

A Comprehensive Framework for the Evaluation of Individual Treatment Rules From Observational Data

François Grolleau, François Petit, Raphaël Porcher

▶ To cite this version:

François Grolleau, François Petit, Raphaël Porcher. A Comprehensive Framework for the Evaluation of Individual Treatment Rules From Observational Data. 2022. hal-03735405

HAL Id: hal-03735405 https://hal.science/hal-03735405

Preprint submitted on 21 Jul 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

A Comprehensive Framework for the Evaluation of Individual Treatment Rules From Observational Data

François Grolleau*

Université Paris Cité, Centre de Recherche Épidémiologie et Statistiques (CRESS-UMR1153), INSERM, INRAE, Paris, France Centre d'Épidémiologie Clinique, Assistance Publique-Hôpitaux de Paris, Hôtel-Dieu, Paris, France,

François Petit

Université Paris Cité, Centre de Recherche Épidémiologie et Statistiques (CRESS-UMR1153), INSERM, INRAE, Paris, France,

and Raphaël Porcher

Université Paris Cité, Centre de Recherche Épidémiologie et Statistiques (CRESS-UMR1153), INSERM, INRAE, Paris, France Centre d'Épidémiologie Clinique, Assistance Publique-Hôpitaux de Paris, Hôtel-Dieu, Paris, France

July 21, 2022

^{*}The authors gratefully acknowledge the Agence Nationale de la Recherche who partially funded this work under grant agreement no. ANR-18-CE36-0010-01. François Petit was supported by the IdEx Université Paris Cité, ANR-18-IDEX-0001. Raphaël Porcher acknowledges the support of the French Agence Nationale de la Recherche as part of the "Investissements d'avenir" program, reference ANR-19-P3IA-0001 (PRAIRIE 3IA Institute).

Abstract

Individualized treatment rules (ITRs) are deterministic decision rules that recommend treatments to individuals based on their characteristics. Though ubiquitous in medicine, ITRs are hardly ever evaluated in randomized controlled trials. To evaluate ITRs from observational data, we introduce a new probabilistic model and distinguish two situations: i) the situation of a newly developed ITR, where data are from a population where no patient implements the ITR, and ii) the situation of a partially implemented ITR, where data are from a population where the ITR is implemented in some unidentified patients. In the former situation, we propose a procedure to explore the impact of an ITR under various implementation schemes. In the latter situation, on top of the fundamental problem of causal inference, we need to handle an additional latent variable denoting implementation. To evaluate ITRs in this situation, we propose an estimation procedure that relies on an expectationmaximization algorithm. In Monte Carlo simulations our estimators appear unbiased with confidence intervals achieving nominal coverage. We illustrate our approach on the MIMIC-III database, focusing on ITRs for dialysis initiation in patients with acute kidney injury.

Keywords: Personalized medicine, Causal inference, Off-policy evaluation, Mixture of experts, Expectation-maximization algorithm

1 Introduction

Individualized treatment rules (ITRs) are decision rules that recommend treatments to individuals based on their observed characteristics to maximize favorable outcomes on average. ITRs are widespread in medicine. In fact, most guidelines as well as the recently released computerized clinical decision support tools can be viewed as ITRs (Sutton et al. 2020). Notable examples include decision tools for revascularization strategies in patients with coronary artery disease (Takahashi et al. 2020) and for the personalization of blood pressure targets in hypertensive patients (Basu et al. 2017). For evaluating the impact of an ITR, the gold standard would be to conduct a randomized controlled trial (RCT) comparing the implementation of that ITR to usual care. Yet, there are practical challenges to conducting such RCTs (Tannock & Hickman 2016). As ITRs often recommend treatments similar to usual care, the expected population-level effect is likely small and necessitating very large sample sizes. Moreover, health agency oversight is less stringent for the implementation of ITRs than for drug compounds and so, both the incentives and funding opportunities for conducting RCTs of ITRs are scarce. In practice, these RCTs remain rare. As a result, many ITRs are being implemented despite the lack of evidence supporting their benefit. In this paper, we develop a framework to evaluate from observational data the impact of ITRs.

To fit with most ITRs available in medicine (e.g., computerized clinical decision support tools, or guidelines), we view ITRs as deterministic maps recommending one of two treatment options. To make inference accounting for real-life prescription of treatment by physicians, we consider that deterministic ITRs are stochastically implemented with a probability of implementation depending on patient characteristics. Critically, we then distinguish two situations: i) the ITR was just released and treatment prescription was never based on it in the population, or ii) the ITR was available and, for some patients in the population, treatment prescription was based on it. We term these two situations the *new ITR* and the *partially implemented ITR* situations, respectively. In the former situation, we propose to numerically explore the benefit an ITR may have under different implementation schemes. In the latter situation, inference is more challenging as we are typically given observational data where we do not know which patients had implemented the ITR. That is, on top of the fundamental problem of causal inference, we need to handle an additional latent variable denoting implementation. To address this situation, we develop a new probability model and rely on a mixture of experts fitted via an EM algorithm for inference.

ITR estimation and evaluation has been considered in the literature of both statistics (Qian & Murphy 2011, Zhao et al. 2012, Luedtke & van der Laan 2016) and machine learning (Kallus 2018, Thomas & Brunskill 2016). Works most related to ours include the evaluation of stochastic rules (Díaz & van der Laan 2013), biomarker performance (Janes et al. 2014), and ITR value accounting for the number of treated units (Imai & Li 2021). To our knowledge, no work has focused on data originating from a *partially implemented ITR situation*, nor pursued to develop a comprehensive framework for the evaluation of ITRs from observational data.

This article is organized as follows. In the evaluation metrics section, we introduce our causal model as well as our three estimands of interest: the Average Rule Effect (ARE), the Average Implementation Effect (AIE), and the Maximal Implementation Gain (MIG). In the inference section, we provide a method to estimate the ARE, AIE, and MIG and compute their standard error in both the *new ITR* and *partially implemented ITR* situations. In the simulation section, we study the properties of our estimators in the more challenging *partially implemented ITR situation*. Finally, in the application section, we illustrate our approach on the MIMIC-III database, focusing on ITRs for dialysis initiation in patients with acute kidney injury. We evaluate two ITRs corresponding to the *new ITR* and the *partially implemented ITR* situations. The computer code for simulation studies and data applications is available at https://github.com/fcgrolleau/ITReval.

2 Setup and evaluation metrics

Following Neyman-Rubin causal model, we consider that a patient with observed outcome Y has two potential outcomes $Y^{a=0}$ and $Y^{a=1}$ representing the outcome s/he would achieve if, possibly contrary to fact, s/he had received treatment option A = 0 or A = 1 respectively (Neyman 1923, Rubin 1974). Without loss of generality, we consider A = 1 indicates that a patient received a specific treatment, and A = 0 indicates s/he received a control. Additionally, we consider for each patient, a vector of pre-treatment covariates X with values in \mathcal{X} .

We assume that we are given an ITR that is, a deterministic map $r: \mathcal{X} \to \{0, 1\}$ which assigns a treatment option to each patient with covariates x. We model the implementation of the rule by the binary random variable S where S = 1 indicates that, based on the ITR, the physician prescribed the recommended treatment. For the rest of this paper, we term S = 1 as implementing the ITR. On the contrary, S = 0 indicates that the physician did not use the ITR to prescribe the treatment and therefore did not *implement* the ITR. Note that when S = 0 the prescribed treatment may still match the treatment recommended by the ITR—i.e., the physician did not base her/his decision on the ITR recommendation but chose the treatment on other grounds. We define the stochastic implementation function as the conditional distribution $\rho(x) = \mathbb{E}[S|X = x].^1$

¹Note that we consider here the stochastic implementation of a deterministic rule r through ρ . This is

We define the propensity score π as the conditional distribution $\pi(x) = \mathbb{E}[A|X = x]$ and the treatment-specific prognostic functions μ_0 , μ_1 as the functions satisfying $\mu_1(x) = \mathbb{E}[Y^{a=1}|X = x]$ and $\mu_0(x) = \mathbb{E}[Y^{a=0}|X = x]$. We denote τ the individual treatment effect (ITE) function i.e., $\tau(x) = \mathbb{E}[Y^{a=1} - Y^{a=0}|X = x] = \mu_1(x) - \mu_0(x)$.

2.1 A probability model for the data generating mechanism

Our goal in this subsection is to introduce a new causal model that allows to determine the causal effect of implementing versus not implementing an ITR. We introduce $A^{s=1}$, the potential treatment that would be given to a patient if her/his physician implemented the ITR i.e., $A^{s=1} = r(X)$, and $A^{s=0}$ the potential treatment s/he would be given if her/his physician did not implement the ITR. Similarly, we define $Y^{s=1}$ and $Y^{s=0}$, patient's potential outcomes when physicians do or do not implement the ITR, respectively. We further imagine the following two situations:

- A. The situation where physicians never implement the ITR. In this situation, we write with superscript $(-)^{s=0}$ the observable random variables of the patients of these physicians. That is, for the patients of these physicians $X = X^{s=0}, S = S^{s=0}, A =$ $A^{s=0}, Y = Y^{s=0}$. From this point onward, we call this situation the *new ITR situation*.
- **B.** The situation where physicians sometimes implement the ITR to prescribe treatment. In this situation, we write with superscript $(-)^*$ the observable random variables of the patients of these physicians. That is, for the patients of these physicians $X = X^*, S = S^*, A = A^*, Y = Y^*$. For the remainder of this paper, we refer to this situation as the *partially implemented ITR situation*.

different from defining a function $\mathcal{X} \to [0; 1]$ which would assign to each value x a probability to allocate treatment A = 1. This would correspond to what we call a stochastic rule — which r is not.

We consider that implementing the ITR has no impact on X, that is, we consider that the equality $X^{s=0} = X^*$ holds. For clarity, we thus drop superscripts on X. To identify causal effects, we rely on the subsequent assumptions of consistency, exchangeability and overlap.

Assumption 1 (consistency). The effect of the ITR on the outcome Y is only mediated through the treatment, that is,

$$Y^{s=1} = A^{s=1}Y^{a=1} + (1 - A^{s=1})Y^{a=0},$$
(1)

$$Y^{s=0} = A^{s=0}Y^{a=1} + (1 - A^{s=0})Y^{a=0},$$
(2)

$$A^* = S^* A^{s=1} + (1 - S^*) A^{s=0}, (3)$$

$$Y^* = A^* Y^{a=1} + (1 - A^*) Y^{a=0}.$$
(4)

Assumption 2 (exchangeability). All confounders and variables causing implementation are measured, that is,

$$\{Y^{a=1}, Y^{a=0}\} \perp A^{s=0} | X,$$
 (5)

$$\{Y^{a=1}, Y^{a=0}\} \perp A^* | X,$$
 (6)

$$A^{s=0} \perp \!\!\!\perp S^* | X. \tag{7}$$

Assumption 3 (overlap). Within all realistic levels of covariates, the patients could receive either treatment—including in the absence of ITR implementation. That is, denoting $\pi^{s=0}$ the propensity score in the absence of ITR implementation, i.e., $\pi^{s=0}(x) = \mathbb{E}[A^{s=0}|X = x]$,

$$\forall x \in \mathcal{X}, \quad 0 < \pi(x) < 1, \text{ and } 0 < \pi^{s=0}(x) < 1.$$

Observing the overlap assumption, we see that as the propensity score functions π and $\pi^{s=0}$ can never be deterministic rules, they are thus stochastic rules. We define two additional stochastic rules: the propensity score under stochastic implementation π^* that is $\pi^*(x) = \mathbb{E}[A^*|X = x]$, and the stochastic implementation function under implementation ρ^* i.e., $\rho^*(x) = \mathbb{E}[S^*|X = x]$.

Summarizing the consistency equations (1), (2), (3), (4) and the exchangeability equations (5), (6), (7), the data generating mechanism in the *new ITR* and the *partially implemented ITR* situations can be represented by the probabilistic graphical models in Figure 1A and 1B, respectively.



Figure 1: The probabilistic graphical model associated with the data generating mechanism in the *new ITR situation* (Panel A) and the *partially implemented ITR situation* (Panel B).

In our setup, it will prove convenient to define q_1 , the prognostic function under ITR implementation, as $q_1(x) = \mathbb{E}[Y^{s=1}|X=x]$ and q_0 , the prognostic function in the absence of ITR implementation, as $q_0(x) = \mathbb{E}[Y^{s=0}|X=x]$. Conditioning equation (1) with respect to X leads to $q_1(x) = r(x)\mu_1(x) + \{1 - r(x)\}\mu_0(x)$.

2.2 Estimands of interest

We now introduce three estimands. First, the Average Rule Effect (ARE) of an ITR r:

$$\Delta(r) = \mathbb{E}[Y^{s=1} - Y^{s=0}].$$

This represents the population-level effect of the ITR on outcome Y in a randomized trial comparing a group where patients are systematially given the treatment recommended by ITR to usual care in the absence of ITR implementation.

Second, we define the Average Implementation Effect (AIE) of r as

$$\Lambda(r, \rho^*) = \mathbb{E}[Y^* - Y^{s=0}].$$

This represents the population-level effect of the ITR on outcome Y in a randomized trial comparing a group where physicians are provided with the ITR's treatment recommendation to usual care under no implementation. We may thus consider that in the experimental group the treatment A is prescribed according to a stochastic implementation ρ^* of r. Note that the ARE can be considered as a special case of the AIE where the stochastic implementation is perfect ($\rho^* \equiv 1$). We however single it out because this representation allows to assess whether the ITR has a potential population-level benefit, or if it is instead poorly designed.

Last, we define the Maximal Implementation Gain (MIG) of r:

$$\Gamma(r, \rho^*) = \mathbb{E}[Y^{s=1} - Y^*].$$

This represents the difference in average outcome between a full implementation of the ITR and the current or future partial implementation of the rule. From these definitions, it follows that

$$\Delta(r) = \Lambda(r, \rho^*) + \Gamma(r, \rho^*).$$

3 Inference

We assume we are given a single random sample of n independent and identically distributed (i.i.d.) units $(X_i^T, A_i, Y_i)_{1 \le i \le n}$ from a target population. As in the previous section, we distinguish between data originating from *new* and *partially implemented* ITR situations:

- A. In data originating from the *new ITR situation*, we have $S_i = S_i^{s=0} = 0$, $A_i = A_i^{s=0}$, and $Y_i = Y_i^{s=0}$. Note that the second equality implies $\pi = \pi^{s=0}$. In this situation, we drop all $(-)^{s=0}$ superscripts and use π rather than $\pi^{s=0}$ for clarity. Clearly, estimation of the AIE and MIG is not possible from data alone in this situation. Nonetheless, in the next section, we propose to explore their behaviour by hypothesizing implementation schemes.
- **B.** In data originating from the *partially implemented ITR situation*, we have $S_i = S_i^*$, $A_i = A_i^*$, and $Y_i = Y_i^*$. Note that the first two equalities imply $\rho = \rho^*$ and $\pi = \pi^*$. In this situation, we thus drop all $(-)^*$ superscripts for clarity. Because we expect that S_i will not have been collected in this situation, we treat it as a latent variable.

3.1 New individualized treatment rule situation

3.1.1 Average rule effect

Using consistency (1), exchangeability (2), and positivity (3), in the *new ITR situation*, we have

$$\Delta(r) = \mathbb{E}\left[q_1(X) - Y\right]$$
$$= \mathbb{E}\left[\left\{\frac{A}{\pi(X)} + \{1 - r(X)\}\frac{1 - A}{1 - \pi(X)} - 1\right\}Y\right].$$

These equations suggest the following two estimators for $\Delta(r)$

$$\widehat{\Delta}_Q(r) = n^{-1} \sum_{i=1}^n r(X_i) \widehat{\mu}_1(X_i) + (1 - r(X_i)) \widehat{\mu}_0(X_i) - Y_i,$$
(8)

$$\widehat{\Delta}_{IPW}(r) = n^{-1} \sum_{i=1}^{n} \left[r(X_i) \frac{A_i}{\widehat{\pi}(X_i)} + \{1 - r(X_i)\} \frac{1 - A_i}{1 - \widehat{\pi}(X_i)} - 1 \right] Y_i,$$
(9)

where as for all estimators proposed hereafter, $\mu_0(\cdot)$, $\mu_1(\cdot)$, and $\pi(\cdot)$ can be estimated via any supervised learning method from observations $(X_i^{\mu}, Y_i)_{i:A_i=0}$, $(X_i^{\mu}, Y_i)_{i:A_i=1}$, and $(X_i^{\pi}, A_i)_{1 \le i \le n}$ respectively.² An augmented counterpart of these estimators can be derived from Zhang et al. (2012):

$$\widehat{\Delta}_{AIPW}(r) = n^{-1} \sum_{i=1}^{n} \left[\frac{\mathcal{C}_{i}^{r} Y_{i}}{\widehat{\pi}(X_{i}) \mathcal{C}_{i}^{r} + \{1 - \widehat{\pi}(X_{i})\}(1 - \mathcal{C}_{i}^{r})} - \frac{\mathcal{C}_{i}^{r} - [\widehat{\pi}(X_{i}) \mathcal{C}_{i}^{r} + \{1 - \widehat{\pi}(X_{i})\}(1 - \mathcal{C}_{i}^{r})]}{\widehat{\pi}(X_{i}) \mathcal{C}_{i}^{r} + \{1 - \widehat{\pi}(X_{i})\}(1 - \mathcal{C}_{i}^{r})} \widehat{q}_{1}(X_{i}) - Y_{i} \right]$$
(10)

where we set $C_i^r = \mathbb{1}\{r(X_i) = A_i\}$ and $\hat{q}_1(X_i) = r(X_i)\hat{\mu}_1(X_i) + \{1 - r(X_i)\}\hat{\mu}_0(X_i)$ for clarity. We refer the reader to Tsiatis et al. (2019, section 3.3.3 p. 61) for an extensive study of this specific case and the derivation of approximate large sample distribution. Using the ITE, the ARE can also be reformulated as

$$\Delta(r) = \mathbb{E}\left[\{r(X) - \pi(X)\}\tau(X)\right]$$
(11)

(a proof is given in Appendix A). This leads to the following estimator

$$\widehat{\Delta}_{ITE}(r) = n^{-1} \sum_{i=1}^{n} \{ r(X_i) - \widehat{\pi}(X_i) \} \widehat{\tau}(X_i).$$

Though the latter estimator requires to estimate the ITE τ , and hence may be less practical than estimators (8), (9), or (10), the equation (11) makes explicit the respective contribution of τ , r and π to the ARE.

²Here, X_i^{μ} and X_i^{π} denote two subsets of the relevant variables contained in X_i .

3.1.2 AIE, MIG under the modeling of the stochastic implementation functions

When the ITR is new and has never been deployed, the way in which it will be implemented is unpredictable. Hence, the AIE and the MIG cannot be estimated from data alone. However, it can be interesting to study numerically how the AIE and MIG would vary under different stochastic implementation schemes as this can provide information about the appropriateness of future ITR deployment. In the *new ITR situation*, it is possible to show that

$$\Lambda(r,\rho^*) = \mathbb{E}\left[\left\{\pi^*(X) - \pi(X)\right\}\tau(X)\right]$$
(12)

$$= \mathbb{E}\left[\rho^*(X)\{r(X) - \pi(X)\}\tau(X)\right],\tag{13}$$

and

$$\Gamma(r, \rho^*) = \mathbb{E}\left[\{r(X) - \pi^*(X)\}\tau(X)\right]$$
$$= \mathbb{E}\left[\{1 - \rho^*(X)\}\{r(X) - \pi(X)\}\tau(X)\right]$$

(a proof is given in Appendix B). Hence, given an estimate $\hat{\tau}$ of the ITE function (see Jacob 2021, for a review of the available estimation methods), an estimate $\hat{\pi}$ of the propensity score and a numerical model of ρ^* for a future stochastic implementation function, estimates of the AIE and the MIG are computable via

$$\widehat{\Lambda}_{ITE}(r,\rho^*) = n^{-1} \sum_{i=1}^n \rho^*(X_i) \{ r(X_i) - \widehat{\pi}(X_i) \} \widehat{\tau}(X_i),$$

and

$$\widehat{\Gamma}_{ITE}(r,\rho^*) = n^{-1} \sum_{i=1}^n \{1 - \rho^*(X_i)\} \{r(X_i) - \widehat{\pi}(X_i)\} \widehat{\tau}(X_i)$$

Below, we propose, three schemes that model the form the stochastic implementation function may take in future deployment of the ITR: • The random implementation scheme, where we model $\rho^*(\cdot)$ as

$$\rho_{rd,\alpha}^*(x) = \alpha$$

with $\alpha \in [0; 1]$ a parameter modelling the random implementation such that, uniformly for all patients, higher values of α are associated with higher probabilities of following the rule. This model of ρ^* describes a situation where patients are treated according to an implementation of the rule at random with probability α regardless of their characteristics.

• The cognitive bias scheme, where we model $\rho^*(\cdot)$ as

$$\rho_{cb,\alpha}^*(x) = \{1 - |r(x) - \pi(x)|\}^{\frac{1}{2}\log\frac{\alpha+1}{1-\alpha}}$$

with $\alpha \in [0; 1]$ a cognitive bias parameter such that higher values of α are associated with lower probabilities of following the rule for a given gap between the recommendation from the ITR and usual care under no implementation. This implementation scheme describes a situation where physicians follow the ITR recommendation more often when recommendations are similar to current practices and this trend to resist change increases as α increases.

The confidence level scheme, where we assume that the ITR was constructed from estimated ITEs, τ(x), as in for instance r(x) = 1[τ(x) < 0]. For this scheme to be actionable, τ(x) and their standard errors se_{τ(x)} must be provided along the ITR they helped build. Under such conditions, we model ρ*(·) as

$$\rho_{cl,\alpha}^*(x) = \mathbb{1}[\{\tilde{\tau}(x) - q_{1-\alpha/2}se_{\tilde{\tau}(x)}\}\{\tilde{\tau}(x) + q_{1-\alpha/2}se_{\tilde{\tau}(x)}\} > 0]$$

for ITEs provided on an absolute scale (i.e., individual absolute risk difference) with $\alpha \in [0; 1]$ a type I error parameter such that smaller values of α lead to wider confidence intervals for $\tilde{\tau}(x)$. This scheme describes a situation where physicians follow the ITR recommendation only when there is evidence that $\tau(x) \neq 0$ at significance level α .

3.1.3 Illustrative examples

In this section, we aim to provide a sense of what our method is trying to achieve when applied in the *new ITR situation*. For that purpose, in this subsection, we provide a toy model. Observing Formula (12), we see that the AIE of a new ITR r gets far off from zero as increases the difference between current treatment allocation and future treatment allocation under a stochastic implementation of r. More precisely, observing Formula (13), we note that the AIE of a new ITR r gets far off from zero as patients with common levels of covariates x have i) a high probability $\rho^*(x)$ of implementing the rule, and/or ii) a difference $r(x) - \pi(x)$ between recommendation from the ITR and usual care under no implementation far off from zero, and/or iii) large ITEs $\tau(x)$.

For illustration purposes, we imagined a disease for which only one patient characteristic, the age x, is relevant to treatment decision-making. In a population of patients with mean age 50 (standard deviation 15), we wish to evaluate the effectiveness of an ITR r with respect to the occurrence of an unfavorable binary outcome (i.e., 10-year mortality). In our two examples, ground truth is such that the treatment is beneficial for patients aged 40 to 60, detrimental for patients aged 60 to 80, and has almost no effect outside these ranges. For the sake of simplicity, we suppose that in both examples r is $r(X) = \mathbb{1}[\tau(X) < 0]$ that is, the rule is optimal (Figure 2 Panels A and B).

In our first example (Figure 2A), the usual care under no implementation is such that younger patients are treated more often while in our second example (Figure 2B), older patients are treated more often. In the random implementation schemes (Figure 2 Panels C and D), physicians follow the ITR's recommendation at random with probability 1/3 (red lines) or 2/3 (green lines). In the cognitive bias schemes (Figure 2 Panels E and F), physicians follow the ITR's recommendation more often when the ITR's recommendation tracks the usual care under no implementation. Cognitive bias parameter is 2/3 (red lines) or 1/3 (green lines), and higher parameter values are associated with lower probabilities of complying with the ITR. In the confidence level schemes (Figure 2 Panels G and H), physicians follow the ITR's recommendation only when confidence intervals for the predicted ITEs do no cross zero. Type I error parameters for the confidence intervals are 0.05 (red lines) or 0.45 (green lines) with higher values associated with tighter confidence intervals and therefore higher probabilities of implementing the ITR.

Despite the fact that both examples relied on the implementation of an identical ITR based on the true ITE function, the population-level benefit of this ITR is different between examples for all schemes. The ARE of the deterministic rule was -0.16 in the population from example 1 and -0.24 in the population from example 2, indicating an 8% greater benefit of implementing the ITR in population 2 than in population 1 if physicians always followed the ITR's recommendation. Similarly, in the stochastic implementation schemes, the population-level benefit of the ITR differ in the two example populations. In the random implementation scheme, AIEs are -0.11 and -0.05 in population 1 (Figure 2C) versus -0.16 and -0.08 in population 2 (Figure 2D) for random implementation parameters 2/3 and 1/3 respectively. We find similar differences in AIEs between population 1 and population 2 in the cognitive bias scheme (Figure 2 Panels E and F) and confidence level scheme (Figure 2 Panels G and H).



Figure 2: Panel A displays our first illustrative example where under no implementation usual care is such that younger patients are treated more often. Panel B displays our second illustrative example where under no implementation usual care is such that older patients are treated more often. The random implementation scheme is given panels C and D for both examples, respectively, the cognitive bias scheme on panels E and F, and the confidence level scheme panels G and H. The dotted lines correspond to the ITEs plus/minus its standard errors. AIEs are reported for each of implementation schemes and lower values indicate greater benefit from ITR implementation. We denote p_X the probability density of X (we re-scale $p_X(x)$ by a factor 35 for illustration purposes).

3.2 Partially implemented individualized treatment rule situation

In this subsection, our aim is to estimate the ARE, AIE and MIG with data sampled from a population where the ITR r is partially implemented. Two cases have to be distinguished depending on whether the variable S is collected. If S is collected, estimating the ARE, AIE and MIG can be achieved by using a suitable adaptation of the IPW/AIPW estimators for the average treatment effect (Lunceford & Davidian 2004). However, in practice, we expect that the variable S will not have been collected. Hence, we focus our attention on the case where S needs to be regarded as a latent variable. Recall that in this subsection, we are dealing with a partially implemented ITR where the observed treatment follows either from the ITR being implemented or from the physicians disregarding the ITR to make treatment decisions. Inference in the partially implemented ITR situation relies on assumptions (1-3). Because in this situation estimation of the MIG is more straightforward than estimation of the ARE and AIE, we distinguish between the two cases.

3.2.1 Maximal implementation gain

We start by studying the MIG, as neither S nor ρ play a role for this estimand in the *partially* implemented ITR situation. In fact in this situation, using consistency (1), exchangeability (2), and positivity (3), we have

$$\Gamma(r,\rho) = \mathbb{E}\left[q_1(X) - Y\right]$$
$$= \mathbb{E}\left[\left\{r(X)\frac{A}{\pi(X)} + \{1 - r(X)\}\frac{1 - A}{1 - \pi(X)} - 1\right\}Y\right].$$

This suggests the estimators

$$\widehat{\Gamma}_Q(r,\rho) = n^{-1} \sum_{i=1}^n r(X_i)\widehat{\mu}_1(X_i) + \{1 - r(X_i)\}\widehat{\mu}_0(X_i) - Y_i,$$

and

$$\widehat{\Gamma}_{IPW}(r,\rho) = n^{-1} \sum_{i=1}^{n} \left[r(X_i) \frac{A_i}{\widehat{\pi}(X_i)} + \{1 - r(X_i)\} \frac{1 - A_i}{1 - \widehat{\pi}(X_i)} - 1 \right] Y_i.$$

The derivation is similar to that of the ARE in the *new ITR situation* (equations 8 and 9). We refer the reader to section 3.1.1, equation (10) for an augmented version of this estimator.

3.2.2 Average rule effect and average implementation effect

Note that the MIG estimand is distinct from the ARE and AIE in that is does not involve the expectation term $\mathbb{E}(Y^{s=0})$. In contrast, the ARE and AIE depend on the pairs of expectations $\mathbb{E}(Y^{s=1})$, $\mathbb{E}(Y^{s=0})$ and $\mathbb{E}(Y)$, $\mathbb{E}(Y^{s=0})$ respectively. The quantity $\mathbb{E}(Y)$ is straightforward to estimate. The expectation $\mathbb{E}(Y^{s=1})$ can be estimated by various means, for instance by taking the expectation of $\hat{q}_1(X)$ as in section 3.2.1.

Estimation of $\mathbb{E}(Y^{s=0})$ is more challenging than that of $\mathbb{E}(Y^{s=1})$ because, substitution of $Y^{s=0}$ by its definition in (2) involves the potential outcome $A^{s=0}$ which is not identifiable from equation (3) as S is a latent variable. Our approach to estimate $\mathbb{E}(Y^{s=0})$ relies on the following result.

Lemma 1. In the partially implemented ITR, the following relations holds

(i)
$$\pi(x) = \rho(x)r(x) + \{1 - \rho(x)\}\pi^{s=0}(x),$$
 (14)

(*ii*)
$$q_0(x) = \mu_1(x)\pi^{s=0}(x) + \mu_0(x)\{1 - \pi^{s=0}(x)\}.$$
 (15)

A proof of the lemma is given in Appendix C. Lemma 1 shows that the functions π and q_0 can be represented by mixtures of experts (Jordan & Jacobs 1994). In particular, observing equation (14), we see that π can be viewed as a mixture of the known expert r and the unknown expert $\pi^{s=0}$, while the component of the mixture depends on the unknown gating network ρ . Equation (15) suggest to rewrite $\Delta(r)$ and $\Lambda(r, \rho)$ as

$$\Delta(r) = \mathbb{E}\left[q_1(X) - q_0(X)\right]$$
$$= \mathbb{E}\left[\left\{r(X) - \pi^{s=0}(X)\right\}\tau(X)\right]$$

and

$$\Lambda(r,\rho) = \mathbb{E} \left[Y - q_0(X) \right]$$

= $\mathbb{E} \left[Y - \mu_1(X) \pi^{s=0}(X) + \mu_0(X) \{ 1 - \pi^{s=0}(X) \} \right].$

Because $\pi^{s=0}$ is unknown, we propose to estimate it via the procedure detailed in Algorithm 1. This procedure details an EM algorithm, based on the fitting algorithm of Xu & Jordan (1993) and Jordan & Jacobs (1994). Figure 3 depicts the graphical model for the approach to estimating $\pi^{s=0}(\cdot)$. Estimating $\mu_0(\cdot)$, $\mu_1(\cdot)$, and $\tau(\cdot)$ as in section 3.1, mixture of experts estimators for the ARE and AIE are given by

$$\widehat{\Delta}_{ME}(r) = n^{-1} \sum_{i=1}^{n} \{ r(X_i) - \widehat{\pi}^{s=0}(X_i) \} \widehat{\tau}(X_i),$$

and

$$\widehat{\Lambda}_{ME}(r,\rho) = n^{-1} \sum_{i=1}^{n} Y_i - \widehat{\mu}_1(X_i) \widehat{\pi}^{s=0}(X_i) + \widehat{\mu}_0(X_i) \{1 - \widehat{\pi}^{s=0}(X_i)\}.$$

We propose to estimate the variance of the $\widehat{\Delta}_{ME}(r)$ and $\widehat{\Lambda}_{ME}(r,\rho)$ estimators via the boostrap. We assess the validity of this strategy in Monte Carlo simulations.



Figure 3: The graphical representation of the mixture of experts fitted by Algorithm 1. Note that we consider r as a known deterministic expert network, while both the expert network $\pi^{s=0}$ and the gating network ρ are unknown stochastic rules.

Algorithm 1 The EM procedure for estimating $\pi^{s=0}(\cdot)$ in the partially implemented ITR situation.

Input: The ITR $r: \mathcal{X} \to \{0, 1\}$, and data $(X_i^T, A_i)_{1 \le i \le n}$ where $X_i^{\pi^{s0}}$ and X_i^{ρ} are two relevant subsets of

the variables contained in X_i .

Initialize the prior probabilities associated with the nodes of the tree as $g_{0,i} \leftarrow 0.5$, and $g_{1,i} \leftarrow 0.5$.

Initialize the parameters ζ of the expert $\pi^{s=0}(\cdot)$ at random e.g., $\zeta \sim \mathcal{N}(0, D)$ with D a diagonal matrix.

Compute individual contributions to r's likelihood as $P_{1,i} \leftarrow r(X_i)^{A_i} \{1 - r(X_i)\}^{1-A_i}$.

Compute individual predictions from the initiated expert network $\pi^{s=0}(\cdot)$ as $p_{0,i} \leftarrow \operatorname{expit}(\zeta^T X_i^{\pi^{s0}})$.

Iterate until convergence on the parameters ζ :

Compute individual contributions to $\pi^{s=0}$'s likelihood as $P_{0,i} \leftarrow p_{0,i}^{A_i} (1-p_{0,i})^{1-A_i}$.

Compute the posterior probabilities associated with the nodes of the tree as

 \triangleright E-step

$$h_{0,i} \leftarrow \frac{g_{0,i}P_{0,i}}{g_{0,i}P_{0,i}+g_{1,i}P_{1,i}}$$
 and $h_{1,i} \leftarrow \frac{g_{1,i}P_{1,i}}{g_{0,i}P_{0,i}+g_{1,i}P_{1,i}}$

For the gating network $\rho(\cdot)$ estimate parameters γ by solving the IRLS problem \triangleright M-step

$$\gamma \leftarrow \underset{\gamma}{\arg\max} \sum_{i=1}^{n} h_{1,i} \ln \left\{ \operatorname{expit}(\gamma^{T} X_{i}^{\rho}) \right\} + (1 - h_{1,i}) \ln \left\{ 1 - \operatorname{expit}(\gamma^{T} X_{i}^{\rho}) \right\}$$

For the expert network $\pi^{s=0}(\cdot)$ estimate parameters ζ by solving the IRLS problem

$$\zeta \leftarrow \underset{\zeta}{\operatorname{arg\,max}} \sum_{i=1}^{n} h_{0,i} \Big[A_i \ln \big\{ \operatorname{expit}(\zeta^T X_i^{\pi^{s0}}) \big\} + (1 - A_i) \ln \big\{ 1 - \operatorname{expit}(\zeta^T X_i^{\pi^{s0}}) \big\} \Big]$$

Update the prior probabilities associated with the nodes of the tree as

$$g_{1,i} \leftarrow \operatorname{expit}(\gamma^T X_i^{\rho}) \quad \text{and} \quad g_{0,i} \leftarrow 1 - g_{1,i}$$

Update the predictions from the expert network $\pi^{s=0}(\cdot)$ as $p_{0,i} \leftarrow \operatorname{expit}(\zeta^T X_i^{\pi^{s=0}})$.

Return: $\hat{\pi}^{s=0}(x) = \operatorname{expit}(\zeta^T x).$

4 Simulations

4.1 Setup

In this section, we study the properties of the MIG, ARE and AIE estimators in the partially implemented ITR situation. To this end, we simulate data analysis in a setting where an ITR was partially implemented. We generate synthetic datasets comprising six Bernoulli, log-normally or normally distributed covariates $X = (X_1, X_2, \ldots, X_6)$ as follows.

Step 1 We randomly generate intermediate covariates X'_1, X'_2, \ldots, X'_6 from a multivariate gaussian distribution

$$(X_1', X_2', \ldots, X_6')^T \sim \mathcal{N}(0, \Sigma).$$

To generate Σ , we chose 6 eigenvalues $(\lambda_1, \lambda_2, \dots, \lambda_6) = (1, 1.2, 1.4, 1.6, 1.8, 2)$, and sample a random orthogonal matrix O of size 6×6 . The covariance matrix Σ is obtained via

$$\Sigma = O \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & \lambda_6 \end{bmatrix} O^T.$$

Step 2 To allow for the Bernoulli or log-normal distribution of covariates, we generate

$$X_1, X_2, \dots, X_6$$
 as follows $(X_1, X_2) = (\mathbb{1}\{X'_1 < 0\}, \mathbb{1}\{X'_2 < 0\}), (X_3, X_4, X_5) = (\exp(X'_3), \exp(X'_4), \exp(X'_5)), X_6 = X'_6$. We add $X_0 \equiv 1$ to allow for intercepts.

Step 3 We generate data from the covariates in this manner:

$$S|X = x \sim \text{Bernouilli}(\exp((\gamma^{T}x))), \qquad Y^{a=0}|X = x \sim \text{Bernouilli}(\exp((\alpha^{T}x))),$$

$$r(X) = \mathbb{1}\{\delta^{T}x < 0\}, \qquad Y^{a=1}|X = x \sim \text{Bernouilli}(\exp((\beta^{T}x))),$$

$$A^{s=0}|X = x \sim \text{Bernouilli}(\exp((\zeta^{T}x))), \qquad Y^{s=0} = A^{s=0}Y^{a=1} + (1 - A^{s=0})Y^{a=0},$$

$$A^{s=1} = r(X), \qquad Y^{s=1} = A^{s=1}Y^{a=1} + (1 - A^{s=1})Y^{a=0},$$

$$A = SA^{s=1} + (1 - S)A^{s=0}, \qquad Y = AY^{a=1} + (1 - A)Y^{a=0}$$

with

$$\gamma = (0, 0, 0, 0, 0, 0, 1)^T, \qquad \delta = (0.05, -0.5, 0.5, -0.5, 0.5, 0, 0)^T,$$

$$\alpha = (0, -0.3, -0.05, 0.5, -0.15, -0.2, 0)^T, \quad \beta = (0, -0.2, 0.05, 0.3, -0.1, -0.1, 0)^T$$

and we vary ζ .

In scenario A, we set $\zeta = \delta$ which corresponds to a situation where treatment allocation in the absence of the ITR is different from the ITR. In scenario B, we set $\zeta = (0, 0, 0, 0, 0, 0, 0)$ which corresponds to a situation where treatment allocation in the absence of the ITR is random with probability 0.5. In scenario C, we set $\zeta = -\delta$ which corresponds to a situation where treatment allocation in the absence of the ITR resembles the ITR. In each scenario we generate a target population of two million individuals from which we approximate ground truth for our estimands and drew random samples. We vary the sample size: n =200, 800, 2000. The potential outcomes as well as the variable S are regarded as unobserved variables. Models for μ_0 and μ_1 are correctly specified with $X^{\mu} = (X_2, X_3, X_4, X_5, X_6)$ as explanatory variables. We fit the mixture of expert in equation (14) with Algorithm 1, specifying the gating network ρ with $X^{\rho} = X_7$ and the expert network $\pi^{s=0}$ with variables $X^{\pi s0} = (X_1, X_2, X_3, X_4, X_5)$. For each scenario/sample size combination, we implement 1000 simulation iterations and 999 bootstrap replications to generate confidence intervals.

4.2 Results

The results of our simulations are reported in Table 1 and Figure 4. The MIG estimator $\widehat{\Gamma}_Q(r,\rho)$ which does not rely on an EM procedure exhibits, as expected, the properties of unbiasness and consistency. The ARE estimator $\widehat{\Delta}_{ME}(r)$ and the AIE estimator $\widehat{\Lambda}_{ME}(r,\rho)$ appear unbiased and consistent as well. Their standard error is comparable to that of the MIG estimator $\widehat{\Gamma}_Q(r,\rho)$. Ninety five percent bootstrap confidence intervals achieve close to nominal coverage for all three estimators.



Figure 4: Absolute bias and Root Mean Square Error (RMSE) for the MIG, ARE and AIE estimators across 1000 simulation iterations in nine scenario/sample size combinations. Absolute bias is the darker portion of each bar ; RMSE corresponds to the total bar size. The MIG, ARE and AIE estimators are from $\widehat{\Gamma}_Q(r, \rho), \widehat{\Delta}_{ME}(r), \widehat{\Lambda}_{ME}(r, \rho)$ respectively.

5 Applications on the MIMIC-III database

The Multi-Parameter Intelligent Monitoring in Intensive Care III (MIMIC-III) database is a publicly available electronic health record that contain data from 53,423 patients hos-

Table 1: Simulation results for the MIG, AIE and ARE estimators under all nine scenario/sample size combinations. The MIG, ARE and AIE estimators are from $\widehat{\Gamma}_Q(r,\rho), \widehat{\Delta}_{ME}(r), \widehat{\Lambda}_{ME}(r,\rho)$ respectively. SE: standard error; RMSE: root mean squared error; CI: 95% confidence interval. Coverage probabilities are for 95% confidence intervals.

	Scenario A			Scenario B			Scenario C		
n	MIG	ARE	AIE	MIG	ARE	AIE	MIG	ARE	AIE
True value	-0.013	-0.026	-0.013	-0.008	-0.016	-0.008	-0.004	-0.007	-0.003
Relative bias									
200	0.102	0.028	0.057	0.000	0.000	0.000	0.192	0.163	0.365
800	0.012	-0.044	-0.090	0.082	0.002	0.005	0.037	0.090	0.201
2000	-0.009	-0.027	-0.055	0.029	-0.016	-0.032	-0.013	0.069	0.155
Bias									
200	-0.001	-0.001	0.001	0.000	0.000	0.000	-0.001	-0.001	0.000
800	0.000	0.001	0.001	-0.001	0.000	0.001	0.000	-0.001	0.000
2000	0.000	0.001	0.001	0.000	0.000	0.000	0.000	0.000	-0.001
Empirical SE									
200	0.024	0.037	0.024	0.020	0.027	0.020	0.017	0.016	0.017
800	0.012	0.018	0.012	0.010	0.013	0.010	0.008	0.008	0.008
2000	0.007	0.011	0.008	0.006	0.008	0.006	0.005	0.005	0.005
RMSE									
200	0.024	0.037	0.024	0.020	0.027	0.020	0.017	0.016	0.017
800	0.012	0.018	0.012	0.010	0.013	0.010	0.008	0.008	0.008
2000	0.007	0.011	0.008	0.006	0.008	0.006	0.005	0.005	0.005
Coverage									
200	0.964	0.954	0.937	0.953	0.976	0.941	0.936	0.979	0.944
800	0.945	0.953	0.949	0.945	0.956	0.947	0.953	0.957	0.949
2000	0.959	0.951	0.934	0.948	0.946	0.950	0.962	0.946	0.947
CI width									
200	0.099	0.153	0.096	0.080	0.113	0.082	0.067	0.073	0.067
800	0.047	0.071	0.046	0.039	0.050	0.039	0.033	0.031	0.033
2000	0.029	0.044	0.029	0.024	0.031	0.024	0.021	0.019	0.021

pitalized in intensive care at Beth Israel Deaconess Medical Center from 2001 to 2012 (Johnson et al. 2016). From these, we include the 3,748 intensive care unit adult patients with severe acute kidney injury who had received either invasive mechanical ventilation or vasopressor infusion. We report the full inclusion/exclusion criteria in the Appendix D and the inclusion flow-diagram in the Appendix E. The patients we include are eligible for recommendation from both ITRs described below. For the sake of focusing on the estimators in our methodology, we handle patients with missing data by conducting a single imputation using chained equations (White et al. 2011).

In this section, we consider two example ITRs. In the first example, we evaluate a new ITR for dialysis initiation.³ This last ITR was not available at the time of data collection and decision to initiate dialysis never followed from its implementation. In the second example, we evaluate an ITR that was partially implemented at the time of data collection. Specifically, we evaluate the impact of an ITR that recommends initiating dialysis in the most severe patients based on the Sequential Organ Failure Assessment (SOFA) score (Vincent et al. 1996).

5.1 New ITR: dialysis initiation based on a combination of six biomarkers

Grolleau et al. (2022) have recently developed a new ITR for dialysis initiation in the intensive care unit using data from two RCTs. Briefly, this new ITR recommends initiating dialysis within 24 hours only in specific patients based on a combination of six biomarkers (SOFA score, pH, potassium, blood urea nitrogen, weight and, the prescription of im-

³In this section, we use the term "dialysis" loosely to refer to all kidney support therapies suitable for acute kidney injury patients (i.e., including but not limited to intermittent hemodialysis and continuous hemofiltration).

munosuppressive drug). Following the methodology detailed in section 3.1, we evaluate the impact of this new ITR on 60-day mortality. We estimate the ARE using a double robust estimator as detailed in section 3.1.1. We explore the impact of various degrees of implementation under either cognitive bias or confidence level schemes. The estimated values of AIE are given in Figure 5. Estimation of the ARE shows a trend for benefit from the implementation of the new ITR ($\hat{\Delta}_{AIPW}(r) = -0.02$; 95% confidence interval [-0.06 to 0.01]). Note that the ARE estimate is not equal to estimation of the AIE under full implementation as, contrary to the ARE case, for the AIE we use an ITE model. The variables included in each model are reported in Appendix G.



Figure 5: Evaluation of the impact of a new ITR (i.e., dialysis initiation within 24 hours only in specific patients based on a combination of six biomarkers) on 60-day mortality using the MIMIC-III observational database. Ninety-five percent confidence intervals are from the bootstrap. Blue diamonds are for the cognitive bias scheme; orange diamonds are for the confidence level scheme. Panels A depict the AIE for different values of implementation parameter α , Panels B depict the AIE as a function of the proportion of (future) patients implementing the new ITR: $n^{-1} \sum_{i=1}^{n} \rho_{\cdot,\alpha}^*(X_i)$. More negative values of the AIE indicate greater benefit from ITR implementation. Ninety-five percent confidence intervals are from the bootstrap.

5.2 Partially implemented ITR: dialysis initiation based on SOFA scores

We evaluate the impact on 60-day mortality of a partially implemented, yet never evaluated, ITR that recommended initiating dialysis within 24 hours only in the patients with a Sequential Organ Failure Assessment (SOFA) score greater than 11. Following the methodology of section 3.2, we posit models for the treatment-specific prognosis functions, the propensity score in the absence of ITR implementation (expert network) and the stochastic implementation function (gating network). Specification of propensity score in the absence of ITR implementation include the variables thought to have caused treatment initiation while specification of the stochastic implementation function include all variables thought to be associated with ITR implementation. The variables included in each model are reported in Appendix F. The estimates of the MIG, ARE and AIE are given in Figure 6. Estimation of the MIG shows evidence of harm from further implementing the ITR $(\widehat{\Gamma}_Q(r,\rho) = 4.7\%; 95\%$ confidence interval [2.7% to 6.8%]). Similarly, estimation of the ARE shows a trend for harm in implementing the ITR in all patients versus in no one $(\widehat{\Delta}_Q(r) = 1.8\%; 95\%$ confidence interval [-0.1% to 3.9%]) indicating that the ITR may be poorly designed. However, estimation of the AIE shows that the withdrawal of the ITR would on average yield outcomes worse than in the current situation ($\widehat{\Lambda}_Q(r,\rho) = -2.9\%$; 95% confidence interval [-3.1\% to -2.6\%]). This suggest that even though the ITR may be poorly designed, physicians identify correctly the patients who benefit from ITR implementation. In sum, these results indicate that neither full nor null implementation of the ITR would improve patient outcomes (at the population level). Rather, either one of these changes in ITR implementation, our analysis suggest, would worsen patient outcomes (at the population level).



Figure 6: Evaluation of the impact of a partially implemented ITR (i.e., dialysis initiation within 24 hours only in the patients with a SOFA score greater than 11) on 60-day mortality using the MIMIC-III observational database. Ninety-five percent confidence intervals are from the bootstrap. MIG=Maximal Implementation Gain. ARE=Average Rule Effect. AIE=Average Implementation Effect.

6 Discussion

Our goal was to construct an ecosystem for the evaluation of ITRs that will ultimately benefit patients. We believe that the probability model and inferential approach we introduced in this paper provide actionable tools to move this agenda forward. Below, we discuss some limitations of our approach.

In the *new ITR situation*, the exploration of the AIE relies on assuming future implementation schemes. Though sensible, the three implementation schemes we propose are subjective. Other realistic implementation schemes can be assumed and readily implemented in our methodology.

In the *partially implemented ITR situation*, inference relies on assuming a new probabilistic model. This model is largely inspired by the Neyman-Rubin causal model. As in the original model, our model requires assuming that the effect of an ITR is mediated only by the treatment prescribed by physicians to their patients. This consistency assumption may not hold in some specific settings. In cases where the decision to implement the ITR is taken by the patients—not the physicians—, it is possible that the mere "human-ITR interaction" affects outcomes. For instance, if an ITR recommends a patient treatment A, this patient may choose treatment B and compensate for not implementing the ITR by taking another effective treatment say C. This indirect effect of the ITR through treatment C would not be accounted for in our framework. With respect to the exchangeability assumption, our methodology relies on expert knowledge of the variables causing ITR implementation. For some research questions, such expert knowledge may be lacking thereby, limiting the usefulness of our approach. However, we do not believe that any statistical method can provide helpful workarounds under such conditions. Finally, as in the conventional average treatment effect, one may be tempted to estimate the prognostic function q_0 directly, rather than estimating the propensity score π from a mixture model. As $Y_i^{s=0}$ are not observed, this would require to posit a hierarchical mixture of experts model. Though compelling at first glance, this approach may be impractical as hierarchical mixture of experts were shown to have likelihoods with arbitrary bad local maxima yielding EM algorithms sensible to initialization conditions. In contrast, for mixtures of two experts, the EM algorithm with a random initialization is believed to be successful with high probability (Jin et al. 2016). Notwithstanding the empirical evidence, to our knowledge, there is no formal proof of this last statement. In a related vein, our simulations suggest that fitting a mixture of two experts with an EM algorithm, the resulting ARE and AIE estimators are unbiased and consistent with bootstrap confidence intervals achieving nominal coverage. We note that $\widehat{\Delta}_{ME}(r)$ and $\widehat{\Lambda}_{ME}(r,\rho)$ are M-estimators. Thus, assuming models for $\mu_0(\cdot)$, $\mu_1(\cdot), \tau(\cdot)$, and $\pi(\cdot)$ are correctly specified with parameters estimated for instance via maximum likelihood, $\widehat{\Delta}_{ME}(r)$ and $\widehat{\Lambda}_{ME}(r, \rho)$ have unbiased estimating equations. Appealing to M-estimation theory, one could derive these estimators' consistency, asymptotic normality,

and variance. We hope to study mixture of experts estimators' theoretical properties and explicit variance in future work.

Other future directions for our work may include the adaptation of algorithm 1 to nonparametric networks and the extension of the proposed framework to non-binary treatments and dynamic ITRs.

Acknowledgements

We thank Dr. Viet Thi Tran and Prof. Stéphane Gaudry for their insightful comments on epidemiological and critical care applications.

References

- Basu, S., Sussman, J. B., Rigdon, J., Steimle, L., Denton, B. T. & Hayward, R. A. (2017), 'Benefit and harm of intensive blood pressure treatment: Derivation and validation of risk models using data from the SPRINT and ACCORD trials', *PLoS Medicine* 14(10), e1002410.
- Díaz, I. & van der Laan, M. J. (2013), 'Assessing the causal effect of policies: an example using stochastic interventions', *The International Journal of Biostatistics* **9**(2), 161–174.
- Grolleau, F., Porcher, R., Barbar, S., Hajage, D., Bourredjem, A., Quenot, J.-P., Dreyfuss,
 D. & Gaudry, S. (2022), 'Personalization of renal replacement therapy initiation: a secondary analysis of the AKIKI and IDEAL-ICU trials', *Critical Care* 26(1), 64.
- Imai, K. & Li, M. L. (2021), 'Experimental Evaluation of Individualized Treatment Rules', Journal of the American Statistical Association 0(0), 1–15. DOI: 10.1080/01621459.2021.1923511.

Jacob, D. (2021), 'CATE meets ML', *Digital Finance* **3**(2), 99–148.

- Janes, H., Brown, M. D., Pepe, M. S. & Huang, Y. (2014), 'An Approach to Evaluating and Comparing Biomarkers for Patient Treatment Selection', *The International Journal* of Biostatistics 10(1), 99–121.
- Jin, C., Zhang, Y., Balakrishnan, S., Wainwright, M. J. & Jordan, M. I. (2016), Local maxima in the likelihood of Gaussian mixture models: structural results and algorithmic consequences, *in* 'Proceedings of the 30th International Conference on Neural Information Processing Systems', NIPS'16, Curran Associates Inc., Red Hook, NY, USA, pp. 4123– 4131.
- Johnson, A. E., Pollard, T. J., Shen, L., Lehman, L.-W. H., Feng, M., Ghassemi, M., Moody, B., Szolovits, P., Celi, L. A., Mark, R. G. & et al. (2016), 'MIMIC-III, a freely accessible critical care database', *Scientific Data* 3(1).
- Jordan, M. I. & Jacobs, R. A. (1994), 'Hierarchical Mixtures of Experts and the EM Algorithm', *Neural Computation* **6**(2), 181–214.
- Kallus, N. (2018), Balanced policy evaluation and learning, in 'Proceedings of the 32nd International Conference on Neural Information Processing Systems', NIPS'18, Curran Associates Inc., Red Hook, NY, USA, pp. 8909–8920.
- Luedtke, A. R. & van der Laan, M. J. (2016), 'Statistical inference for the mean outcome under a possibly non-unique optimal treatment strategy', *The Annals of Statistics* 44(2), 713–742. Publisher: Institute of Mathematical Statistics.
- Lunceford, J. K. & Davidian, M. (2004), 'Stratification and weighting via the propensity score in estimation of causal treatment effects: a comparative study', *Statistics in Medicine* 23(19), 2937–2960.

- Neyman, J. (1923), 'On the application of probability theory to agricultural experiments.
 essay on principles. section 9 (translation published in 1990)', *Statistical Science* 5, 472–480. Publisher: Institute of Mathematical Statistics.
- Qian, M. & Murphy, S. A. (2011), 'Performance guarantees for individualized treatment rules', *The Annals of Statistics* **39**(2), 1180–1210. Publisher: Institute of Mathematical Statistics.
- Rubin, D. B. (1974), 'Estimating causal effects of treatments in randomized and nonrandomized studies', *Journal of Educational Psychology* 66(5), 688–701. Place: US Publisher: American Psychological Association.
- Sutton, R. T., Pincock, D., Baumgart, D. C., Sadowski, D. C., Fedorak, R. N. & Kroeker,
 K. I. (2020), 'An overview of clinical decision support systems: benefits, risks, and
 strategies for success', NPJ digital medicine 3, 17.
- Takahashi, K., Serruys, P. W., Fuster, V., Farkouh, M. E., Spertus, J. A., Cohen, D. J., Park, S.-J., Park, D.-W., Ahn, J.-M., Kappetein, A. P., Head, S. J., Thuijs, D. J., Onuma, Y., Kent, D. M., Steyerberg, E. W. & van Klaveren, D. (2020), 'Redevelopment and validation of the SYNTAX score II to individualise decision making between percutaneous and surgical revascularisation in patients with complex coronary artery disease: secondary analysis of the multicentre randomised controlled SYNTAXES trial with external cohort validation', *The Lancet* 396(10260), 1399–1412.
- Tannock, I. F. & Hickman, J. A. (2016), 'Limits to Personalized Cancer Medicine', The New England Journal of Medicine 375(13), 1289–1294.
- Thomas, P. S. & Brunskill, E. (2016), Data-efficient off-policy policy evaluation for reinforcement learning, *in* 'Proceedings of the 33rd International Conference on International

Conference on Machine Learning - Volume 48', ICML'16, JMLR.org, New York, NY, USA, pp. 2139–2148.

- Tsiatis, A. A., Davidian, M., Holloway, S. T. & Laber, E. B. (2019), Dynamic Treatment Regimes: Statistical Methods for Precision Medicine, CRC Press.
- Vincent, J. L., Moreno, R., Takala, J., Willatts, S., Mendonça, A. D., Bruining, H., Reinhart, C. K., Suter, P. M. & Thijs, L. G. (1996), 'The SOFA (Sepsis-related Organ Failure Assessment) score to describe organ dysfunction/failure', *Intensive Care Medicine* 22(7), 707–710.
- White, I. R., Royston, P. & Wood, A. M. (2011), 'Multiple imputation using chained equations: Issues and guidance for practice', *Statistics in Medicine* **30**(4), 377–399.
- Xu, L. & Jordan, M. I. (1993), 'EM Learning on A Generalized Finite Mixture for Combining Multiple Classifiers', World Congress on Neural Networks 4, 227–230.
- Zhang, B., Tsiatis, A. A., Laber, E. B. & Davidian, M. (2012), 'A robust method for estimating optimal treatment regimes', *Biometrics* 68(4), 1010–1018.
- Zhao, Y., Zeng, D., Rush, A. J. & Kosorok, M. R. (2012), 'Estimating Individualized Treatment Rules Using Outcome Weighted Learning', *Journal of the American Statistical* Association 107(449), 1106–1118.

Appendix

A Proof for the ARE formula in the new ITR situation

$$\begin{aligned} \Delta(r) &= \mathbb{E} \left[Y^{s=1} - Y^{s=0} \right] \\ &= \mathbb{E} \left[r(X) Y^{a=1} + \{1 - r(X)\} Y^{a=0} - A Y^{a=1} - (1 - A) Y^{a=0} \right] \\ &= \mathbb{E} \left[\{r(X) - A\} \{Y^{a=1} - Y^{a=0}\} \right] \\ &= \mathbb{E} \left[\mathbb{E} \left[\{r(X) - A\} \{Y^{a=1} - Y^{a=0}\} | X \right] \right] \\ &= \mathbb{E} \left[\{r(X) - \pi(X)\} \tau(X) \right] \end{aligned}$$
(A1)

where equality in (A1) uses equation (5) from assumption 2 and the fact that in the new ITR situation $A = A^{s=0}$.

B Proof for the AIE/MIG formulas in the new ITR situation

$$\Lambda(r, \rho^{*}) = \mathbb{E}\left[\{Y^{*} - Y^{s=0}\}\right]$$

$$= \mathbb{E}\left[A^{*}Y^{a=1} + (1 - A^{*})Y^{a=0} - A^{s=0}Y^{a=1} - (1 - A^{s=0})Y^{a=0}\right]$$

$$= \mathbb{E}\left[\{A^{*} - A\}\{Y^{a=1} - Y^{a=0}\}\right]$$

$$= \mathbb{E}\left[\mathbb{E}\left[\{A^{*} - A\}\{Y^{a=1} - Y^{a=0}\}|X\right]\right]$$

$$= \mathbb{E}\left[\{\pi^{*}(X) - \pi(X)\}\tau(X)\right]$$
(B1)
$$= \mathbb{E}\left[\rho^{*}(X)\{r(X) - \pi(X)\}\tau(X)\right]$$
(B2)

$$= \mathbb{E}\left[\rho^*(X)\left\{r(X) - \pi(X)\right\}\tau(X)\right]$$
(B2)

where equality in (B1) uses equations (5) and (6) from assumption 2 and the fact that in

the new ITR situation $A = A^{s=0}$. The equality in (B2) follows from

$$\pi^*(X) = \mathbb{E}[A^*|X]$$

= $\mathbb{E}[S^*A^{s=1} + (1 - S^*)A^{s=0}|X]$
= $\rho^*(X)r(X) + \{1 - \rho^*(X)\}\pi(X)$

where the last equality uses the fact that $A^{s=1} = r(X)$ and equation (7) from assumption 2. The proof for the $\Gamma(r, \rho^*)$ formula follows a similar argument.

C Proof of lemma 1 in the partially implemented ITR situation

(i)
$$\pi(X) = \mathbb{E}[A|X]$$

 $= \mathbb{E}[Sr(X) + (1 - S)A^{s=0}|X]$
 $=r(X)\mathbb{E}[S|X] + \mathbb{E}[(1 - S)A^{s=0}|X]$
 $=r(X)\rho(X) + \mathbb{E}[1 - S|X]\mathbb{E}[A^{s=0}|X]$ (C1)
 $=r(X)\rho(X) + \{1 - \rho(X)\}\pi^{s=0}(X)$

Equality in (C1) relies on equation (7) from assumption 2 and the fact that in the partially implemented ITR situation $S = S^*$.

(*ii*)
$$q_0(X) = \mathbb{E}[Y^{s=0}|X]$$

 $= \mathbb{E}[A^{s=0}Y^{a=1} + (1 - A^{s=0})Y^{a=0}|X]$
 $= \mathbb{E}[A^{s=0}Y^{a=1}|X] + \mathbb{E}[(1 - A^{s=0})Y^{a=0}|X]$
 $= \mathbb{E}[A^{s=0}|X]\mathbb{E}[Y^{a=1}|X] + \mathbb{E}[(1 - A^{s=0})|X]\mathbb{E}[Y^{a=0}|X]$ (C2)
 $= \pi^{s=0}(X)\mu_1(X) + \{1 - \pi^{s=0}(X)\}\mu_0(X)$

Equality in (C2) relies on equation (5) from assumption 2.

D Inclusion and exclusion criteria from the MIMIC-III analysis

Inclusion criteria were (all needed be fulfilled):

- Admission to an intensive care unit
- Age greater than 18 years on the day of intensive care admission
- Evidence of severe acute kidney injury (stage 3 in the Kidney Disease Improving Global Outcomes classification) during the stay in intensive care
- Initiation of either mechanical ventilation or intravenous vasopressors during the stay in intensive care, prior to severe acute kidney injury

Exclusion criteria were (all needed be absent):

- End-stage renal kidney disease at intensive care admission
- Renal replacement therapy initiated prior to severe acute kidney injury
- Patients included in the study for an earlier episode of severe acute kidney injury in intensive care
- Patients expected to die within three days

E Patient inclusion diagram from the MIMIC-III database



Figure E1: Patient inclusion diagram from the MIMIC-III database. ICU=Intensive Care Unit. ESRD=End-Stage Kidney Renal Disease. AKI=Acute Kidney Injury. RRT=Renal Replacement Therapy.

F Model specification in the MIMIC-III analysis for a partially implemented ITR: dialysis initiation based on SOFA scores

For this analysis, we specify the models as follows.

• Treatment-specific prognosis functions

 $\hat{\mu}_0(X^{\mu}) = \mathbb{P}(\text{Death_at_day_60}|\text{Patient_characteristics})$

 $\sim~{\rm Age} + {\rm SOFA}$

 $\hat{\mu}_1(X^{\mu}) = \mathbb{P}(\text{Death_at_day_60}|\text{Patient_characteristics})$

 $\sim \text{Age} + \text{SOFA}$

- Propensity score in the absence of ITR implementation (expert network) $\hat{\pi}^{s=0}(X^{\pi^{s0}}) = \mathbb{P}(\text{Dialysis_initiation_within_24h}|\text{Patient_characteristics}, S = 0)$ $\sim \text{Age} + \text{Weight} + \text{BUN} + \text{pH} + \text{Potassium} + \text{SOFA}$
- Stochastic implementation function (gating network) $\hat{\rho}(X^{\rho}) \sim \text{Age} + \text{BUN} + \text{pH} + \text{Potassium}$

BUN=Blood Urea Nitrogen. SOFA=Sequential Organ Failure Assessment.

G Model specification in the MIMIC-III analysis for the new ITR: dialysis initiation based on a combination of six biomarkers

For this analysis, we specify the models as follows.

• Propensity score model

 $\hat{\pi}(X^{\pi}) = \mathbb{P}(\text{Dialysis_initiation_within_24h}|\text{Patient_characteristics})$

- $\sim {\rm Age} + {\rm Weight} + {\rm BUN} + {\rm pH} + {\rm Potassium} + {\rm SOFA} + {\rm Immunosuppressive_drug}$
- Prognosis model
 - $\hat{\mathbb{E}}[Y|X, A] = \mathbb{P}(\text{Death_at_day_60}|\text{Patient_characteristics}, \text{Dialysis_initiation_within_24h})$
 - \sim Age + Weight + BUN + pH + Potassium + SOFA

+ Dialysis_initiation_within_24h $\times \left[Age + Weight + BUN + pH + Potassium + SOFA \right]$

• ITE model

 $\hat{\tau}(X) = \hat{\mathbb{E}}[Y|X, A = 1] - \hat{\mathbb{E}}[Y|X, A = 0]$

• Prognostic function under the ITR

 $\hat{q}_1(X) = r(X)\hat{\mathbb{E}}[Y|X, A = 1] + \{1 - r(X)\}\hat{\mathbb{E}}[Y|X, A = 0]$

BUN=Blood Urea Nitrogen. SOFA=Sequential Organ Failure Assessment. ITE=Individualized Treatment Effect. ITR=Individualized Treatment Rule.