

A pragmatic approach to estimating average treatment effects from EHR data: the effect of prone positioning on mechanically ventilated COVID-19 patients

Adam Izdebski¹, Patrick J. Thorat, MD, EDIC², Robbert C.A. Lalisang, MD¹, Dean M. McHugh³, Diederik Gommers, MD, PhD⁴, Olaf L. Cremer, MD, PhD⁵, Rob J. Bosman, MD⁶, Sander Rigter, MD⁷, Evert-Jan Wils, MD, PhD⁸, Tim Frenzel, MD, PhD⁹, Dave A. Dongelmans, MD, PhD¹⁰, Remko de Jong, MD¹¹, Marco A.A. Peters, MD¹², Marlijn J.A. Kamps, MD¹³, Dharmanand Ramnarain, MD¹⁴, Ralph Nowitzky, MD¹⁵, Fleur G.C.A. Nooteboom, MD¹⁶, Wouter de Ruijter, MD, PhD¹⁷, Louise C. Urlings-Strop, MD, PhD¹⁸, Ellen G.M. Smit, MD¹⁹, D. Jannet Mehagnoul-Schipper, MD, PhD²⁰, Tom Dormans, MD, PhD²¹, Cornelis P.C. de Jager, MD, PhD²², Stefaan H.A. Hendriks, MD²³, Sefanja Achterberg, MD, PhD²⁴, Evelien Oostdijk, MD, PhD²⁵, Auke C. Reidinga, MD²⁶, Barbara Festen-Spanjer, MD²⁷, Gert B. Brunnekreef, MD²⁸, Alexander D. Cornet, MD, PhD²⁹, Walter van den Tempel, MD³⁰, Age D. Boelens, MD³¹, Peter Koetsier, MD³², Judith Lens, MD³³, Harald J. Faber, MD³⁴, A. Karakus, MD³⁵, Robert Entjes, MD³⁶, Paul de Jong, MD³⁷, Thijs C.D. Rettig, MD, PhD³⁸, Sesmu Arbous, MD, PhD³⁹, Lucas M. Fleuren, MD², Tariq A. Dam, MD², Michele Tonutti, MRes¹, Daan P. de Bruin¹, Paul W.G. Elbers, MD, PhD, EDIC², and Giovanni Cinà, PhD¹

¹Pacmed, Amsterdam, The Netherlands

²Department of Intensive Care Medicine, Laboratory for Critical Care Computational Intelligence, Amsterdam Medical Data Science, Amsterdam UMC, Vrije Universiteit, Amsterdam, The Netherlands

³Institute of Logic, Language, and Computation, University of Amsterdam, The Netherlands

⁴Department of Intensive Care, Erasmus Medical Center, Rotterdam, The Netherlands

⁵Intensive Care, UMC Utrecht, Utrecht, The Netherlands

⁶ICU, OLVG, Amsterdam, The Netherlands

⁷Department of Anesthesiology and Intensive Care, St. Antonius Hospital, Nieuwegein, The Netherlands

⁸Department of Intensive Care, Franciscus Gasthuis & Vlietland, Rotterdam, The Netherlands

⁹Department of Intensive Care Medicine, Radboud University Medical Center, Nijmegen, The Netherlands

¹⁰Department of Intensive Care Medicine, Amsterdam UMC, Amsterdam, The Netherlands

¹¹Intensive Care, Bovenij Ziekenhuis, Amsterdam, The Netherlands

¹²Intensive Care, Canisius Wilhelmina Ziekenhuis, Nijmegen, The Netherlands

¹³Intensive Care, Catharina Ziekenhuis Eindhoven, Eindhoven, The Netherlands

¹⁴Department of Intensive Care, ETZ Tilburg, Tilburg, The Netherlands

¹⁵Intensive Care, HagaZiekenhuis, Den Haag, The Netherlands

- ¹⁶Intensive Care, Laurentius Ziekenhuis, Roermond, The Netherlands
- ¹⁷Department of Intensive Care Medicine, Northwest Clinics, Alkmaar, The Netherlands
- ¹⁸Intensive Care, Reinier de Graaf Gasthuis, Delft, The Netherlands
- ¹⁹Intensive Care, Spaarne Gasthuis, Haarlem en Hoofddorp, The Netherlands
- ²⁰Intensive Care, VieCuri Medisch Centrum, Venlo, The Netherlands
- ²¹Intensive care, Zuyderland MC, Heerlen, The Netherlands
- ²²Department of Intensive Care, Jeroen Bosch Ziekenhuis, Den Bosch, The Netherlands
- ²³Intensive Care, Albert Schweitzerziekenhuis, Dordrecht, The Netherlands
- ²⁴ICU, Haaglanden Medisch Centrum, Den Haag, The Netherlands
- ²⁵ICU, Maasstad Ziekenhuis Rotterdam, Rotterdam, The Netherlands
- ²⁶ICU, SEH, BWC, Martiniziekenhuis, Groningen, The Netherlands
- ²⁷Intensive Care, Ziekenhuis Gelderse Vallei, Ede, The Netherlands
- ²⁸Department of Intensive Care, Ziekenhuisgroep Twente, Almelo, The Netherlands
- ²⁹FRCP, Department of Intensive Care, Medisch Spectrum Twente, Enschede, The Netherlands
- ³⁰Department of Intensive Care, Ikazia Ziekenhuis Rotterdam, Rotterdam, The Netherlands
- ³¹Anesthesiology, Antonius Ziekenhuis Sneek, Sneek, The Netherlands
- ³²Intensive Care, Medisch Centrum Leeuwarden, Leeuwarden, The Netherlands
- ³³ICU, IJsselland Ziekenhuis, Capelle aan den IJssel, The Netherlands
- ³⁴ICU, WZA, Assen, The Netherlands
- ³⁵Department of Intensive Care, Diaconessenhuis Hospital, Utrecht, The Netherlands
- ³⁶Department of Intensive Care, Adrz, Goes, The Netherlands
- ³⁷Department of Anesthesia and Intensive Care, Slingeland Ziekenhuis, Doetinchem, The Netherlands
- ³⁸Department of Anesthesiology, Intensive Care and Pain Medicine, Amphia Ziekenhuis, Breda, The Netherlands
- ³⁹Intensivist, LUMC, Leiden, The Netherlands

6th December 2021

Abstract

Despite the recent progress in the field of causal inference, to date there is no agreed upon methodology to glean treatment effect estimation from observational data. The consequence on clinical practice is that, when lacking results from a randomized trial, medical personnel is left without guidance on what seems to be effective in a real-world scenario.

This article proposes a pragmatic methodology to obtain preliminary but robust estimation of treatment effect from observational studies, to provide front-line clinicians with a degree of confidence in their treatment strategy. Our study design is applied to an open problem, the estimation of treatment effect of the proning maneuver on COVID-19 Intensive Care patients.

Keywords: Causal Inference, EHR data, observational study, COVID-19

1 Introduction

A central problem of causal inference is estimating treatment effects from observational data. A vast body of literature has addressed the issue of finding the best method for this task, but there is no established winner [Dorie et al., 2019]. Most model comparisons are performed on synthetic data, and the race for developing new methodologies is ongoing [Bica et al., 2021].

However, clinicians in practice are often faced with the problem of not having guidance on which treatment to use. A clear example of this phenomenon is the COVID-19 pandemic, during which clinical staff had to treat patients with a novel condition for which no randomized controlled trials (RCTs) on treatment effectiveness were available [Guérin et al., 2020]. When there is no RCT data, clinicians are left to improvise based only on their clinical experience. This points to the need for a methodology to achieve preliminary yet robust conclusions from observational data, to guide clinicians’ behaviour while randomized trials are ongoing and to determine future directions for RCTs. This need is highlighted by the increasing attention that regulatory bodies in the US and Europe are placing on the use of observational evidence [Cave et al., 2019, FDA, 2019].

In this paper we propose a procedure to arrive at sensible treatment effect estimation from Electronic

Health Records (EHRs). Our proposal is twofold. First, we suggest to emulate a target RCT, whose results provide a meaningful point of reference [Hernán and Robins, 2016, Gershman et al., 2018]. This is achieved by carefully replicating as much as possible the design of the past trial, i.e. using the same inclusion criteria, the same baseline covariates, and so forth. From a heuristic standpoint, this past RCT can act as a ‘prior’, indicating how much the results obtained from EHR data diverge from what was previously known. Second, we put forward a shortlist of methods to estimate treatment effect, including well-established ones and more recent Machine Learning techniques, with different theoretical guarantees. More than finding a single best-performing model, we intend to measure the level of agreement between the candidate models.

If our procedure is successful in both respects, meaning one can successfully emulate the design of similar RCTs from the past and the array of models provide relatively uniform estimations, we argue that the conclusions can be regarded as robust enough to guide clinical practice. This approach is ‘pragmatic’ in the sense that it can provide sought-after preliminary knowledge when an RCT is not available. To showcase the usefulness and applicability of this strategy, we test it on a real-world, large-scale use case, the estimation of the effect of proning on severely hypoxic mechanically ventilated COVID-19 patients in the Intensive Care unit (ICU). The code to implement this strategy and to obtain the results described in the paper is fully available online.

1.1 Contributions

Our contributions can be summarized as follows:

1. We implement an open-source data processing pipeline that integrates unstructured and noisy data from 25 Dutch intensive care units that shared data within the Dutch Data Warehouse (DDW) collaboration [Fleuren et al., 2021].
2. We use the data to design an observation study that emulates the PROSEVA trial. This facilitates validity of the used causal inference methods [Hernán and Robins, 2016], allows for accurate treatment effect estimation [Forbes and Dahabreh, 2020, Matthews et al., 2021] and puts our estimates in a context.
3. We found that applying prone positioning improved the P/F ratio, with estimates varying

between 14.54 and 20.11 mm Hg (depending on the model) in the time window 2-8 hours after proning and between 13.53 and 15.26 mm Hg in the time window 12-24 hours after proning (see Table 5). This aligns with results reported by previous RCTs and suggests a positive effect of prone positioning on patients with COVID-19 induced ARDS.

2 Related work

Prone positioning, consisting of turning the patient from the supine to prone position, is a commonly used technique for the treatment of severely hypoxemic mechanically ventilated patients with acute respiratory distress syndrome (ARDS). The current definition of ARDS uses the P/F ratio to classify disease severity. This ratio, defined as the partial pressure of O₂ in an arterial blood sample (PaO₂) divided by the fraction of inspired oxygen (FiO₂) ($\frac{PaO_2}{FiO_2}$ in mm Hg), subdivides patients in mild, moderate, and severe subgroups, using upper limits of 300, 200, and 100 mm Hg respectively [Force et al., 2012].

2.1 RCTs on prone positioning

Meta-analysis of eight RCTs reported that prone positioning applied for at least 12 hours a day improved the P/F ratio on average by 23 mm Hg (95% CI: 12, 35) as compared to the supine (not-proned, control) group [Munshi et al., 2017]. In the PROSEVA trial, it was found that applying prone positioning for an average of 17 hours a day improved the P/F ratio on average by 15 mm Hg (95% CI: 3, 27) [Guérin et al., 2013] and therefore it is currently considered a cornerstone in the treatment of patients with ARDS and a P/F ratio of <150 mm Hg [Guérin et al., 2020].

Prone positioning has been used widely as adjuvant therapy for respiratory failure in severe COVID-19 induced pneumonia similar to ARDS. However, increasing evidence suggests that mechanisms underlying COVID-19 respiratory failure may be different from ‘classical’ ARDS due to the frequent occurrence of pulmonary thrombosis. At the time of writing, no RCT has been conducted on the effect of proning on intubated patients with COVID-19 induced ARDS; furthermore, the available observational studies were conducted only on relatively small and single-center cohorts.

2.2 Emulating a target trial

Given the existence of several RCTs on prone positioning, we design our observational study in order to facilitate a comparison with those randomized experiments, as well as to ensure the validity of causal inference assumptions. We note that by emulating a target trial high quality electronic health record data may be used to attempt to answer the clinical question of interest instead. Emulating a target trial requires specifying the protocol of a target trial including e.g. the inclusion criteria, treatment strategies, assignment procedures, outcome of interest, treatment effect of interest, together with a synchronization of treatment assignment and eligibility determination at time zero. Such emulation is consistent with a formal counterfactual theory of causality and allows for accurate treatment effect estimation [Hernán and Robins, 2016, García-Albéniz et al., 2017, Gershman et al., 2018, Forbes and Dahabreh, 2020, Matthews et al., 2021].

2.3 Estimating treatment effects from electronic health record data

Estimating treatment effects from observational data is a fundamental problem transcending disciplinary boundaries. While multiple methods and strategies to arrive at causal quantities of interests do exist [Bica et al., 2021], the task itself is claimed to be notoriously hard. A typical testing ground for treatment effect estimators is synthetic data, i.e. a simulation experiment where the ground truth of the causal effect is known, which is not the case in real-world scenarios making it hard to benchmark and compare different methods on observational studies. The quest for validating causal inference methods, when ground truth causal effects are unknown, is ongoing [Neal et al., 2020, Schuler et al., 2018]. Yet, there seems to be no absolute winner and the best methodology appears to be context-dependent. In this article we take a selection of standard as well as more recent methods and put them to the test on a topical problem, namely the estimation of the effect of proning on COVID-19 patients. We use them to arrive at clinically meaningful conclusions from unprocessed, unstructured and highly noisy EHR data. To our best knowledge our study is the first large-scale observational study investigating the effect of prone positioning on patients with COVID-19 induced ARDS providing valuable guidance for clinicians.

3 Study design

We carefully designed our data set to emulate an RCT, in order to provide the most meaningful term of comparison and ensure validity of our study. Emulating an RCT means using the same inclusion criteria, selecting the same covariates, measuring a similar outcome, and so forth. This strategy allowed us to compare ATE estimates with the literature, and we were able to achieve a reasonable level of confidence in the fulfillment of the identifiability assumptions, which are not testable from observational data (see Appendix 4 for a discussion).

We note that, since COVID-19 is a new disease, emulating an RCT does not mean replicating it. As far as we know, the treatment effect of proning for COVID-19 ARDS could be different from the one for standard ARDS, so the PROSEVA RCT should not be regarded as providing the “true” effects that the models are trying to estimate. Nonetheless, the two pathologies are still close enough that the PROSEVA trial can provide a meaningful term of comparison; indeed this similarity is why proning has been employed for COVID-19 ARDS. For this reason the PROSEVA trial can i) provide a blueprint for the kind of experiment we want to run for proning on COVID-19 patients – in terms of covariates, inclusion criteria, outcome definition, etc – and ii) serve to put the observational ATE estimates in context.

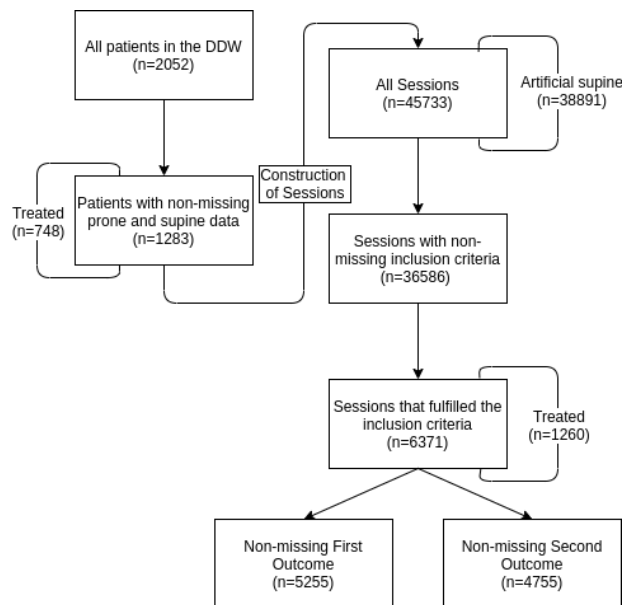


Figure 1: Overview of our study design.

3.1 Constructing Observations

The data contained information of 2052 patients recorded in the DDW who were admitted to the ICU between March and October 2020. All patients were diagnosed with COVID-19. For all patients we extracted data about turning a patient into prone and supine position. The data of each patient was sliced into prone and supine sessions, where a session was defined as a time interval throughout which the patient was in the same position, either prone or supine. This is exemplified in Figure 2. All created sessions were later selected depending on adherence to the inclusion criteria (see Section 3.3).

We took sessions to be the observations of our study, meaning that a single patient can be represented by multiple observations. A session was defined as a time interval throughout which the patient was in the same position, either prone or supine.

Sessions were employed as observations to have a data point for every moment in which a treatment decision of proning could have been made. Artificial supine sessions were created for this reason, since during a supine session there are multiple occasions in which clinicians can decide to turn the patient to prone position. An additional artificial supine session was created if, at least 8h before the end of the original supine session, the patient’s blood gas (PaO_2) and ventilator values (PEEP, FiO_2) were recorded again. The newly created artificial session ended at the same time point at which the original supine session ended. Altogether, we created 45733 sessions of which 2762 were prone, 4080 were original supine sessions and 38891 were artificial supine sessions.

3.2 Constructing Features and Outcomes

28 covariates were constructed from the data collected in the DDW. They correspond to all covariates present in the PROSEVA trial except for sepsis, McCabe and SAPS II scores, which were not available in the DDW. We included two additional covariates, one indicating morbid obesity ($\text{BMI} > 35$) and one for driving pressure, since plateau pressure measurements were missing for the majority of observations, because most patients in Dutch ICU’s were ventilated in a pressure controlled mode. The list of included covariates can be found in Table 1.

For each session we constructed covariates from the DDW by taking the last measurement in the eight hour interval before the start of the session. Since in practice there can be some delay in the recording

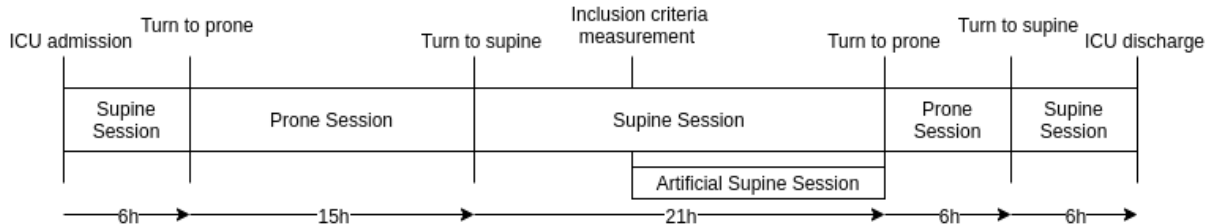


Figure 2: We describe how the data of a fictitious ICU patient would be converted into four supine and two prone sessions. An artificial supine session is created from a long supine session if there is a corresponding measurement allowing to check the inclusion criteria performed at least 8h before the end of the original session.

of such parameters, if a measurement was missing, the first measurement was taken in the 30 minutes interval following the start of the session.

We defined two outcomes. The first outcome was the last P/F ratio measurement recorded in the time interval from two hours up to eight hours from the start of the session but before the end of the session. This outcome corresponded to the measurement taken just after the prone/supine session started (‘early proning effect’). The second outcome was the last P/F ratio measurement recorded in the time interval from 12h up to 24h from the start of the session but before end of the session. Given the average length of prone sessions in this cohort (≈ 19 hours), this outcome corresponded to the measurement taken just before the session ends (‘late proning effect’). Consequently, all sessions shorter than 12h will not have a defined late proning effect.

Because of the way observations were defined, the outcomes measured the average treatment effect of a single prone session on the P/F ratio, as opposed to measuring the treatment effect of repeated proning sessions applied for a fixed time interval (e.g. three days for the PROSEVA trial). Beside the medical relevance of intermediate outcomes, this methodological choice was due to the fact that in practice patients were not prone for an identical amount of hours for a few days in a row, so the exact administering of the proning sessions from the RCT cannot be replicated. Furthermore, the aforementioned way of constructing covariates achieved perfect separation between covariates and outcomes, i.e. covariates were always measured before the outcomes, allowing for valid causal inference.

3.3 Cohort Selection and Inclusion Criteria

All sessions that did not have corresponding PaO_2 , FiO_2 and PEEP baseline measurements were discarded. Since non-invasive mechanical ventilation for COVID-19 patients was rare in the Netherlands (about 1%) [Fleuren et al., 2021], all sessions reporting the combination of PEEP and FiO_2 were considered to be sessions concerning mechanically ventilated patients. Furthermore, we discarded all proning sessions longer than 96 hours, since such sessions were extremely rare, and from a medical perspective were most likely the result of data entry errors or omissions. Following this, the inclusion criteria from the PROSEVA trial were applied for the resulting sessions. These were $\text{P/F ratio} < 150$ mm Hg, $\text{FiO}_2 \geq 60\%$ and $\text{PEEP} \geq 5$ cm of water. We chose not to include the criterion ‘tidal volume about 6 ml/kg’, since the original PROSEVA paper did not specify explicit cutoff values for tidal volume. In addition, current Dutch guidelines for mechanical ventilation of COVID-19 patients already suggest to ventilate COVID-19 patients similarly to ARDS patients with tidal volumes of about 6 ml/kg [Heunks et al., 2020]. For both outcomes, all observations with a missing outcome were discarded (18% for the first and 25% for the second outcome). Consequently, for the late proning effect we automatically discarded all sessions shorter than 12 hours, since they cannot have an outcome in the window from 12 to 24 hours. For all numerical variables we imputed data with mean values. All binary variables were imputed with the value *False*. Imputation was performed in such a way to prevent any form of data leakages between data splits. This selection resulted in 6371 observations (from 745 patients) included in our study, 1260 (from 493 patients) of which were

prone, significantly more than the PROSEVA trial that included 466 observations (each corresponding to a patient), 237 of which were prone.

3.4 Feature Characteristics and comparison to PROSEVA Trial

All the previous steps of our study design were defined in order to create a sub-population from the available EHR data that emulates the PROSEVA trial. The characteristics of all observations that fulfilled the inclusion criteria are summarized in Table 1.

Compared to the PROSEVA trial, the average COVID-19 patient was slightly older, with a higher incidence of diabetes, chronic renal failure and COPD. The lower average SOFA score is in line with the observation that COVID-19 patients suffer primarily from respiratory failure but not from septic shock or multiple organ dysfunction syndrome that typically accompany ARDS patients. On average, observations in our study had higher PEEP at similar or lower FiO_2 . Tidal volume was slightly above the recommended 6 ml/kg PBW with lower respiratory rates compared to the PROSEVA ARDS patients.

In our study higher levels of PEEP were applied and lower P/F ratios measured during prone sessions compared to PROSEVA, which may be explained by delayed proning in overwhelmed ICUs or uncertainty of effects of proning for COVID-19 patients. The most imbalanced variables between the prone and supine sessions of our study were FiO_2 , PaO_2 , P/F ratio, PEEP, SOFA Score and the usage of medications.

The PROSEVA trial similarly reported imbalance with respect to the SOFA Score and medication usage, but not w.r.t. FiO_2 , P/F ratio, PEEP and PaO_2 .

3.5 Data Sharing

Within the boundaries of privacy laws and ethics, access to the DDW can be requested on the portal of the consortium: www.icudata.nl. The code used to process the data, implement the causal inference models and perform the experiments is available online.

Table 1: Characteristics of observations that fulfilled the inclusion criteria. We report mean values with standard deviation for numerical variables and frequency for binary variables, ignoring missing values.

CHARACTERISTIC	PRONE SESSIONS	SUPINE SESSIONS
AGE - YR.	63.6 ± 10.6	63.5 ± 10.3
MALE SEX - %	73%	77%
BMI	28.3 ± 5.2	28.8 ± 5.6
SOFA SCORE	7.5 ± 3.2	8.0 ± 3.4
DIABETES - %	25%	25%
RENAL FAILURE - %	19%	20%
HEPATIC DISEASE - %	5%	6%
CORONARY ARTERY DISEASE - %	5%	7%
CANCER - %	8%	9%
COPD - %	19%	16%
IMMUNODEFICIENCY - %	14%	15%
MORBID OBESITY - %	9%	13%
VASOPRESSORS - %	64%	56%
NEUROMUSCULAR BLOCKERS - %	46%	25%
RENAL-REPLACEMENT THERAPY - %	7%	11%
GLUCOCORTICOIDS - %	5%	10%
TIDAL VOLUME - ML	454.7 ± 124.0	473.3 ± 134.4
TIDAL VOLUME - ML KG OF PBW	6.7 ± 1.8	6.9 ± 1.9
RESPIRATORY RATE	23.2 ± 8.7	24.3 ± 9.2
PEEP (CM H ₂ O)	13.1 ± 3.6	12.4 ± 3.7
FiO_2 - %	79.0 ± 14.8	71.0 ± 12.5
PLATEAU PRESSURE - CM H ₂ O	25.8 ± 6.6	26.7 ± 6.7
DRIVING PRESSURE - CM H ₂ O	13.8 ± 5.6	14.0 ± 5.5
PaO_2 - MM HG	69.4 ± 12.8	71.5 ± 11.7
P/F RATIO - MM HG	91.0 ± 22.9	103.3 ± 22.0
PaCO_2 - MM HG	56.7 ± 15.2	56.6 ± 16.4
ARTERIAL PH	7.3 ± 0.1	7.4 ± 0.1
LUNG COMPLIANCE STATIC	42.7 ± 31.0	43.4 ± 32.0

4 Methods

We frame our problem of estimating the average treatment effect of proning using the potential outcome framework [Rubin, 2005]. Let X be a vector of observed *covariates (features)* and T a binary *treatment* variable, where $T = 1$ means treated. In the potential outcomes framework we introduce two new random variables Y^1, Y^0 called *potential outcomes*. We call the observed potential outcome the *factual outcome* Y and the unobserved one the *counterfactual outcome* Y^* . We assume that the data is generated from a fixed and unknown distribution $p(X, T, (Y^1, Y^0))$ such that the standard *identifiability* criteria hold, namely consistency, positivity and ignorability. In our observational setting, for each sample we observe only the factual outcome, hence the data is of the form $\{(X_i, T_i, Y_i)\}_{i=1}^n$.

We define the average treatment effect as $ATE = \mathbb{E}[Y^1 - Y^0]$ and the conditional average treatment effect (CATE) as $\tau(x) = \mathbb{E}[Y^1 - Y^0 \mid X = x]$. Under the identifiability criteria we can treat observational data as a RCT on covariate values and estimate CATE from observational data as

$$\tau(x) = \mathbb{E}[Y \mid X = x, T = 1] - \mathbb{E}[Y \mid X = x, T = 0].$$

4.1 Model evaluation

In the task of estimating the average treatment effect of proning we are interested in assessing how close our estimate \widehat{ATE} is to the true effect ATE, as measured by the error $\epsilon_{ATE} = |ATE - \widehat{ATE}|$. Similarly for the conditional effect $\epsilon_{CATE}(\hat{\tau}) = \mathbb{E}_x[L(\tau(x), \hat{\tau}(x))]$, where L is an arbitrary loss function, $\hat{\tau}$ is the estimated conditional effect and τ is the true conditional effect. In our setting we never observe the true effect, hence we cannot use standard causal inference evaluation metrics ϵ_{ATE} and ϵ_{CATE} used on synthetic and semi-synthetic data sets like IHDP [Hill, 2011]. There is no widely established way to evaluate causal inference models on real-world observational data. We can however leverage two other markers to understand whether models are providing sensible estimates.

First, we compare ATEs estimated by our models to the ATE estimated by a randomized control trial that our study emulates. By emulating the RCT’s design and by ensuring that the identifiability criteria are satisfied in the observational study, the estimated ATEs should be comparable to the RCT’s estimate. Second, as the goal of causal inference methods is not only to provide an estimate close to

the RCT’s estimate but also to accurately predict the outcomes, we additionally evaluate the predictive performance by calculating the root-mean-squared-error (RMSE) of each method. For each model f , we calculate RMSE for factual outcomes y_i hence $RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (f(x_i, t_i) - y_i)^2}$.

4.2 Models

In order to perform treatment effect estimation, we selected a suite of standard and well-known causal inference models that rely on a diverse set of modeling assumptions. We employed the following models to estimate the treatment effect of proning on oxygenation:

- Linear Regression (LR)
- Propensity score model: Doubly Robust Inverse Probability Weighting (DR-IPW)
- Propensity score model: Blocking
- Bayesian Additive Regression Trees (BART)
- Treatment-Agnostic Representation Network (TAR-Net)
- Counterfactual Regression (CfR)

For details of the models as well as justification for their choice we refer the reader to the Appendix C. Clarifications on the notation can be found in Appendix B, while the specification on the models’ hyperparameters are reported in Appendix D.

5 Results

The goal of the study was to estimate the average treatment effect (ATE) of prone positioning on P/F ratio as measured after a session starts (early proning effect) and just before an average prone session ends (late proning effect). For each outcome the following procedure was instated. To estimate the ATE we performed a 80/20 train/test split. Each model was trained on 100 bootstrap re-samples of the train set (we repeatedly sampled 95% of the train set with replacements) and with each model an ATE estimate was calculated on the test set. We obtained the ATE estimate (Table 5) by taking the average of ATE estimates for every bootstrap sample of the train set. The 95% confidence intervals are calculated as bootstrap percentile intervals (as in Section 8 of Wasserman [2013]). Prior to calculating ATE estimates, a

hyper-parameter search was performed for some of the models. This search, as well as the strategy followed for the estimation of the propensity scores, is described in Appendix D. As we could not measure the accuracy of the estimated ATE, we evaluated the predictive performance for each model with RMSE on the test set for each of the 100 different bootstrap samples and reported the average. The 95% confidence intervals were calculated in the same way as for the ATE estimates.

The aim of the experiments was to assess whether the different frameworks for causal inference examined in Appendix C provided uniform ATE estimates. Moreover, we checked whether the estimates were close to the ATE of 15 mm Hg reported by the PROSEVA trial. Comparing the models’ ATE estimates with the result of the PROSEVA trial allowed us to gain further confidence in the robustness of the models’ estimates.

Table 2: Estimated early (top) and late (bottom) proning effects. ATE and RMSE are reported with mean and 95% CIs.

Model	ATE	RMSE
LR	15.31(11.69, 19.80)	58.48(57.91, 59.11)
DR-IPW	14.54 (10.29, 19.00)	59.09 (58.15, 60.10)
Blocking	15.59 (11.44, 20.29)	60.02 (58.27, 61.71)
BART	20.11 (14.17, 27.50)	48.62 (45.12, 56.81)
TARNet	17.70 (8.80, 25.60)	51.79 (50.74, 53.53)
CfR	18.14 (9.28, 27.35)	51.82 (50.83, 53.52)

Model	ATE	RMSE
LR	14.47(10.34, 19.34)	53.88(53.41, 54.58)
DR-IPW	13.53 (8.88, 19.45)	54.14 (53.58, 54.96)
Blocking	14.87 (9.39, 19.34)	53.22 (51.88, 55.14)
BART	13.99 (6.70, 20.34)	50.44 (47.35, 55.40)
TARNet	15.13 (5.66, 25.97)	50.61 (48.80, 52.01)
CfR	15.26 (5.54, 24.94)	50.59 (48.58, 51.54)

5.1 Average Treatment Effect Estimation

The estimated mean ATEs ranged from 14.54 to 20.11 mm Hg for the early proning effect, and from 13.53 to 15.26 mm Hg for the late proning effect. Confidence intervals for both outcomes being strictly above zero indicated a positive treatment effect of proning on

P/F ratio. A summary of all estimates can be found in Table 5 and Figure 3.

5.1.1 Early Proning Effect

The unadjusted treatment effect, meaning the difference in the mean outcome in the treated and control group, was 12.89 (95% CI: 7.88, 17.22). Linear Regression (LR) and propensity score models (DR-IPW, blocking) provided the narrowest 95% CIs for the ATE while deep learning methods (TARNet, CfR) the widest ones. This was probably due to the fact that the latter models are more prone to overfitting on the different bootstrap samples. Linear Regression, DR-IPW weighting and Blocking reported the ATE estimates closest to the PROSEVA trial, with BART, TARNet and CfR giving slightly higher estimates.

The best performing model with respect to predictive performance was BART (RMSE: 48.62), slightly outperforming TARNet and CfR (RMSE: 51.72 and 51.82 respectively). The performance of the selected models was in line with experiments on synthetic data (Table 1 in Shalit et al. [2017]), where BART, CfR and TARNet reported comparable errors with respect to ATE estimation (ϵ_{ATE}) in out-of-sample testing on both IHDP and Jobs data sets. Linear Regression and propensity models were outperformed by more recent techniques with respect to predictive performance, probably due to the greater flexibility of the latter methods. This might explain also why the ATE estimates were slightly higher for the non-linear models.

5.1.2 Late Proning Effect

The unadjusted treatment effect was equal to 11.88 (95% CI: 7.88, 16.91). Compared to the early proning effect, all models provided lower ATE estimates. The estimates reported for the two outcomes differed significantly, with BART reporting the wider gap (20.11 vs. 13.99) as opposed to Linear Regression giving more similar estimates (15.31 vs. 14.47). Compared to the early proning effect, all estimates were closer to the treatment effect from the PROSEVA trial.

The predictive performance of BART, TARNet and CfR remained the strongest with similar RMSE as for the early proning effect. Obtaining ATE estimates for the late proning effect seemed to be an easier task for the Linear Regression and propensity models as they reported smaller RMSE compared to the early proning effect. While non-linear models still outperformed linear ones, the difference was much

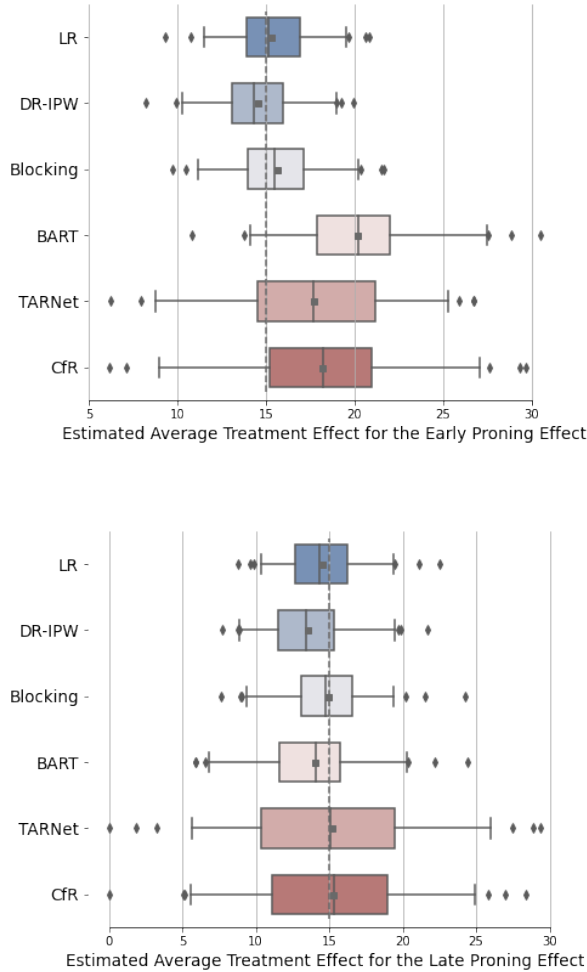


Figure 3: Distribution of the average treatment effects estimated on the test set for the early proning effect (Up) and the late proning effect (Down). Each boxplot displays the beginning of the 95% CI, the first quartile, the median, the third quartile and the end of the 95% CI of all obtained ATE estimates. The black square indicates the sample mean, which is the final ATE estimate reported in Table 5. The vertical, dotted line indicates the ATE from PROSEVA trial (=15).

smaller in RMSE, indicating that perhaps non-linear effects played less of a role for this outcome.

Models of a similar kind, both Linear Regression and propensity models, and CFR and TARNet reported comparable ATE estimates and predictive performances. This was somewhat surprising, since one would expect that the correction that propensity score models and CFR provide on top of Linear Regression and TARNet (respectively) would improve their performance and possibly modify their ATE. One of the possible explanations is that our data was rather balanced, with high treatment assignment variation. As there was no observable metric to measure the accuracy of CATE estimates and RCTs do not provide CATE estimates to compare with, it remains an open question to what extent these models could be used for CATE estimation.

6 Discussion

In this article we addressed an urgent problem, namely the lack of guidance for clinicians when RCT results are not available, by describing an approach to obtain information on treatment effectiveness from EHR data. We suggested that a multi-viewpoint approach, where different models are tested in parallel, coupled with an RCT-like experimental design can be employed to reach a reasonable level of confidence in the result of an observational study. This approach was tested on an important use case, namely treatment effect of the proning maneuver on COVID-19 patients on the ICU. Our experiments showed that our diverse set of models reached a relatively uniform conclusion. Moreover, our findings suggested that mechanically ventilated COVID-19 patients benefit from proning to a similar extent as established by the PROSEVA trial. This provides important medical knowledge in a situation when RCTs on proning for mechanically ventilated COVID-19 patients are not yet available and observational evidence is limited to single-center and relatively small cohorts [Pan et al., 2020, Perier et al., 2020, Berrill, 2021, Weiss et al., 2021, Shelhamer et al., 2021].

The extent to which results such as the ones presented here can be used to directly affect clinical decision making, or can be contrasted with RCTs, should be a matter of further discussion, as well as scrutiny from a regulatory perspective [Eichler et al., 2020]. These results cannot be considered definitive and a randomized trial is required to formally establish the effectiveness of the treatment. Repetition of these experiments

on another observational COVID-19 cohort may also help in assessing the robustness of the estimates; the fully documented and open-source code released with this article facilitates and encourages such follow-up analyses. Despite their potential brittleness and biases, observational studies such as the one presented here can help provide temporary answers for clinicians while RCTs are developed, as well as accelerate and guide prospective studies [Beaulieu-Jones et al., 2020]. Lastly, further improvement to the models’ predictive performance could warrant the investigation of the effect of proning on subgroups or of the threshold on P/F ratio for which proning is beneficial.

Limitations In conclusion, our results provide further evidence on the effectiveness of proning for the treatment of mechanically ventilated COVID-19 patients. The design of our study, implemented in fully open-source code, allows for reproducibility on other COVID-19 cohorts and provides a blueprint for treatment effect estimation in scenarios where RCT data is lacking.

Regarding the limitations that should be considered when weighing the result of this study. First, we do not compare randomized experiments and observational studies with identical designs, such as in Shadish et al. [2008], since from the DDW data it is clear that in practice the proning maneuver is not administered exactly as in the RCTs. This also holds for the outcome: the difference in treatment led us to consider oxygenation after proning, which differs from the oxygenation outcome after three days of repeated proning as defined in the PROSEVA trial. It should be noted however that RCTs already display a great deal of variability. The proning maneuver itself is performed differently in all RCTs, and the outcomes are also measured at different time points, thus the discrepancy between our study and the RCTs is arguably not greater than the one existing between RCTs. It also worth noting that the usage of sessions introduces dependence between samples from the same patient and might lead to over-representation of some individuals; sicker patients, for example, might stay longer in the ICU and generate more sessions. However, we do not register any additional imbalance in the distribution of baseline covariates, as reported in Table 1. Furthermore, we sought to mitigate this problem by employing bootstrap samples. We note that the construction of multiple data points per patients opens the possibility to study repeated or time-varying treatment by means of causal inference

methods for longitudinal data; this line of enquiry is left for future work.

Second, all models assume the absence of unmeasured confounders. To facilitate comparison with previous RCTs, we decided to limit the set of covariates to those in the study by Guérin et al. [2013]. While this expert-picked list ensures that most relevant variables are considered, hence giving some measure of reassurance that the assumption is fulfilled, it is already known that it might be relevant to add a few more covariates to have a more complete picture of the patient, e.g. fluid balance. We believe that for our case study it is possible to minimize the impact of this assumption, since the ICU routinely monitors patients in great detail. It is therefore reasonable to assume that the vast majority, if not all, of relevant parameters are recorded in EHRs, meaning that a future iteration of this study on a larger set of covariates could provide more robust results.

Third, it could be argued that a different or wider selection of models is more appropriate. The guiding principle of our choice of models was to represent both traditional methods used in the medical community as well as top performers from previous causal inference competitions, to arrive at modern Machine Learning techniques with strong theoretical guarantees.

Finally, a further point of improvement for our approach is the enhancement of the comparison between the different models. We adopted simple and well-known metrics such as RMSE, but more advanced techniques are becoming available for observational studies, such as those by Alaa and Van Der Schaar [2019], allowing for a more fine-grained analysis of the models’ differences.

References

- A. Alaa and M. Van Der Schaar. Validating causal inference models via influence functions. In *International Conference on Machine Learning*, pages 191–201. PMLR, 2019.
- B. K. Beaulieu-Jones, S. G. Finlayson, W. Yuan, R. B. Altman, I. S. Kohane, V. Prasad, and K.-H. Yu. Examining the use of real-world evidence in the regulatory process. *Clinical Pharmacology & Therapeutics*, 107(4):843–852, 2020.
- M. Berrill. Evaluation of oxygenation in 129 proning sessions in 34 mechanically ventilated covid-19 patients. *Journal of Intensive Care Medicine*, 36(2): 229–232, 2021.
- I. Bica, A. M. Alaa, C. Lambert, and M. van der Schaar. From real-world patient data to individualized treatment effects using machine learning: current and future methods to address underlying challenges. *Clinical Pharmacology & Therapeutics*, 109(1):87–100, 2021.
- G. Blattenberger and R. Fowles. Avalanche forecasting: Using bayesian additive regression trees (bart). In *Demand for Communications Services—Insights and Perspectives*, pages 211–227. Springer, 2014.
- L. Buitinck, G. Louppe, M. Blondel, F. Pedregosa, A. Mueller, O. Grisel, V. Niculae, P. Prettenhofer, A. Gramfort, J. Grobler, R. Layton, J. VanderPlas, A. Joly, B. Holt, and G. Varoquaux. API design for machine learning software: experiences from the scikit-learn project. In *ECML PKDD Workshop: Languages for Data Mining and Machine Learning*, pages 108–122, 2013.
- A. Cave, X. Kurz, and P. Arlett. Real-world data for regulatory decision making: challenges and possible solutions for europe. *Clinical Pharmacology & Therapeutics*, 106(1):36–39, 2019.
- H. A. Chipman, E. I. George, R. E. McCulloch, et al. Bart: Bayesian additive regression trees. *The Annals of Applied Statistics*, 4(1):266–298, 2010.
- 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.
- H.-G. Eichler, F. Koenig, P. Arlett, H. Enzmann, A. Humphreys, F. Pétavy, B. Schwarzer-Daum, B. Sepodes, S. Vamvakas, and G. Rasi. Are novel, nonrandomized analytic methods fit for decision making? the need for prospective, controlled, and transparent validation. *Clinical Pharmacology & Therapeutics*, 107(4):773–779, 2020.
- FDA. Us food and drug administration. real-world evidence, 2019. URL <https://www.fda.gov/science-research/science-and-research-special-topics/real-world-evidence>.
- L. M. Fleuren, D. P. de Bruin, M. Tonutti, R. C. Lalisang, and P. W. Elbers. Large-scale icu data sharing for global collaboration: the first 1633 critically ill covid-19 patients in the dutch data warehouse. *Intensive Care Medicine*, 2021. doi: 10.1007/s00134-021-06361-x.
- S. P. Forbes and I. J. Dahabreh. Benchmarking observational analyses against randomized trials: a review of studies assessing propensity score methods. *Journal of general internal medicine*, pages 1–9, 2020.
- A. D. T. Force, V. Ranieri, G. Rubenfeld, B. Thompson, N. Ferguson, E. Caldwell, et al. Acute respiratory distress syndrome. *Jama*, 307(23):2526–2533, 2012.
- X. García-Albéniz, J. Hsu, and M. A. Hernán. The value of explicitly emulating a target trial when using real world evidence: an application to colorectal cancer screening. *European journal of epidemiology*, 32(6):495–500, 2017.
- B. Gershman, D. P. Guo, and I. J. Dahabreh. Using observational data for personalized medicine when clinical trial evidence is limited. *Fertility and Sterility*, 109(6):946–951, 2018. ISSN 0015-0282. doi: <https://doi.org/10.1016/j.fertnstert.2018.04.005>. URL <https://www.sciencedirect.com/science/article/pii/S0015028218303303>.
- C. Guérin, J. Reignier, J.-C. Richard, P. Beuret, A. Gacouin, T. Boulain, E. Mercier, M. Badet, A. Mercat, O. Baudin, et al. Prone positioning in severe acute respiratory distress syndrome. *New England Journal of Medicine*, 368(23):2159–2168, 2013.
- C. Guérin, R. K. Albert, J. Beitler, L. Gattinoni, S. Jaber, J. J. Marini, L. Munshi, L. Papazian, A. Pesenti, A. Vieillard-Baron, et al. Prone position

- in ards patients: why, when, how and for whom. *Intensive care medicine*, pages 1–12, 2020.
- M. A. Hernán and J. M. Robins. Using big data to emulate a target trial when a randomized trial is not available. *American journal of epidemiology*, 183(8):758–764, 2016.
- M. A. Hernán and J. M. Robins. *Causal Inference: What If*. Boca Raton: Chapman & Hall/CRC, 2020.
- L. M. A. Heunks, R. Endeman, and J. van der Hoeven. Mechanical ventilation in covid-19. dutch society for intensive care medicine (nvc) guideline. Oct 2020.
- J. Hill. Bayesian nonparametric modeling for causal inference. *Journal of Computational and Graphical Statistics*, 20:217–240, 03 2011.
- G. W. Imbens. Nonparametric estimation of average treatment effects under exogeneity: A review. *Review of Economics and statistics*, 86(1):4–29, 2004.
- G. W. Imbens and D. B. Rubin. *Causal inference in statistics, social, and biomedical sciences*. Cambridge University Press, 2015.
- F. D. Johansson, U. Shalit, N. Kallus, and D. Sontag. Generalization bounds and representation learning for estimation of potential outcomes and causal effects, 2020.
- J. D. Kang, J. L. Schafer, et al. Demystifying double robustness: A comparison of alternative strategies for estimating a population mean from incomplete data. *Statistical science*, 22(4):523–539, 2007.
- A. Kapelner and J. Bleich. bartMachine: Machine learning with Bayesian additive regression trees. *Journal of Statistical Software*, 70(4):1–40, 2016.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- A. A. Matthews, K. Szummer, I. J. Dahabreh, B. Lindahl, D. Erlinge, M. Fetchting, T. Jernberg, A. Berglund, and M. A. Hernán. Comparing effect estimates in randomized trials and observational studies from the same population: an application to percutaneous coronary intervention. *medRxiv*, 2021.
- L. Munshi, L. Del Sorbo, N. K. Adhikari, C. L. Hodgson, H. Wunsch, M. O. Meade, E. Uleryk, J. Mancebo, A. Pesenti, V. M. Ranieri, et al. Prone position for acute respiratory distress syndrome. a systematic review and meta-analysis. *Annals of the American Thoracic Society*, 14(Supplement 4): S280–S288, 2017.
- B. Neal, C.-W. Huang, and S. Raghupathi. Realcause: Realistic causal inference benchmarking. *arXiv preprint arXiv:2011.15007*, 2020.
- C. Pan, L. Chen, C. Lu, W. Zhang, J.-A. Xia, M. C. Sklar, B. Du, L. Brochard, and H. Qiu. Lung recruitability in covid-19-associated acute respiratory distress syndrome: a single-center observational study. *American journal of respiratory and critical care medicine*, 201(10):1294–1297, 2020.
- F. Perier, S. Tuffet, T. Maraffi, G. Alcalá, M. Victor, A.-F. Haudebourg, N. De Prost, M. Amato, G. Carreaux, and A. Mekontso Dessap. Effect of positive end-expiratory pressure and proning on ventilation and perfusion in covid-19 acute respiratory distress syndrome. *American Journal of Respiratory and Critical Care Medicine*, 202(12):1713–1717, 2020.
- 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.
- D. B. Rubin. Causal inference using potential outcomes: Design, modeling, decisions. *Journal of the American Statistical Association*, 100(469):322–331, 2005.
- A. Schuler, M. Baiocchi, R. Tibshirani, and N. Shah. A comparison of methods for model selection when estimating individual treatment effects. *arXiv preprint arXiv:1804.05146*, 2018.
- S. Seabold and J. Perktold. statsmodels: Econometric and statistical modeling with python. In *9th Python in Science Conference*, 2010.
- W. R. Shadish, M. H. Clark, and P. M. Steiner. Can nonrandomized experiments yield accurate answers? a randomized experiment comparing random and nonrandom assignments. *Journal of the American statistical association*, 103(484):1334–1344, 2008.
- U. Shalit, F. D. Johansson, and D. Sontag. Estimating individual treatment effect: generalization bounds and algorithms, 2017.

- M. C. Shelhamer, P. D. Wesson, I. L. Solari, D. L. Jensen, W. A. Steele, V. G. Dimitrov, J. D. Kelly, S. Aziz, V. P. Gutierrez, E. Vittinghoff, et al. Prone positioning in moderate to severe acute respiratory distress syndrome due to covid-19: A cohort study and analysis of physiology. *Journal of Intensive Care Medicine*, 36(2):241–252, 2021.
- C. Villani. *Optimal transport: old and new*, volume 338. Springer Science & Business Media, 2008.
- L. Wasserman. *All of statistics: a concise course in statistical inference*. Springer Science & Business Media, 2013.
- T. T. Weiss, F. Cerda, J. B. Scott, R. Kaur, S. Sungurlu, S. H. Mirza, A. A. Alolaiwat, A. E. Augustynovich, and J. Li. Prone positioning for patients intubated for severe acute respiratory distress syndrome (ards) secondary to covid-19: a retrospective observational cohort study. *British journal of anaesthesia*, 126(1):48–55, 2021.
- Q. Zhou and J. S. Liu. Extracting sequence features to predict protein–dna interactions: a comparative study. *Nucleic acids research*, 36(12):4137–4148, 2008.

A Dutch Data Warehouse

The DDW was constructed by combining, unifying, and structuring EHR data from 25 different ICUs in the Netherlands. A dedicated team of ICU clinicians, data scientists, and IT staff put great effort in the medical validation of the data by assessing data quality at several stages in the extensive data validation process [Fleuren et al., 2021]. Since raw EHR data contains all routinely and frequently captured clinical information, this whole endeavour eventuated in a unique dataset containing highly granular data of COVID-19 patients throughout their entire ICU stay.

The DDW includes data on demographics, comorbidities, monitoring and life support devices, laboratory results, clinical observations, medications, fluid balance, and outcomes. While the DDW is continuously growing, it contained over 400 million data points from 2052 COVID-19 patients at the moment of accessing.

B Notation

Throughout the text we associate each observation with an index $i, j, k \in \mathbb{N}$. The covariate (feature) space $\mathcal{X} \subseteq \mathbb{R}^d$ and outcome space $\mathcal{Y} \subseteq \mathbb{R}$ are defined as usual. We denote random variables with capital letters X, Y, Z . We define the covariate vector X , outcome vector Y and treatment indicator T and write X_i, Y_i, T_i , if they correspond to an observation i . We denote realizations of random variables with corresponding lower case letters x, y, z . If a realization of a random variable X is associated with the observation i , then we add a subscript x_i . When not necessary, we omit the corresponding subscripts. When a random variable X follows a distribution p , we will write $\mathbb{E}[X]$ instead of $\mathbb{E}_{x \sim p(x)}[x]$ to denote the expected value.

C Models

We choose models that rely on a diverse set of modeling assumptions, with the goal of estimating the average treatment effect of prone positioning. We introduce traditional methods for causal inference: linear regression and propensity score based models in Section C.0.1 and Section C.0.2, the popular Bayesian regression model BART in Section C.0.3 and a recent deep-learning based framework in Section C.0.4. The way we choose hyper-parameters for each of the models, as well as details on implementation, is described in Appendix D.

C.0.1 Linear Regression

Linear regression (LR) is the simplest parametric approach to estimating average treatment effect from observational data [Hernán and Robins, 2020]. We add it to the shortlist of models since it is widely used in the literature, and it is a valid way to adjust for confounders under the following assumption [Hernán and Robins, 2020]. Linear regression assumes that the conditional expectation of the outcome Y can be expressed as

$$\mathbb{E}[Y \mid X = x, T = t] = \alpha_0 + \alpha_1 \cdot t + \alpha_2 \cdot x \tag{1}$$

We estimate weights in Equation 1 by fitting a linear regression model and use the coefficient of the treatment effect indicator α_1 as the estimated average treatment effect. This is because by Equation 1 it holds that

$$\begin{aligned} \tau(x) &= \mathbb{E}[Y \mid X = x, T = 1] - \mathbb{E}[Y \mid X = x, T = 0] \\ &= \alpha_1 \end{aligned}$$

Therefore, not only is the outcome Y a linear combination of the treatment T and covariates X , but also the treatment effect is constant across different values of covariates X .

C.0.2 Models Based on Propensity Score

A propensity score, introduced by Rosenbaum and Rubin [1983], is defined as the conditional probability of receiving the treatment $e(x) = \mathbb{P}(T = 1 \mid X = x)$.

In our setting the true propensity score is unknown and needs to be estimated from data using a *propensity score model*. Such model is used to correct the linear regression model in two ways: blocking and inverse probability weighting (IPW) [Imbens, 2004].

First, we use blocking (sometimes referred to as ‘stratification’) to partition the observational data into subsets (blocks) containing observations with an estimated propensity score in a given range. We estimate the treatment effect within each block by fitting a linear regression to the data within the block. We obtain the final ATE estimate by taking the average of block-wise estimates, weighted by the proportion of data points in each block.

Second, following Kang et al. [2007] we define a doubly robust estimator that combines inverse probability weighting with linear regression (DR-IPW) as follows. For each observation with covariate value x and treatment indicator t , we use the estimated propensity score $\hat{e}(x)$ to define a weight

$$w(x, t) = \left[\hat{e}(x)^t \cdot (1 - \hat{e}(x))^{(1-t)} \right]^{-1}$$

We then fit a linear regression weighted by w to our data, giving more importance to observations with an unlikely treatment assignment. The ATE estimate is the treatment indicator’s coefficient of the weighted model, as for the linear regression.

The idea behind both models is to fit a linear regression to a modified version of the original data, using the propensity scores to balance the difference between the treated and control group and to correct the fact that treatment is not independent of covariates. This can be beneficial in a variety of theoretical scenarios, for a detailed discussion of the above models and their statistical properties [Imbens, 2004, Kang et al., 2007].

C.0.3 BART

Bayesian Additive Regression Trees (BART) is a nonparametric, Bayesian method for estimating functions using regression trees [Chipman et al., 2010]. It is a popular method for causal inference [Dorie et al., 2019], used often as a baseline comparison for state-of-the-art deep learning methods. The model’s popularity is based on its outstanding performance in multiple real world scenarios [Zhou and Liu, 2008, Blattenberger and Fowles, 2014], as well as simulated ones [Dorie et al., 2019].

BART assumes that the conditional expectation of the outcome Y can be expressed as

$$\mathbb{E}[Y \mid X = x, T = t] = g(x, t) + \epsilon; \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

where g is an arbitrary function, ϵ is an error term and σ^2 is a constant. BART approximates function g with a sum $\sum_{m=1}^M \mathcal{T}_m(x, t)$ of M -many regression trees $\mathcal{T}_1, \dots, \mathcal{T}_M$. We obtain a CATE estimate by taking the difference $\hat{\tau}(x) = \sum_{m=1}^M (\mathcal{T}_m(x, 1) - \mathcal{T}_m(x, 0))$. The ATE estimate is obtained by averaging over CATE estimates computed on the test set $\widehat{ATE} = \frac{1}{n} \sum_{i=1}^n \hat{\tau}(x_i)$.

C.0.4 Counterfactual Regression

Counterfactual Regression (CfR) is a general method for estimating CATE that is based on a theoretical framework developed by Shalit et al. [2017]. CfR is a regularized minimization procedure that simultaneously tries to predict the correct observed outcomes and to find a representation that balances the treated and control group. Thus CfR fits a representation $\Phi : \mathcal{X} \rightarrow \mathcal{X}$ and a regression function $f : \mathcal{X} \times \{0, 1\} \rightarrow \mathcal{Y}$, where \mathcal{X} is the covariate space and \mathcal{Y} is the outcome space, in order to minimize the objective

$$\min_{f, \Phi: \|\Phi\|=1} \frac{1}{n} \sum_{i=1}^n w_i \cdot (f(\Phi(x_i), t_i) - y_i)^2 + \lambda \cdot \mathcal{R}(f) + \alpha \cdot \text{Wass}(\{\Phi(x_i)\}_{t_i=1}, \{\Phi(x_i)\}_{t_i=0}) \quad (2)$$

where $w_i = \frac{1}{2} \left(\frac{t_i}{u} + \frac{1-t_i}{1-u} \right)$ is a weight accounting for the treatment imbalance, $u = \sum_{i=1}^n \frac{t_i}{n}$ is the fraction of observations being treated, α and $\lambda \cdot \mathcal{R}(f)$ are arbitrary regularization terms and Wass is the Wasserstein distance defined by Villani [2008]. Theorem 3 of Johansson et al. [2020] guarantees that by minimizing

Equation 2 we obtain f, Φ such that $\hat{\tau}(x) = f(\Phi(x), 1) - f(\Phi(x), 0)$ is a consistent CATE estimator. We also employ a special case of CfR, called TARNet, whose objective is given by Equation 2 with $\alpha = 0$. We follow CfR’s and TARNet’s implementation proposed by Shalit et al. [2017]. We use CfR and TARNet to obtain the ATE estimate by averaging over CATE estimates computed on the test set $\widehat{ATE} = \frac{1}{n} \sum_{i=1}^n \hat{\tau}(x)$.

D Algorithmic Details

D.0.1 Propensity Score Models

We evaluate three different strategies to estimate the propensity score. We fit a logistic regression model to 1) confounders chosen by ICU doctors, 2) all covariates, 3) all covariates with interactions. The second model is preferable to the first as adjusting with the second model results in lower imbalance, when measured by the absolute value of the normalized mean difference defined by Imbens and Rubin [2015]. We choose the second model over the third, despite a slightly lower imbalance, because the propensity scores estimated by the third model tend to extreme values, making it less useful for calculating IPW weights. We use a scikit-learn implementation for the logistic regression with no regularization and with weights adjusting for class imbalance [Buitinck et al., 2013]. After estimating the propensity scores (Figure 4) we perform propensity score clipping ($= 0.1$) to account for the lack of overlap in the region with values below 0.1. For the linear regression models we use the implementation proposed in *statsmodels* package [Seabold and Perktold, 2010].

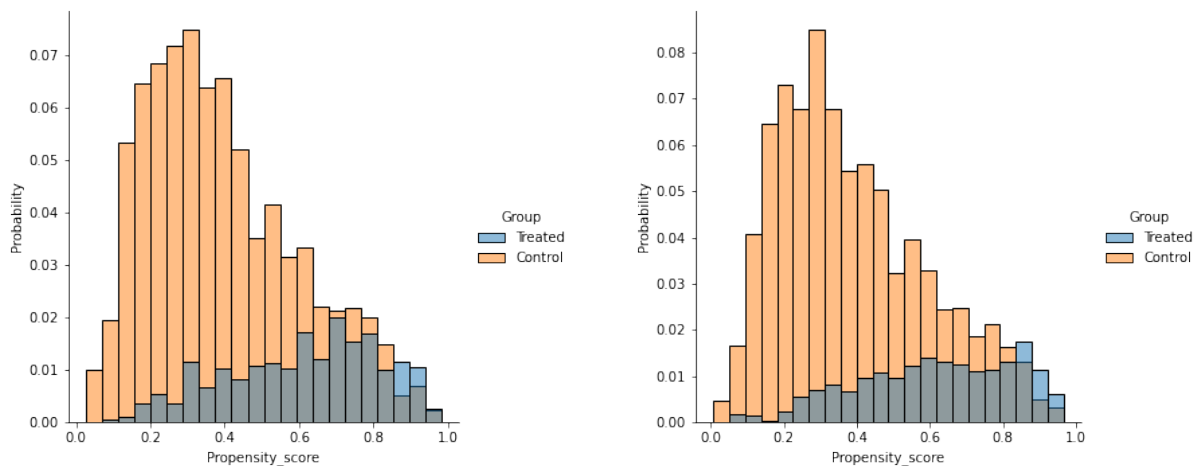


Figure 4: Distribution of the estimated propensity scores for the data set corresponding to the early proning effect (Up) and the late proning effect (Down). The plots show a significant overlap for the estimated propensity scores above 0.1.

D.0.2 BART

We use BART as implemented by Kapelner and Bleich [2016]. We train BART with $m = 50$ trees and use 5-fold CV (with respect to the lowest RMSE) on the train data to choose hyperparameters: $k = 2$, $\nu = 3$, $q = 0.9$ controlling the model regularization. These hyperparameters are reported as default by Chipman et al. [2010]. Additionally, we use default parameters $\alpha = 0.95$, $\beta = 2$ specifying the prior over trees’ structure.

D.0.3 CFR and TARNet

We follow CFR's and TARNet's implementation proposed by Shalit et al. [2017]. They implement CFR and TARNet as a feed-forward neural network with three fully-connected layers and with ELU activation functions for both the representation Φ and for the regression functions $f_1 = f(x, 1)$ and $f_0 = f(x, 0)$. We use default hyperparameters with layer size of 200 for the representation and 100 for outcome networks. The model is trained using Adam optimizer and ℓ_2 weight decay for outcome networks [Kingma and Ba, 2014]. We choose default regularization hyper-parameters: $\lambda = 0.0001$ and $\alpha = 1$ ($\alpha = 0$ for TARNet). We perform early stopping w.r.t. surrogate mean-squared-error defined by Shalit et al. [2017] and calculated on a held-out validation set. We perform a grid search evaluated on the validation set in order to assess whether to use different than default hyperparameters for TARNet and CFR. We decided to use the default hyperparameters for all models as improvement w.r.t. their predictive performance was negligible.

E Medical Abbreviations