

1. Background

- Mechanical ventilation with positive end-expiratory pressure (PEEP) is an essential intervention in the management of critically ill patients in the ICU.
- Determining the optimal PEEP level remains a challenge due to conflicting evidence from clinical studies.
- Our research leverages machine learning methods that take confounding into account to estimate individualized treatment effects (ITE) of high PEEP on survival outcomes.

MIMIC-IV Dataset:

We will utilize the MIMIC-IV [1] dataset, a comprehensive collection of ICU patient data. As an observational dataset, MIMIC-IV contains inherent confounding variables, which must be carefully addressed to ensure accurate analysis and interpretation.

- Confounding variables are variables that affect both the choice of treatment and the outcome, and therefore distort the true effect of the treatment.

2. Research Question

How can the DR-learner, a machine learning-based method, be used to predict survival outcomes in ICU patients under different PEEP regimes based on individual characteristics, and how does this method compare to other CATE estimators when evaluated on an RCT dataset?

4. Results on simulated data

Performance comparison of the meta learners using average MSE of 10 iterations

- Gradient Boosting served as the base learner for S- and T- learners and the first stage of the DR-learner, while Linear Regression was used for the final stage of the DR-learner

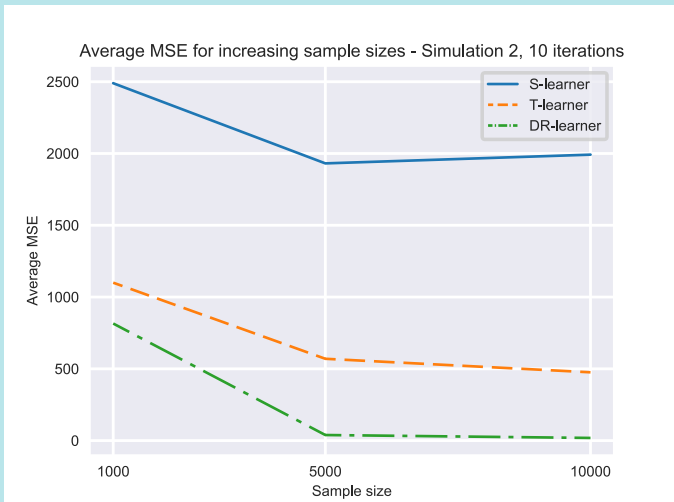


Figure 1: Simulation 2 (Complex Linear) - different linear response functions are applied across the feature space, with a propensity score of 0.5.

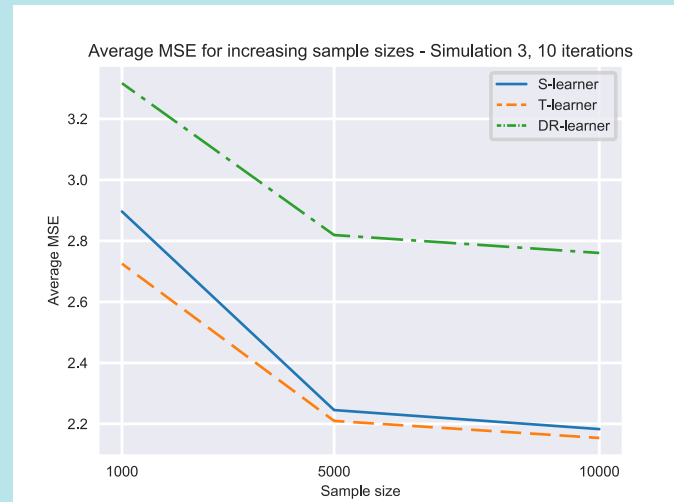


Figure 2: Simulation 3 (Complex Non-Linear) - non-linear response functions

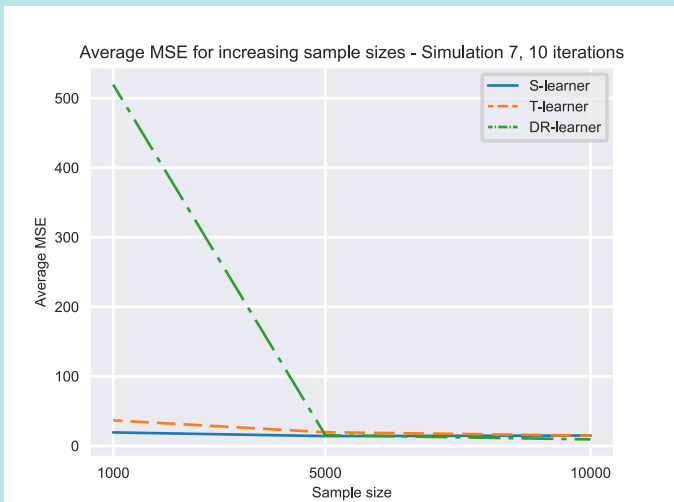


Figure 3: Simulation 7 - 12% of units receive treatment, with a simple CATE function to estimate.

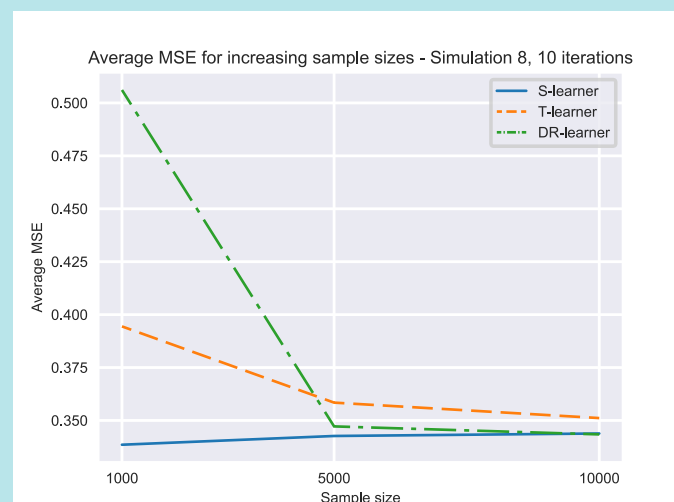


Figure 4: Simulation 8 (Beta Confounded) - Modification of Sim. 6, response functions differ and are dependent on covariates.

Simulation Results:

- S-learner: Performed well in most scenarios except different linear response functions (Figure 1).
- T-learner: Generally outperformed by other learners.
- DR-learner: Excelled in handling unbalanced and confounded data, demonstrating robustness in complex conditions (Figures 3 and 4). However, in non-linear scenarios (Figure 2) did not perform as well, likely due to its reliance on Linear Regression in the final stage

3. Methodology

The S-, T-, and (**Doubly-Robust**) DR- meta-learners [2][3] were used to estimate the Conditional Average Treatment Effect (CATE) of high PEEP on patient outcomes.

S-learner

- Uses a single model combining treatment indicator and patient features to predict outcomes.

T-learner

- Uses separate models for treated and control groups (response functions).

DR-learner

- Estimates propensity score and outcome regression.
- Generates pseudo outcomes.
- Regresses on pseudo outcomes to estimate the CATE.

Advantages of the DR-learner

- Reduced bias with two models
- Increased resilience to model errors
- Flexibility in application

Approach

- Generate 8 simulated datasets to compare the meta learners
 - Six simulations followed the methodology of Künzel et al. [2]
 - Two additional simulations for further comparison.
- Pre-process the MIMIC-IV dataset, impute missing data, and identify confounding variables
- Select underlying models
- Train and evaluate meta learners on the pre-processed MIMIC-IV dataset
- Evaluate the estimated CATE on an RCT dataset

5b. Results on MIMIC-IV and RCT

Qini curves of meta-learners with selected underlying models [4]

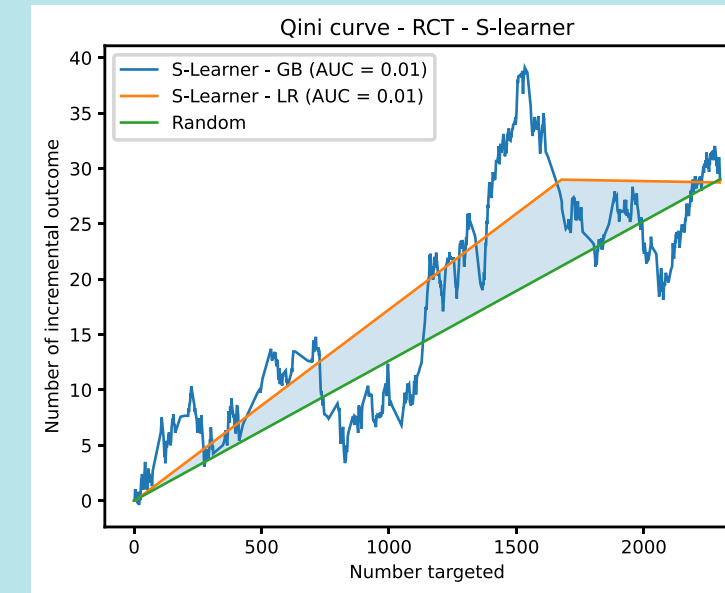


Figure 5: S-learner performance comparison with the largest area highlighted

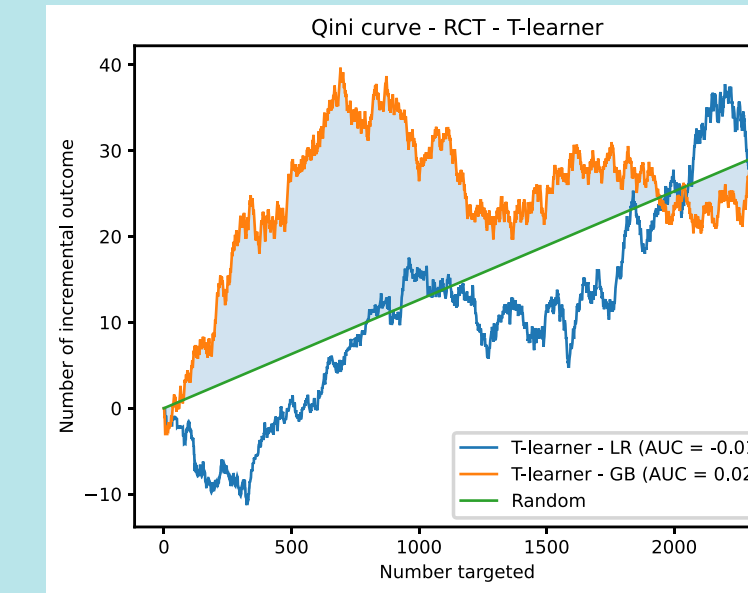
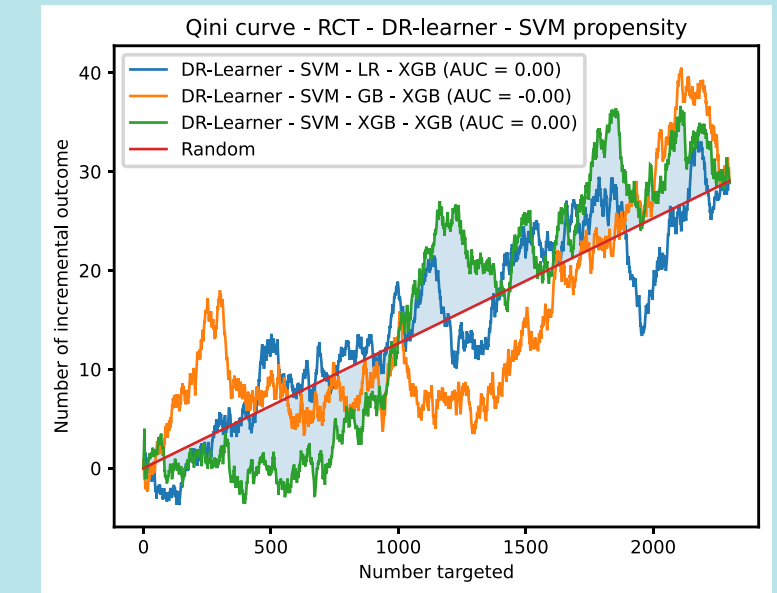
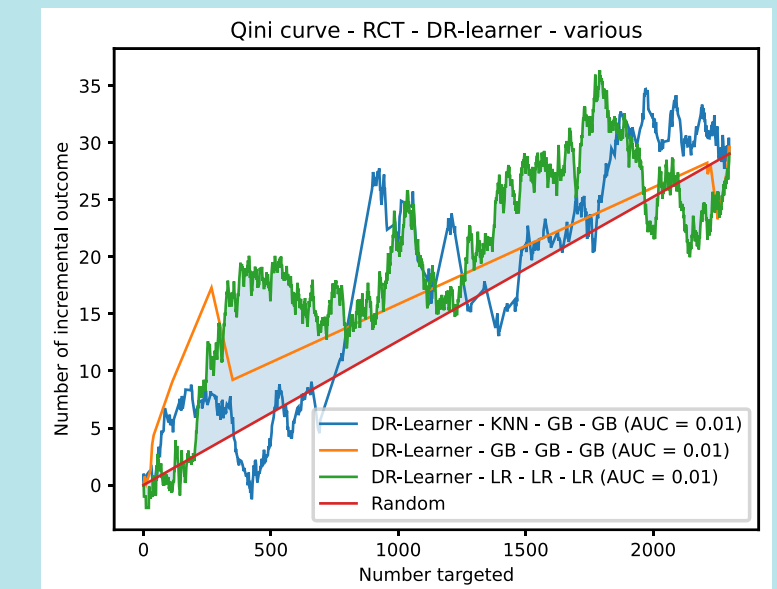


Figure 6: T-learner performance comparison with the largest area highlighted



(a): DR-learner with SVM for the propensity model



(b): The rest of the DR-learners

Figure 7: DR-learner performance comparison with the largest area highlighted

MIMIC-IV Dataset:

- Performance varied with train-test splits.
- Non-linear models tended to overfit; linear models underperformed.
- DR-learner showed potential but faced challenges with real-world data due to overfitting and data limitations.
- Difficulties to estimate the treatment effect reliably

RCT Evaluation:

- Evaluated meta-learners trained on MIMIC-IV with features available in the RCT dataset.
- Also unable to reliably predict the treatment effect
- Results were inconsistent across learners with generally low AUQC.

5a. Results on MIMIC-IV and RCT

Average area under the Qini curve (AUQC) of the trained meta learners

Learner	Model	MIMIC-IV		RCT Dataset
		Train Set	Test Set	
S	Gradient Boosting	0.091742	-0.021610	0.005332
	Linear Regression	-0.004268	-0.00602	0.007350
T	Gradient Boosting	0.432664	0.009727	0.019314
	Linear Regression	0.048424	0.016777	-0.007351
DR	SVM Propensity and XGBoost Response with XGBoost Regressor Final	0.275965	0.016025	0.001503
	SVM Propensity and Logistic Regression Response with XGBoost Regressor Final	0.281178	0.010721	0.000868
	SVM Propensity and Gradient Boosting Response with XGBoost Regressor Final	0.281015	0.019345	-0.000503
	Logistic Regression Propensity and Logistic Regression Response with Linear Regression Final	0.022966	0.025802	0.009918
	Gradient Boosting Propensity and Gradient Boosting Response with Gradient Boosting Final	0.129593	0.018670	0.006207
	K-Nearest Neighbors Propensity and Gradient Boosting Response with Gradient Boosting Final	0.263782	0.042217	0.005886

Table 1: S-, T-, and DR-learners - Average AUQC of 10 train-test splits of MIMIC-IV, and RCT AUQC Performance Comparison

6. Limitations

- Small sample size of MIMIC-IV dataset (3,941 samples)
- Computational limitations
- Potential misidentification of confounding variables
- Overfitting in non-linear models

7. Conclusion

- The DR-learner shows promise in simulation scenarios, especially with confounded and unbalanced data, but struggles with real-world data due to overfitting and data limitations.
- Both MIMIC-IV and RCT datasets demonstrated poor performance in estimating the treatment effect, indicating the inherent difficulty in this task.
- When evaluated on RCT data, the DR-learner performs similarly to the S- and T-learners.
- The challenges highlight the complexity of reliably estimating ITEs in clinical settings.

8. Future Research

- Train on datasets with more samples
- Explore advanced methods like Neural Networks.
- Fine-tune model parameters
- Reassess and validate variable selection
- Mitigate overfitting.

References

- A. E. W. Johnson et al., "MIMIC-IV, a freely accessible electronic health record dataset," Scientific Data, vol. 10, no. 1, Jan. 2023, doi: 10.1038/s41597-022-01899-x.
- S. R. Künzel, J. S. Sekhon, P. J. Bickel, and B. Yu, "Metalearners for estimating heterogeneous treatment effects using machine learning," Proceedings of the National Academy of Sciences of the United States of America, vol. 116, no. 10, pp. 4156-4165, Feb. 2019, doi: 10.1073/pnas.1804597116.
- E. H. Kennedy, "Towards optimal doubly robust estimation of heterogeneous causal effects," arXiv (Cornell University), Jan. 2020, doi: 10.48550/arxiv.2004.14497.
- F. Devriendt, T. Guns, and W. Verbeke, "Learning to rank for uplift modeling," arXiv, 2020. doi: 10.48550/ARXIV.2002.05897.