options [60]. It can be extended to multiple treatments with associated dosage parameters [119]. The details can be found in https://archive.ics.uci.edu/ml/datasets/bag+of+words.

**MVICU**. The Mechanical Ventilation in the Intensive Care Unit (MVICU) benchmark is used to estimate individual dose-response curves for a couple of treatments with an associated dosage parameter [119]. This dataset includes patients' responses to different configurations of mechanical ventilation in the intensive care unit. The data was sourced from the publicly available MIMIC III database which documents a diverse and very large population of ICU patient stays and contains comprehensive and detailed clinical data. [114].

**TCGA**. The Cancer Genome Atlas (TCGA) is the world's largest and richest collection of genomic data. This dataset is used to estimate individual dose-response curves for a couple of treatments with an associated dosage parameter [119]. The TCGA project collected gene expression data from various types of cancers in 9659 individuals [146]. The treatment options are medication, chemotherapy and surgery. The outcome is the risk of cancer recurrence after receiving the treatment. TCGA data (controlled access and open access data) are available via the Genomic Data Commons (GDC) https://gdc.cancer.gov/.

**Saccharomyces cerevisiae (yeast) cell cycle gene expression dataset**. This is one time series dataset. A time series with the length T = 57 was created by combining four short time series that were measured in different microarray experiments [27].

**THE**. The Tennessee Student/Teacher Achievement Ratio (STAR) experiment is a four-year longitudinal class-size study funded by the Tennessee General Assembly and conducted by the State Department of Education to measure the influence of class size (small class, regular class and regular-with-aide class) on student achievement tests and non-achievement measures [2]. Because this is one randomized controlled experiment, CATE estimates are unbiased due to unconfoundedness. Confounders are artificially introduced by selectively removing a biased subset of samples [63].

**FERTIL2**. This dataset aims to study the impact of more than or exactly 7 years of education for a woman on the number of children in the family [147]. Several observed confounders are included in the dataset, such as age, whether the family has a TV, whether the woman lives in the city. The instrumental variable is a binary indicator of whether the woman was born in the first half of the year. This dataset is used for research about instrumental variables [35].

---

[6]https://www.cdc.gov/nchs/nvss/linked-birth.htm

## 5.2 Codes/Packages

In this part, we summarize the available codes or tool-boxes for causal inference. The codes for methods that mentioned in Section 3 are provided in Table 2 and Table 3, where Table 2 lists the tool-boxes with their supported methods and languages, and Table 3 lists the open-source code of one specific method.

Table 2. Available Tool-boxes for Causal Inference

| Tool-box | Supporting methods | Language | Link |
|---|---|---|---|
| Dowhy [124] | Propensity-based Stratification, PSM, IPW, Regression | Python | https://github.com/microsoft/dowhy |
| Causal ML | Tree-based algorithms, X/T/X/R-learner | Python | https://github.com/uber/causalml |
| EconML [100] | Doubly Robust Learner, Orthogonal Random Forests, Meta-Learners, Deep Instrumental Variables | Python | https://github.com/microsoft/EconML#blogs-and-publications |
| causalToolbox | BART, Causal Forest, T/X/S-learner with BART/RF as base learner | R | https://github.com/soerenkuenzel/causalToolbox |

Table 3. Available Codes of Methods in Section 3

| Method | Language | Link |
|---|---|---|
| IPW | R | https://cran.r-project.org/web/packages/ipw/index.html |
| DR | R | fastDR: https://github.com/gregridgeway/fastDR<br>DR for High dimension: https://github.com/gregridgeway/fastDR |
| Principal Stratification | R | https://cran.r-project.org/web/packages/sensitivityPStrat/index.html |
| Stratification | R | https://cran.r-project.org/web/packages/stratification/ |
| PSM<br>overlap weight<br>trapezoidal weight | Python | https://cran.r-project.org/web/packages/PSW/ |
| Matching based Alg.:<br>exact matching,<br>full matching,<br>genetic matching,<br>nearest neighbor matching,<br>optimal matching,<br>subclassification | R | https://cran.r-project.org/web/packages/MatchIt/ |
| PSM | Python | https://github.com/akelleh/causality |
| | | Continued on next page |

**Table 3 – continued from previous page**

| Method | Language | Link |
|---|---|---|
| Perfect Match | Python | https://github.com/d909b/perfect_match |
| optimal matching | R | https://cran.r-project.org/web/packages/Matching/ |
| CEM | R | https://cran.r-project.org/web/packages/cem/ |
| TMLE | R | https://cran.r-project.org/web/packages/tmle/index.html |
| CMGP [3] | Python | https://bitbucket.org/mvdschaar/mlforhealthlabpub/src/ baa0aa33a6af3fe490484c9e11e3a158968ae56a/ alg/causal_multitask_gaussian_processes_ite/ |
| BART | R Python | https://cran.r-project.org/web/packages/BayesTree/index.html https://github.com/JakeColtman/bartpy |
| GANITE [152] | Python | https://bitbucket.org/mvdschaar/mlforhealthlabpub/src/ baa0aa33a6af3fe490484c9e11e3a158968ae56a/alg/ganite/ |
| BNN [60], CFR-MMD [122], CFR-WASS [122] | Python | https://github.com/clinicalml/cfrnet |
| CEVAE | Python | https://github.com/AMLab-Amsterdam/CEVAE |
| SITE [149] | Python | https://github.com/Osier-Yi/SITE |
| grf | R | https://cran.r-project.org/web/packages/grf/index.html |
| R-learner | R | https://github.com/xnie/rlearner/blob/master/R/xlearner.R |
| Residual Balancing | R | https://github.com/swager/balanceHD |
| CBPS | R | https://github.com/kosukeimai/CBPS |
| dragonnet | Python | github.com/claudiashi57/dragonnet |
| Entropy Balancing | R | https://cran.r-project.org/web/packages/ebal/ |
| DRNets [119] | Python | https://github.com/d909b/drnet |
| Network Deconfounder [46] | Python | https://github.com/rguo12/network-deconfounder-wsdm20 |
| Network Embeddings [138] | Python | https://github.com/vveitch/causal-network-embeddings |
| RMSN [79] | Python | https://github.com/vveitch/causal-network-embeddings |
| TMLE [96] | R | https://github.com/joshuaschwab/ltmle |
| LCVA [99] | Python | https://github.com/rguo12/CIKM18-LCVA |

## 6 APPLICATIONS

Causal inference has a variety of applications in real-world scenarios. In general, the applications of causal inference can be categorized into three directions:

(1) Decision evaluation. This is a natural application of treatment effect estimation as it is consistent with the objective of treatment effect estimation.

(2) Counterfactual estimation. Counterfactual learning greatly helps the areas related to decision making, as it can provide the potential outcomes of different decision choices (or policies).
(3) Dealing with selection bias. In many real-world applications, the records appear in the collected dataset are not representative of the whole population that is interested. Without appropriately handling the selection bias, the generalization of the trained model would be hurt.

In this section, we will discuss how causal inference benefits various real-world applications in detail.

## 6.1 Advertising

Properly measuring the effect of an advertising campaign can answer critical marketing questions such as whether a new advertisement increases the clicks, or whether a new campaign increases sales, etc. Since conducting randomized experiments is expensive and time-consuming, estimating the advertisement effect from the observational data is attracting increasing attention in both industry and research communities [132, 142]. In [78], the randomized nearest neighbor matching method is proposed to estimate the treatment effect of digital marketing campaigns. In [39], the covariate balancing generalized propensity score (CBGPS), discussed in Section 3.1.1, is applied to analyze the efficacy of political advertisements.

However, in the online advertisement area, it is often required to deal with complex advertisement treatments, which could be a discrete or continuous, uni- or multi-dimensional treatment [132]. One advertisement is set as the baseline treatment, and the treatment effect is obtained by comparing the potential outcome of the treatment with different values and the baseline treatment. To estimate the potential outcome of treatment with multi-dimensional values, tree-based method [142] and sparse additive model based method [132] are proposed to enable the comparison between potential treatments and the baseline treatment.

In addition to purely observational data, in the real-world scenarios, it is often the case that dataset is comprised of large samples from control condition(i.e., the old treatment) and small samples (possibility unrepresentative) from a randomized trial which contains both the control condition and the new treatment. In [109], the small randomized trial dataset is connected with the large control dataset using the minimal set of modeling assumptions, which implies the models to predict the control and treated outcome to be similar. Under this assumption, the proposed method jointly learns the control and treated outcome predictor and regularizes the difference between the parameters of two predictors.

The above discussions show the potential applications of the treatment effect estimation in decision evaluation: measuring the effect of the advertisement campaign. Another important application is to handle the selection bias. Due to the existing selection mechanism in the advertising systems, there is a distribution discrepancy between the displayed and non-displayed events [153]. Ignoring such bias would make the advertisement click prediction inaccurate, which would cause a loss of revenue. To handle the selection bias, similar to the doubly robust estimation mentioned in Section 3.1.1, doubly robust policy learning is proposed in [37]. It contains two sub-estimators: direct method estimator obtained from the observed samples, and IPS estimator with the propensity score as the sample weight.

Furthermore, some works notice the difficulty of propensity score estimation due to the deterministic advertisement display policy in commercial advertisement systems. If the display policy is stochastic, the advertisements with low propensity scores still have a chance to appear in the observational dataset so that IPS can correct the selection bias. However, when the display policy is deterministic, the advertisements with low propensity scores are always absent in the observation, which makes propensity score estimation failed. This challenges motives the work of propensity-free doubly robust method proposed in [153] which improves the original doubly robust method in two folds: (1) Train the direct method on a small but unbiased data obtained under the uniform policy, which, to a certain degree, prevents the selection bias propagating to the non-displayed advertisements. (2) Avoid the propensity score estimation by setting the propensity score of the observed items as 1 and combines IPS with

the direct method. In a nutshell, this propensity score free method relies on the direct method trained on a small unbiased dataset to give an unbiased prediction of the advertisement click.

Apart from the applications discussed above, another important application is the advertisement recommendation, which is merged into the next subsection.

## 6.2   Recommendation

The recommendation is highly correlated with the treatment effect estimation, as exposing the user to an item in the recommendation system can be viewed as applying one specific treatment to a unit [71, 118]. Similar to the dataset used in the treatment effect estimation, the dataset used in the recommendation are usually biased due to the self-selection of the users. For example, in the movie rate dataset, users tend to rate the movies that they like: the horrible movie ratings are mostly made by horror movie fans and less by romantics movie fans. Another example is the advertisement recommendation datasets. The recommendation system would only recommend the advertisements to the users whom the system believes are interested in those advertisements. In the above examples, the records in the datasets are not representative of the whole population, which is the selection bias. The selection bias brings challenges to both recommendation model training and evaluation. Re-weighting samples based on the propensity score is a powerful method to solve the problems that stem from selection bias. The improved performance estimation after propensity score weighting can be calculated as follows:

$$\hat{R}_{\text{IPS}}(\hat{Y}|P) = \frac{1}{U \cdot I} \sum_{(u,i):O_{u,i}=1} \frac{\delta_{u,i}(Y,\hat{Y})}{P_{u,i}}, \tag{32}$$

where $\hat{Y}$ is the value upon which to measure the quality of a recommendation system, $U$ is the number of users, and $I$ is the number of items. $O_{u,i}$ is a binary variable to indicates the interaction of the $u$-th user with the $i$-th item in the observational data. $\delta_{u,i}(\cdot,\cdot)$ can be any classical quality measure of a recommendation, such as cumulative gain(CG), discounted cumulative gain (DCG), and precision at $k$. $P$ is the marginal probability matrix, whose entry is defined as $P_{u,i} = P(O_{u,i} = 1)$. The improved quality measure is an unbiased estimation to the real measurement $R(\hat{Y})$ over the whole population, which is defined as $R(\hat{Y}) = \frac{1}{U \cdot I} \sum_{u=1}^{U} \sum_{i=1}^{I} \delta_{u,i}(Y,\hat{Y})$. Based on the unbiased quality measurement, in [118], the propensity-score empirical risk minimization (ERM) for recommendation is proposed: $\hat{Y} \in \mathcal{H}$ is selected to optimize the following problem: $\hat{Y}^{ERM} = \text{argmin}_{\hat{Y} \in \mathcal{H}} \{\hat{R}_{\text{IPS}}(\hat{Y}|P)\}$, where $\hat{R}_{\text{IPS}}(\hat{Y}|P)$ is defined in Eqn. (32). Later, various works are developed to drawbacks of propensity score weighting, including estimation variance reduction [115, 135], handle data sparsity [115, 135], doubly robust estimation [143].

In addition to using IPS or doubly robust estimation based method to overcome the selection bias, similar to the advertisement domain, some works also adopt a small unbiased dataset to correct selection bias. In this case, the dataset contains a large set of logged feedback records under control policy and a small set of records under the randomized recommendation. CausalEmbed(CausE) [19] is a representative method in this direction, which proposed a new matrix factorization algorithm. In detail, CausalEmbed jointly factorizes the matrix of those two datasets and links these two models by regularizing the difference between treated and control representation.

## 6.3   Medicine

Learning the optimal per-patient treatment rules is one of the promising goals of applying treatment effect estimation methods in the medical domain. When the effect of different available medicines can be estimated, doctors can give a better prescription accordingly. In [121], two challenges are mentioned to fulfill such goal: the existence of confounders and the existence of unobserved confounders. Although analyzing from the randomized experimental dataset is the golden solution, it has the following limitations: (1) The goal of randomized experimental data is to analyze the ATE instead of ITE, so the data size is often small, which limits the capacity to derive personalized treatment rules. (2) As mentioned in section 2, conducting randomized trials is often expensive,

time-consuming, and sometimes maybe unethical. Therefore, deriving personalized treatment rules from the observational dataset or the combination of experimental and observational data are two fruitful directions [121].

For the direction of utilizing observational dataset, various methods derive the personalized treatment rules guided by the estimated ITE under the unconfoundedness assumption, such as deep-treat [7], tiered case-cohort design based method [66]. However, in this area, there are limited works to handle the unobserved confounders, and methods discussed in Section 4.2 have great potentials to explore.

## 6.4 Reinforcement Learning

From the perspective of reinforcement learning, ITE estimation can be viewed as a contextual multi-armed bandit problem with the treatment as the action, the outcome as the reward, and the background variables as the contextual information. Arm exploration and exploitation is similar to randomized trials and observational data. Therefore, these two areas share some similar critical challenges: (1) How to get an unbiased outcome/reward estimation? (2) How to handle either the observed or unobserved confounders that affect both the treatment assignment/action choice and the outcome/reward?

To obtain an unbiased reward estimation, importance sampling weighting [98] is the common method adopted in the offline policy evaluation. The weight is set as the probability between the target policy and the logged (observed) policy, which is analogous to IPW mentioned in Section 3.1.1. However, importance sampling proposed in [98] suffers from high variance and highly relies on the assigned weights. To improve this, similar to the doubly robust method in ATE estimation, doubly robust policy evaluation is proposed in [37]. Later, various methods [8, 16, 64, 76, 134, 135, 137, 155] are proposed to improve those two methods with different settings.

As mentioned above, the second challenge is how to deal with confounders. When all the confounders are observed, we can directly optimize the unbiased reward function mentioned in the previous paragraph. However, when there exist unobserved confounders, it can lead to policies that introduce harm rather than benefit, as is generally the case with observational data [65]. The confounding-robust policy learning framework is proposed in [65], optimizing the policy over an uncertain set for propensity weights so that the unobserved confounders can be controlled.

## 6.5 Other Applications

The applications of causal inference are not limited to the areas mentioned above, and areas related to effectiveness measurement, decision making, or handling selection bias, are all potential applications.

**Education**. In the education area, by comparing the outcome of different teaching methods on the student population, a better teaching method can be decided. Moreover, ITE estimation can enhance personalized learning by estimating the outcome of each student on different teaching methods. For example, ITE estimation is developed to answer the questions "Would this particular student benefit more from the video hint or the text hint when this student cannot solve a problem?", so that an Intelligent Tutor System (ITS) can decide which hint is more suitable for a specific student [154].

**Political decision**. In the politics area, causal inference can provide decision support. For example, various methods [60, 120, 122, 149, 152] have been developed on the jobs dataset aiming to answer the question "who would benefit most from subsidized job training?". Causal inference can also help with political decisions such as whether a policy should be promoted to large population size.

**Improving Machine learning methods**. In addition to the decision support, various balancing methods that can handle the selection bias ( mentioned in Section 3), can also be extended to improve the stability of machine learning methods. In [67], the reweighting method is adopted to improve the generalization ability of learned models for unknown environments (i.e., unknown test data). To be specific, the weight of each sample is added to the prediction loss function as a regularization, which is formulated as: $\sum_{j=1}^{p} || \frac{\phi(X_{.,-j}^T)(R \odot X_{.,j})}{R^T X_{.,j}} - \frac{\phi(X_{.,-j}^T)(R \odot (1 - X_{.,j}))}{R^T(1 - X_{.,j})} ||_2^2$,

where $p$ is the number of total features, $\phi(\cdot)$ is the feature transformation function such as neural network, $X_{\cdot,j}$ is the $j$-th feature in $X$, $X_{\cdot,-j}$ is the features in $X$ except the $j$-th feature, $R \in \mathcal{R}^N$ is the global sample weights with N as the number of total samples. This balancing regularizer extends the CBPS method discussed in 3.1.1, by taking the $j$-th feature as the treatment and the remaining features as the background variables, and then combining all the features to obtain the global balancing weight.

## 7 CONCLUSIONS

Causal inference has been an attractive research topic for a long time as it provides an effective way to uncover causal relationships in real-world problems. Nowadays, the flourishing of machine learning brings new vitality into this area, and meanwhile, the incisive ideas in the causal inference area promote the development of machine learning. In this survey, we provide a comprehensive review of the methods under the well-known potential outcome framework. As the potential outcome framework relies on the three assumptions, the methods are separated into two categories. One category relies on those assumptions, while the other one relaxes some of the assumptions. For each category, we provide thorough discussions, comparisons, and summarization of the reviewed methods. The available benchmark datasets and open-source codes of those methods are also listed. Finally, some representative real-world applications of causal inference are introduced, such as advertising, recommendation, medicine, and reinforcement learning.

## REFERENCES

[1] A. Abadie, D. Drukker, J. L. Herr, and G. W. Imbens. Implementing matching estimators for average treatment effects in stata. *The stata journal*, 4(3):290–311, 2004.
[2] C. Achilles, H. P. Bain, F. Bellott, J. Boyd-Zaharias, J. Finn, J. Folger, J. Johnston, and E. Word. Tennessee's student teacher achievement ratio (star) project'. *URL: http://hdl. handle. net/1902.1/10766*, 2008.
[3] A. M. Alaa and M. van der Schaar. Bayesian inference of individualized treatment effects using multi-task gaussian processes. In *Advances in Neural Information Processing Systems*, pages 3424–3432, 2017.
[4] A. M. Alaa, M. Weisz, and M. Van Der Schaar. Deep counterfactual networks with propensity-dropout. *arXiv preprint arXiv:1706.05966*, 2017.
[5] D. Almond, K. Y. Chay, and D. S. Lee. The costs of low birth weight. *The Quarterly Journal of Economics*, 120(3):1031–1083, 2005.
[6] N. Altman and M. Krzywinski. Points of significance: Association, correlation and causation. *Nature Methods*, 12(10):899–900, 2015.
[7] O. Atan, J. Jordon, and M. van der Schaar. Deep-treat: Learning optimal personalized treatments from observational data using neural networks. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
[8] O. Atan, W. R. Zame, and M. van der Schaar. Learning optimal policies from observational data. *arXiv preprint arXiv:1802.08679*, 2018.
[9] S. Athey and G. Imbens. Recursive partitioning for heterogeneous causal effects. *Proceedings of the National Academy of Sciences*, 113(27):7353–7360, 2016.
[10] S. Athey and G. W. Imbens. Machine learning methods for estimating heterogeneous causal effects. *stat*, 1050(5):1–26, 2015.
[11] P. C. Austin. An introduction to propensity score methods for reducing the effects of confounding in observational studies. *Multivariate behavioral research*, 46(3):399–424, 2011.
[12] H. Bang and J. M. Robins. Doubly robust estimation in missing data and causal inference models. *Biometrics*, 61(4):962–973, 2005.
[13] J. Bather. *Decision theory: An introduction to dynamic programming and sequential decisions*. John Wiley & Sons, Inc., 2000.
[14] S. Ben-David, J. Blitzer, K. Crammer, and F. Pereira. Analysis of representations for domain adaptation. In *Advances in neural information processing systems*, pages 137–144, 2007.
[15] I. Bica, A. M. Alaa, and M. van der Schaar. Time series deconfounder: Estimating treatment effects over time in the presence of hidden confounders. *arXiv preprint arXiv:1902.00450*, 2019.
[16] A. Bietti, A. Agarwal, and J. Langford. Practical evaluation and optimization of contextual bandit algorithms. *stat*, 1050:12, 2018.
[17] A. Bloniarz, H. Liu, C.-H. Zhang, J. S. Sekhon, and B. Yu. Lasso adjustments of treatment effect estimates in randomized experiments. *Proceedings of the National Academy of Sciences*, 113(27):7383–7390, 2016.
[18] C. R. Blyth. On simpson's paradox and the sure-thing principle. *Journal of the American Statistical Association*, 67(338):364–366, 1972.
[19] S. Bonner and F. Vasile. Causal embeddings for recommendation. In *Proceedings of the 12th ACM Conference on Recommender Systems*, pages 104–112, 2018.
[20] L. Breiman. Random forests. *Machine learning*, 45(1):5–32, 2001.
[21] L. Breiman. *Classification and regression trees*. Routledge, 2017.

[22] J. Brooks-Gunn, F.-r. Liaw, and P. K. Klebanov. Effects of early intervention on cognitive function of low birth weight preterm infants. *The Journal of pediatrics*, 120(3):350–359, 1992.

[23] M. Caliendo and S. Kopeinig. Some practical guidance for the implementation of propensity score matching. *Journal of economic surveys*, 22(1):31–72, 2008.

[24] B. Chakraborty. *Statistical methods for dynamic treatment regimes*. Springer, 2013.

[25] B. Chakraborty and S. A. Murphy. Dynamic treatment regimes. *Annual review of statistics and its application*, 1:447–464, 2014.

[26] Y. Chang and J. G. Dy. Informative subspace learning for counterfactual inference. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.

[27] Y. Chikahara and A. Fujino. Causal inference in time series via supervised learning. In *IJCAI*, pages 2042–2048, 2018.

[28] H. A. Chipman, E. I. George, and R. E. McCulloch. Bayesian ensemble learning. In *Advances in neural information processing systems*, pages 265–272, 2007.

[29] H. A. Chipman, E. I. George, and R. E. McCulloch. Bart: Bayesian additive regression trees. *The Annals of Applied Statistics*, 4(1):266–298, 2010.

[30] R. K. Crump, V. J. Hotz, G. W. Imbens, and O. A. Mitnik. Dealing with limited overlap in estimation of average treatment effects. *Biometrika*, 96(1):187–199, 2009.

[31] R. B. D'Agostino Jr. Propensity score methods for bias reduction in the comparison of a treatment to a non-randomized control group. *Statistics in medicine*, 17(19):2265–2281, 1998.

[32] A. D'Amour, P. Ding, A. Feller, L. Lei, and J. Sekhon. Overlap in observational studies with high-dimensional covariates. *arXiv preprint arXiv:1711.02582*, 2017.

[33] R. H. Dehejia and S. Wahba. Causal effects in nonexperimental studies: Reevaluating the evaluation of training programs. *Journal of the American statistical Association*, 94(448):1053–1062, 1999.

[34] R. H. Dehejia and S. Wahba. Propensity score-matching methods for nonexperimental causal studies. *Review of Economics and statistics*, 84(1):151–161, 2002.

[35] P. Ding, T. VanderWeele, and J. Robins. Instrumental variables as bias amplifiers with general outcome and confounding. *Biometrika*, 104(2):291–302, 2017.

[36] V. Dorie, J. Hill, U. Shalit, M. Scott, D. Cervone, et al. Automated versus do-it-yourself methods for causal inference: Lessons learned from a data analysis competition. *Statistical Science*, 34(1):43–68, 2019.

[37] M. Dudík, J. Langford, and L. Li. Doubly robust policy evaluation and learning. *arXiv preprint arXiv:1103.4601*, 2011.

[38] J. Fan, K. Imai, H. Liu, Y. Ning, and X. Yang. Improving covariate balancing propensity score: A doubly robust and efficient approach. Technical report, Technical report, Princeton Univ, 2016.

[39] C. Fong, C. Hazlett, K. Imai, et al. Covariate balancing propensity score for a continuous treatment: application to the efficacy of political advertisements. *The Annals of Applied Statistics*, 12(1):156–177, 2018.

[40] C. E. Frangakis and D. B. Rubin. Principal stratification in causal inference. *Biometrics*, 58(1):21–29, 2002.

[41] S. Glazerman, D. M. Levy, and D. Myers. Nonexperimental versus experimental estimates of earnings impacts. *The Annals of the American Academy of Political and Social Science*, 589(1):63–93, 2003.

[42] I. Good, Y. Mittal, et al. The amalgamation and geometry of two-by-two contingency tables. *The Annals of Statistics*, 15(2):694–711, 1987.

[43] X. S. Gu and P. R. Rosenbaum. Comparison of multivariate matching methods: Structures, distances, and algorithms. *Journal of Computational and Graphical Statistics*, 2(4):405–420, 1993.

[44] R. Guo, L. Cheng, J. Li, P. R. Hahn, and H. Liu. A survey of learning causality with data: Problems and methods. *arXiv preprint arXiv:1809.09337*, 2018.

[45] R. Guo, J. Li, and H. Liu. Counterfactual evaluation of treatment assignment functions with networked observational data. *arXiv preprint arXiv:1912.10536*, 2019.

[46] R. Guo, J. Li, and H. Liu. Learning individual treatment effects from networked observational data. *arXiv preprint arXiv:1906.03485*, 2019.

[47] P. R. Hahn, V. Dorie, and J. S. Murray. Atlantic causal inference conference (acic) data analysis challenge 2017. *arXiv preprint arXiv:1905.09515*, 2019.

[48] P. R. Hahn, J. S. Murray, and C. Carvalho. Bayesian regression tree models for causal inference: regularization, confounding, and heterogeneous effects. *arXiv preprint arXiv:1706.09523*, 2017.

[49] B. B. Hansen. The prognostic analogue of the propensity score. *Biometrika*, 95(2):481–488, 2008.

[50] J. Hartford, G. Lewis, K. Leyton-Brown, and M. Taddy. Deep iv: A flexible approach for counterfactual prediction. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1414–1423, 2017.

[51] N. Hassanpour and R. Greiner. Counterfactual regression with importance sampling weights. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pages 5880–5887, 2019.

[52] J. J. Heckman, H. Ichimura, and P. Todd. Matching as an econometric evaluation estimator. *The review of economic studies*, 65(2):261–294, 1998.

[53] J. L. Hill. Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, 20(1):217–240, 2011.

[54] M. G. Hudgens and M. E. Halloran. Toward causal inference with interference. *Journal of the American Statistical Association*, 103(482):832–842, 2008.

[55] K. H. Hullsiek and T. A. Louis. Propensity score modeling strategies for the causal analysis of observational data. *Biostatistics*, 3(2):179–193, 2002.

[56] S. M. Iacus, G. King, and G. Porro. Causal inference without balance checking: Coarsened exact matching. *Political analysis*, 20(1):1–24, 2012.

[57] K. Imai and M. Ratkovic. Covariate balancing propensity score. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 76(1):243–263, 2014.

[58] G. W. Imbens. Nonparametric estimation of average treatment effects under exogeneity: A review. *Review of Economics and statistics*, 86(1):4–29, 2004.

[59] G. W. Imbens and D. B. Rubin. *Causal inference in statistics, social, and biomedical sciences.* Cambridge University Press, 2015.

[60] F. Johansson, U. Shalit, and D. Sontag. Learning representations for counterfactual inference. In *International conference on machine learning*, pages 3020–3029, 2016.

[61] F. D. Johansson, N. Kallus, U. Shalit, and D. Sontag. Learning weighted representations for generalization across designs. *arXiv preprint arXiv:1802.08598*, 2018.

[62] Judea Pearl. Judea pearl on potential outcomes. http://causality.cs.ucla.edu/blog/index.php/2012/12/03/judea-pearl-on-potential-outcomes/, 2012.

[63] N. Kallus, A. M. Puli, and U. Shalit. Removing hidden confounding by experimental grounding. In *Advances in Neural Information Processing Systems*, pages 10888–10897, 2018.

[64] N. Kallus and M. Uehara. Intrinsically efficient, stable, and bounded off-policy evaluation for reinforcement learning. *arXiv preprint arXiv:1906.03735*, 2019.

[65] N. Kallus and A. Zhou. Confounding-robust policy improvement. In *Advances in Neural Information Processing Systems*, pages 9269–9279, 2018.

[66] R. C. Kessler, R. M. Bossarte, A. Luedtke, A. M. Zaslavsky, and J. R. Zubizarreta. Machine learning methods for developing precision treatment rules with observational data. *Behaviour research and therapy*, 120:103412, 2019.

[67] K. Kuang, P. Cui, S. Athey, R. Xiong, and B. Li. Stable prediction across unknown environments. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 1617–1626, 2018.

[68] K. Kuang, P. Cui, B. Li, M. Jiang, and S. Yang. Estimating treatment effect in the wild via differentiated confounder balancing. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 265–274, 2017.

[69] K. Kuang, P. Cui, B. Li, M. Jiang, S. Yang, and F. Wang. Treatment effect estimation with data-driven variable decomposition. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.

[70] S. R. Künzel, J. S. Sekhon, P. J. Bickel, and B. Yu. Metalearners for estimating heterogeneous treatment effects using machine learning. *Proceedings of the National Academy of Sciences*, 116(10):4156–4165, 2019.

[71] A. Lada, A. Peysakhovich, D. Aparicio, and M. Bailey. Observational data for heterogeneous treatment effects with application to recommender systems. In *Proceedings of the 2019 ACM Conference on Economics and Computation*, pages 199–213, 2019.

[72] P. W. Lavori and R. Dawson. A design for testing clinical strategies: biased adaptive within-subject randomization. *Journal of the Royal Statistical Society: Series A (Statistics in Society)*, 163(1):29–38, 2000.

[73] B. K. Lee, J. Lessler, and E. A. Stuart. Weight trimming and propensity score weighting. *PloS one*, 6(3):e18174, 2011.

[74] C. Lee, N. Mastronarde, and M. van der Schaar. Estimation of individual treatment effect in latent confounder models via adversarial learning. *arXiv preprint arXiv:1811.08943*, 2018.

[75] F. Li, K. L. Morgan, and A. M. Zaslavsky. Balancing covariates via propensity score weighting. *Journal of the American Statistical Association*, 113(521):390–400, 2018.

[76] L. Li, W. Chu, J. Langford, T. Moon, and X. Wang. An unbiased offline evaluation of contextual bandit algorithms with generalized linear models. In *Proceedings of the Workshop on On-line Trading of Exploration and Exploitation 2*, pages 19–36, 2012.

[77] S. Li and Y. Fu. Matching on balanced nonlinear representations for treatment effects estimation. In *Advances in Neural Information Processing Systems*, pages 929–939, 2017.

[78] S. Li, N. Vlassis, J. Kawale, and Y. Fu. Matching via dimensionality reduction for estimation of treatment effects in digital marketing campaigns. In *Proceedings of the Twenty-Fifth International Joint Conference on Artificial Intelligence*, pages 3768–3774, 2016.

[79] B. Lim. Forecasting treatment responses over time using recurrent marginal structural networks. In *Advances in Neural Information Processing Systems*, pages 7483–7493, 2018.

[80] W.-Y. Loh. Classification and regression trees. *Wiley Interdisciplinary Reviews: Data Mining and Knowledge Discovery*, 1(1):14–23, 2011.

[81] C. Louizos, U. Shalit, J. M. Mooij, D. Sontag, R. Zemel, and M. Welling. Causal effect inference with deep latent-variable models. In *Advances in Neural Information Processing Systems*, pages 6446–6456, 2017.

[82] X. Ma and J. Wang. Robust inference using inverse probability weighting. *Journal of the American Statistical Association*, pages 1–10, 2010.

[83] Y. Mansour, M. Mohri, and A. Rostamizadeh. Domain adaptation: Learning bounds and algorithms. *arXiv preprint arXiv:0902.3430*, 2009.

[84] S. A. Murphy. Optimal dynamic treatment regimes. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 65(2):331–355, 2003.

[85] S. A. Murphy. An experimental design for the development of adaptive treatment strategies. *Statistics in medicine*, 24(10):1455–1481, 2005.

[86] J. A. Myers, J. A. Rassen, J. J. Gagne, K. F. Huybrechts, S. Schneeweiss, K. J. Rothman, M. M. Joffe, and R. J. Glynn. Effects of adjusting for instrumental variables on bias and precision of effect estimates. *American journal of epidemiology*, 174(11):1213–1222, 2011.

[87] X. Nie and S. Wager. Quasi-oracle estimation of heterogeneous treatment effects. *arXiv preprint arXiv:1712.04912*, 2017.

[88] E. L. Ogburn, T. J. VanderWeele, et al. Causal diagrams for interference. *Statistical science*, 29(4):559–578, 2014.

[89] J. Pearl. Causal diagrams for empirical research. *Biometrika*, 82(4):669–688, 1995.

[90] J. Pearl. *Causality: models, reasoning and inference*, volume 29. Springer, 2000.

[91] J. Pearl. Causal inference in statistics: An overview. *Statistics surveys*, 3:96–146, 2009.

[92] J. Pearl. *Causality*. Cambridge university press, 2009.

[93] J. Pearl. On a class of bias-amplifying variables that endanger effect estimates. *arXiv preprint arXiv:1203.3503*, 2012.

[94] J. Pearl. *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Elsevier, 2014.

[95] J. M. Peña. Reasoning with alternative acyclic directed mixed graphs. *Behaviormetrika*, 45(2):389–422, 2018.

[96] M. Petersen, J. Schwab, S. Gruber, N. Blaser, M. Schomaker, and M. van der Laan. Targeted maximum likelihood estimation for dynamic and static longitudinal marginal structural working models. *Journal of causal inference*, 2(2):147–185, 2014.

[97] M. L. Petersen, K. E. Porter, S. Gruber, Y. Wang, and M. J. van der Laan. Diagnosing and responding to violations in the positivity assumption. *Statistical methods in medical research*, 21(1):31–54, 2012.

[98] D. Precup. Eligibility traces for off-policy policy evaluation. *Computer Science Department Faculty Publication Series*, page 80, 2000.

[99] V. Rakesh, R. Guo, R. Moraffah, N. Agarwal, and H. Liu. Linked causal variational autoencoder for inferring paired spillover effects. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, pages 1679–1682, 2018.

[100] M. Research. Econml: A python package for ml-based heterogeneous treatment effects estimation. https://github.com/microsoft/EconML, 2019. Version 0.x.

[101] J. Robins, M. Sued, Q. Lei-Gomez, and A. Rotnitzky. Comment: Performance of double-robust estimators when" inverse probability" weights are highly variable. *Statistical Science*, 22(4):544–559, 2007.

[102] J. M. Robins. Optimal structural nested models for optimal sequential decisions. In *Proceedings of the second seattle Symposium in Biostatistics*, pages 189–326. Springer, 2004.

[103] J. M. Robins, A. Rotnitzky, and L. P. Zhao. Estimation of regression coefficients when some regressors are not always observed. *Journal of the American statistical Association*, 89(427):846–866, 1994.

[104] P. M. Robinson. Root-n-consistent semiparametric regression. *Econometrica: Journal of the Econometric Society*, pages 931–954, 1988.

[105] P. R. Rosenbaum. Model-based direct adjustment. *Journal of the American Statistical Association*, 82(398):387–394, 1987.

[106] P. R. Rosenbaum and D. B. Rubin. The central role of the propensity score in observational studies for causal effects. *Biometrika*, 70(1):41–55, 1983.

[107] P. R. Rosenbaum and D. B. Rubin. Reducing bias in observational studies using subclassification on the propensity score. *Journal of the American statistical Association*, 79(387):516–524, 1984.

[108] P. R. Rosenbaum and D. B. Rubin. Constructing a control group using multivariate matched sampling methods that incorporate the propensity score. *The American Statistician*, 39(1):33–38, 1985.

[109] N. Rosenfeld, Y. Mansour, and E. Yom-Tov. Predicting counterfactuals from large historical data and small randomized trials. In *Proceedings of the 26th International Conference on World Wide Web Companion*, pages 602–609, 2017.

[110] D. B. Rubin. Matching to remove bias in observational studies. *Biometrics*, pages 159–183, 1973.

[111] D. B. Rubin. Estimating causal effects of treatments in randomized and nonrandomized studies. *Journal of educational Psychology*, 66(5):688, 1974.

[112] D. B. Rubin and N. Thomas. Matching using estimated propensity scores: relating theory to practice. *Biometrics*, pages 249–264, 1996.

[113] D. B. Rubin and N. Thomas. Combining propensity score matching with additional adjustments for prognostic covariates. *Journal of the American Statistical Association*, 95(450):573–585, 2000.

[114] M. Saeed, M. Villarroel, A. T. Reisner, G. Clifford, L.-W. Lehman, G. Moody, T. Heldt, T. H. Kyaw, B. Moody, and R. G. Mark. Multiparameter intelligent monitoring in intensive care ii (mimic-ii): a public-access intensive care unit database. *Critical care medicine*, 39(5):952, 2011.

[115] Y. Saito. Eliminating bias in recommender systems via pseudo-labeling. *arXiv preprint arXiv:1910.01444*, 2019.

[116] B. C. Sauer, M. A. Brookhart, J. Roy, and T. VanderWeele. A review of covariate selection for non-experimental comparative effectiveness research. *Pharmacoepidemiology and drug safety*, 22(11):1139–1145, 2013.

[117] D. Scharfstein, A. Rotnitzky, and J. Robins. Comments and rejoinder. *Journal of the American Statistical Association*, 94(448):1121–1146, 1999.

[118] T. Schnabel, A. Swaminathan, A. Singh, N. Chandak, and T. Joachims. Recommendations as treatments: Debiasing learning and evaluation. In *International Conference on Machine Learning*, pages 1670–1679, 2016.

[119] P. Schwab, L. Linhardt, S. Bauer, J. M. Buhmann, and W. Karlen. Learning counterfactual representations for estimating individual dose-response curves. *arXiv preprint arXiv:1902.00981*, 2019.

[120] P. Schwab, L. Linhardt, and W. Karlen. Perfect match: A simple method for learning representations for counterfactual inference with neural networks. *arXiv preprint arXiv:1810.00656*, 2018.

[121] U. Shalit. Can we learn individual-level treatment policies from clinical data? *Biostatistics (Oxford, England)*, 11 2019.

[122] U. Shalit, F. D. Johansson, and D. Sontag. Estimating individual treatment effect: generalization bounds and algorithms. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 3076–3085, 2017.

[123] C. R. Shalizi and A. C. Thomas. Homophily and contagion are generically confounded in observational social network studies. *Sociological methods & research*, 40(2):211–239, 2011.

[124] A. Sharma, E. Kiciman, et al. DoWhy: A Python package for causal inference. https://github.com/microsoft/dowhy, 2019.

[125] E. Sherman and I. Shpitser. Identification and estimation of causal effects from dependent data. In *Advances in Neural Information Processing Systems*, pages 9424–9435, 2018.

[126] Y. Shimoni, C. Yanover, E. Karavani, and Y. Goldschmnidt. Benchmarking Framework for Performance-Evaluation of Causal Inference Analysis. *ArXiv preprint arXiv:1802.05046*, 2018.

[127] I. Shpitser. Segregated graphs and marginals of chain graph models. In *Advances in Neural Information Processing Systems*, pages 1720–1728, 2015.

[128] J. Smith. A critical survey of empirical methods for evaluating active labor market policies. Technical report, Research Report, 2000.

[129] J. Splawa-Neyman, D. M. Dabrowska, and T. Speed. On the application of probability theory to agricultural experiments. essay on principles. section 9. *Statistical Science*, pages 465–472, 1990.

[130] M. Stephen and W. Christopher. *Counterfactuals and Causal Inference: Methods and Principles for Social Research.* Cambridge University Press Cambridge, UK, 2007.

[131] E. A. Stuart. Matching methods for causal inference: A review and a look forward. *Statistical science: a review journal of the Institute of Mathematical Statistics*, 25(1):1, 2010.

[132] W. Sun, P. Wang, D. Yin, J. Yang, and Y. Chang. Causal inference via sparse additive models with application to online advertising. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*, page 297âĂŞ303, 2015.

[133] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, pages 3104–3112, 2014.

[134] A. Swaminathan and T. Joachims. Counterfactual risk minimization: Learning from logged bandit feedback. In *International Conference on Machine Learning*, pages 814–823, 2015.

[135] A. Swaminathan, A. Krishnamurthy, A. Agarwal, M. Dudik, J. Langford, D. Jose, and I. Zitouni. Off-policy evaluation for slate recommendation. In *Advances in Neural Information Processing Systems*, pages 3632–3642, 2017.

[136] E. J. T. Tchetgen and T. J. VanderWeele. On causal inference in the presence of interference. *Statistical methods in medical research*, 21(1):55–75, 2012.

[137] G. Tennenholtz, S. Mannor, and U. Shalit. Off-policy evaluation in partially observable environments. *arXiv preprint arXiv:1909.03739*, 2019.

[138] V. Veitch, Y. Wang, and D. Blei. Using embeddings to correct for unobserved confounding in networks. In *Advances in Neural Information Processing Systems*, pages 13769–13779, 2019.

[139] T. Verma and J. Pearl. *Equivalence and synthesis of causal models.* UCLA, Computer Science Department, 1991.

[140] M. Volodymyr, K. Koray, S. David, A. R. Andrei, and V. Joel. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.

[141] S. Wager and S. Athey. Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523):1228–1242, 2018.

[142] P. Wang, W. Sun, D. Yin, J. Yang, and Y. Chang. Robust tree-based causal inference for complex ad effectiveness analysis. In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining*, pages 67–76, 2015.

[143] X. Wang, R. Zhang, Y. Sun, and J. Qi. Doubly robust joint learning for recommendation on data missing not at random. In *International Conference on Machine Learning*, pages 6638–6647, 2019.

[144] C. Watkins. *Learning From Delayed Rewards.* PhD thesis, King's College, Cambridge, 1989.

[145] C. J. Watkins and P. Dayan. Q-learning. *Machine learning*, 8(3-4):279–292, 1992.

[146] J. N. Weinstein, E. A. Collisson, G. B. Mills, K. R. M. Shaw, B. A. Ozenberger, K. Ellrott, I. Shmulevich, C. Sander, J. M. Stuart, C. G. A. R. Network, et al. The cancer genome atlas pan-cancer analysis project. *Nature genetics*, 45(10):1113, 2013.

[147] J. M. Wooldridge. *Econometric analysis of cross section and panel data*. MIT press, 2010.

[148] J. M. Wooldridge. Should instrumental variables be used as matching variables? *Research in Economics*, 70(2):232–237, 2016.

[149] L. Yao, S. Li, Y. Li, M. Huai, J. Gao, and A. Zhang. Representation learning for treatment effect estimation from observational data. In *Advances in Neural Information Processing Systems*, pages 2633–2643, 2018.

[150] L. Yao, S. Li, Y. Li, M. Huai, J. Gao, and A. Zhang. Ace: Adaptively similarity-preserved representation learning for individual treatment effect estimation. In *2019 IEEE International Conference on Data Mining*, pages 1432–1437, 2019.

[151] L. Yao, S. Li, Y. Li, H. Xue, J. Gao, and A. Zhang. On the estimation of treatment effect with text covariates. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence*, pages 4106–4113, 2019.

[152] J. Yoon, J. Jordon, and M. van der Schaar. GANITE: estimation of individualized treatment effects using generative adversarial nets. In *6th International Conference on Learning Representations*, 2018.

[153] B. Yuan, J. Hsia, M. Yang, H. Zhu, C. Chang, Z. Dong, and C. Lin. Improving ad click prediction by considering non-displayed events. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management, 2019*, pages 329–338, 2019.

[154] S. Zhao and N. Heffernan. Estimating individual treatment effects from educational studies with residual counterfactual networks. In *10th International Conference on Educational Data Mining*, 2017.

[155] H. Zou, K. Kuang, B. Chen, P. Chen, and P. Cui. Focused context balancing for robust offline policy evaluation. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pages 696–704, 2019.