# Interpretable Semiparametric Regression For Treatment-Covariates Interaction Learning: A Dual-Score System

Muyan Jiang, Yunkai Zhang, and Anil Aswani

*Abstract*— Accurately estimating the effects of continuous treatment variables on binary outcomes is challenging. Traditional logistic regression models assume linearity and struggle with dependencies among treatment effects and patient outcomes. Such limitations hinder optimizing treatments, essential for data-driven decision-making and its various applications (e.g., precision medicine). We introduce a semiparametric regression model that blends the interpretability of parametric models with the adaptability of nonparametric approaches, featuring a dual-score system for assessing patient prognosis and determining an optimal treatment level. We demonstrate the potential of our approach by conducting numerical simulations that suggest convergence occurs, and then we apply our approach to a case study using the International Warfarin Pharmacogenomics Consortium (IWPC) dataset.

## I. INTRODUCTION

Addressing the challenge of estimating treatment effects, especially with binary outcomes and continuous treatment variables, requires models to provide accurate and explicable estimates and account for treatment-covariate interactions. While suitable for binary outcomes, traditional logistic regression techniques struggle with the causal relationship between treatment and covariates. Machine learning techniques such as neural networks can capture complex relationships but lack transparency, limiting their clinical utility. Semiparametric regression models offer a solution by blending nonparametric methods' adaptability with parametric models' interpretability, effectively addressing treatment effects on outcomes [1], [2].

Single index models exemplify this family of models [3]. The model takes the form $\mathbb{E}(Y|X) = g(\xi^T X) + \epsilon$, where $\xi^T X$ is the index, $g(\cdot)$ is an unknown function, and $\epsilon$ is the noise. It has been applied in diverse fields such as economics, environmental science, and healthcare [4], [5]. [6] used this model to study optimizing treatment rules in clinical trials. Moreover, a recent application to COVID-19 treatment [7] used it to develop the "treatment benefit index" to guide individualized treatment recommendations.

Regression models can incorporate continuous treatment variables and offer nuanced insights into treatment-response dynamics, which is crucial for optimizing treatments. Within this context, [8] introduced a novel method to fine-tune personalized dosing strategies by optimizing a local approximation of the value function through outcome-weighted learning.

[9] presented a methodology for evaluating continuous treatment policies, employing kernel-enhanced inverse probability weighting and doubly robust techniques. Additionally, [10] proposed a semiparametric change-plane model to identify and analyze subgroups that exhibit variable treatment effects.

However, many existing methods often prioritize the accuracy of estimating individualized dosing rules over quantitative interpretability [11]. To achieve high accuracy, they use complex statistical models that are elusive for practitioners, which deters their utility in clinical diagnosis [8]. Furthermore, some models cannot be extended to handle classification tasks [12].

This paper introduces a model that leverages semiparametric regression to estimate treatment effects using a partially linear model. This method considers the interaction effects between covariates and a continuous treatment variable, particularly in the setting of binary outcomes. Our methodology aims to balance the relationship between treatment and outcome while emphasizing the interpretability crucial for many applications such as clinical decision-making. An overview of the proposed framework is shown in Fig. 1.

To summarize, our contributions are threefold:

1) We develop a partially linear model to estimate the optimal individualized treatment strategies, preserving quantitative interpretability.
2) We design a dual-score system: one score reflects prognosis and the other offers a reference for an optimal treatment, informed by the interplay among covariates.
3) We apply our method to an anticoagulation study, which shows the potential of our approach for optimizing
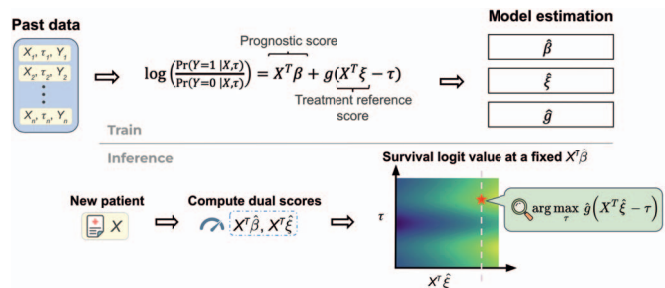


Fig. 1. Overview of our proposed pipeline. During training, the model learns $\hat{\beta}$ and $\hat{\xi}$ associated prognosis and treatment effectiveness, respectively. During inference, as the prognostic score $X^T\hat{\beta}$ is fixed conditioning on the patient's attributes, we can directly visualize outcome density using a heatmap with estimated $\hat{g}$. The optimal treatment value is then extrapolated at the maximum point along the y-axis corresponding to a fixed $X^T\hat{\xi}$.

MJ, YZ, and AA are with the Department of Industrial Engineering & Operations Research, University of California, Berkeley, CA 94720, USA {muyan_jiang, yunkai_zhang, aaswani}@berkeley.edu

treatments based upon the model.

## II. MODEL SETUP AND ESTIMATION

In many settings, a binary variable $Y \in \{0, 1\}$ denotes outcomes, $X \in \mathbb{R}^p$ represents predictive traits, and $\tau \in \mathbb{R}$ is a continuous treatment variable. Usually, $\tau$ is included as an additional covariate in logistic regression models. However, since treatment allocation often correlates significantly with other covariates, it's crucial to develop a model that captures the interaction between treatment and covariates and remains interpretable in terms of diagnostic indices, similar to traditional logistic regression.

### A. Model Formulation

Given $Y \in \{0, 1\}$, $X \in \mathbb{R}^p$, $\tau \in \mathbb{R}$, we consider the following logistic regression model:

$$Pr(Y = 1|X, \tau) = \sigma(X^T\beta + g(X^T\xi - \tau)) \qquad (1)$$

where $\sigma(\cdot) = \frac{1}{1+exp(\cdot)}$ is the logistic function, $g(\cdot) : \mathbb{R} \to \mathbb{R}$ is an unspecified function that models the effect of the treatment, and $\beta, \xi$ are learnable parameters. The first linear score $X^T\beta$ models the main linear effect based on the baseline trait of a patient. The second linear score $X^T\xi$ is a score that linearly interacts with the treatment variable, where the difference $X^T\xi - \tau$ term is the argument of an unknown function $g$ that encapsulates the covariates-treatment interaction. Furthermore, we require that $\xi \in \{\xi \in \mathbb{R}^p : \|\xi\|_2 = 1, \xi_1 \geq 0\}$ for identifiable purposes due to the unknown nature of the function $g$.

We consider the "odds ratio" (OR), a measure of association between an exposure and an outcome, and obtain the following:

$$\log \frac{Pr(Y = 1|X, \tau)}{Pr(Y = 0|X, \tau)} = X^T\beta + g(X^T\xi - \tau), \qquad (2)$$

which reduces to a partially linear model with an unknown nonlinear part.

In a practical application, the initial linear score $X^T\beta$ highlights the patient's pre-treatment condition and the subsequent score $X^T\xi$ is analyzed in conjunction with the function $g$, which accounts for the additive impact of the treatment. Given data $\{X_i, Y_i, \tau_i\}_{i=1,...,n}$, denote $\bar{Y}_i = \log \frac{Pr(Y_i=1|X_i,\tau_i)}{Pr(Y_i=0|X_i,\tau_i)}$, $Z_i(\xi) = X_i^T\xi - \tau_i$ for a fixed $\xi$. Our partially linear model becomes

$$\bar{Y}_i = X_i^T\beta + g(Z_i(\xi)) + e_i, \qquad (3)$$

where $e_i$ is the error term with $\mathbb{E}(e_i|Z_i(\xi)) = 0$.

### B. Proposed Estimation

In the above formulation, it remains to estimate parameters $\beta, \xi$ and the unknown function $g$, which we denote as $\hat{\beta}, \hat{\xi}$, and $\hat{g}$ respectively.

Condition (3) on $Z_i(\xi)$, we obtain

$$\mathbb{E}(\bar{Y}_i|Z_i(\xi)) = \mathbb{E}(X_i^T\beta|Z_i(\xi)) + g(Z_i(\xi)). \qquad (4)$$

Taking the difference of (3) and (4), we obtain the following least-squares regression:

$$e_{yi} = e_{xi}^T\beta + e_i, \qquad (5)$$

where $e_{yi} = \bar{Y}_i - \mathbb{E}(\bar{Y}_i|Z_i(\xi))$ and $e_{xi} = X_i - \mathbb{E}(X_i|Z_i(\xi))$.

Based on (5), we can estimate $\beta$ by estimating the unknown conditional means using the Nadaraya-Watson (NW) estimator and obtain:

$$\hat{\beta} = \left(\sum_i^n \hat{e}_{xi}\hat{e}_{xi}^T\right)^{-1} \sum_i^n \hat{e}_{xi}\hat{e}_{yi}^T, \qquad (6)$$

where

$$\hat{e}_{xi} = X_i - \frac{\sum_{j\neq i}^n K_{h_i}(Z_j(\xi) - Z_i(\xi)) \cdot X_j}{\sum_{j\neq i}^n K_{h_i}(Z_j(\xi) - Z_i(\xi))}, \qquad (7)$$

$$\hat{e}_{yi} = \bar{Y}_i - \frac{\sum_{j\neq i}^n K_{h_i}(Z_j(\xi) - Z_i(\xi)) \cdot \bar{Y}_j}{\sum_{j\neq i}^n K_{h_i}(Z_j(\xi) - Z_i(\xi))}. \qquad (8)$$

Here, $K(\cdot) : \mathbb{R} \to \mathbb{R}^+$ is any kernel function that satisfies $\int K(t)dt = 1, \int tK(t)dt = 0$, and $0 < \int t^2K(t)dt < \infty$ and $h_i$ is the selected bandwidth.

An estimate of $\hat{\beta}$ depends on a given $\xi$, so we first estimate $\hat{\xi}$ by minimizing residuals as follows:

$$\hat{\xi} = \arg\min_\xi \sum_i^n (r_i(\xi))^2, \qquad (9)$$

where $r_i(\xi) = \hat{e}_{yi} - \hat{e}_{xi}^T\hat{\beta}$.

We then use the obtained $\hat{\xi}$ to compute the corresponding $\hat{\beta}$ from (6). Finally, we estimate $\hat{g}$:

$$\hat{g}(Z) = \frac{\sum_i^n K_h(Z_i(\hat{\xi}) - Z) \cdot (\bar{Y}_i - X_i^T\hat{\beta})}{\sum_i^n K_h(Z_i(\hat{\xi}) - Z)}. \qquad (10)$$

Furthermore, for applications in high-stakes fields, we often want to trade off some predictive accuracy for better interpretability via LASSO regularization [13]. Therefore, on top of the constrained version of (9), we add a Lasso penalty to $\xi$ and optimize the following:

$$\hat{\xi} = \arg\min_\xi \sum_i^n (r_i(\xi))^2 + \lambda \|\xi\|_1 \qquad (11)$$

$$\text{s.t.} \quad \|\xi\|_2 = 1, \xi_1 > 0.$$

The bandwidth $h$ and Lasso penalty $\lambda$ are hyperparameters tunable via cross-validation.

## III. NUMERICAL SIMULATION STUDIES

In our simulation study, we designed four scenarios of varying complexity. Table I summarizes the design for all scenarios.

For scenarios 1 and 2, $X$ are generated from a multivariate Gaussian distribution with random mean $\mu$ and covariance matrix $\Sigma$ where $\mu \sim Unif(-1, 1)^p$ and $\Sigma = AA^T$ with $A \sim Unif(0, 1)^{p \times p}$. The treatment variable is generated using a standard Gaussian distribution $\tau \sim \mathcal{N}(0, 1)$. For the parameters, we simulate $\beta \sim \text{Unif}(-1, 1)^4$ and $\xi \sim \text{Unif}(\mathbb{Q}_1)$ where $\mathbb{Q}_1$ denotes the first quadrant on the unit sphere. The

TABLE I
SETUP OF SIMULATED SCENARIOS.

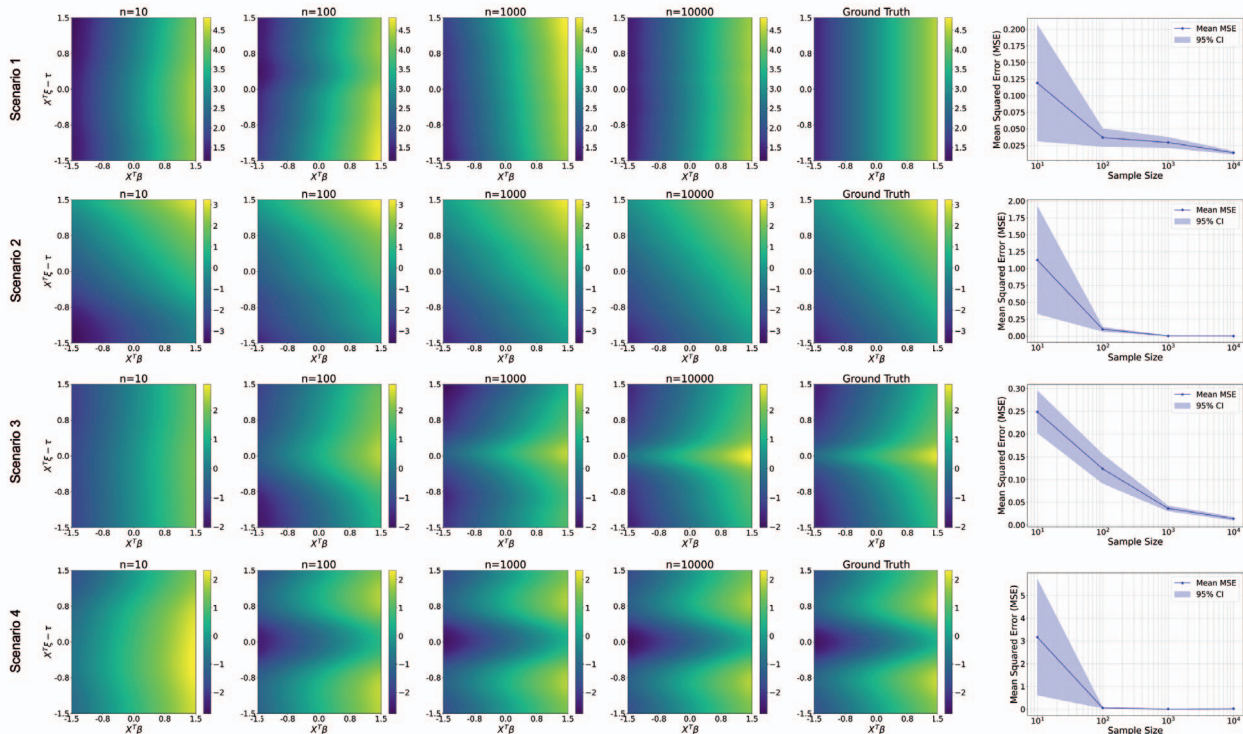| | Covariates Type | Partially Linear Interaction | Treatment-Covariates Function | Intervention Setting |
|---|---|---|---|---|
| Scenario 1 | Continuous | $\bar{Y} = X^T\beta + 3$ | Constant | RCTs |
| Scenario 2 | Continuous | $\bar{Y} = X^T\beta + (X^T\xi - \tau)$ | Linear | RCTs |
| Scenario 3 | Continuous | $\bar{Y} = X^T\beta - 0.5\log(|X^T\xi - \tau|)$ | Unimodal | Observational |
| Scenario 4 | Categorical & Continous | $\bar{Y} = X^T\beta - 1.2\cos(\pi \cdot X^T\xi - \tau) \cdot e^{-(X^T\xi-\tau)^2}$ | Multimodal | RCTs |



Fig. 2.  Convergence of the surface estimation. These heatmaps represent function value estimations with respect to the arguments $X^T\beta$ and $X^T\xi - \tau$ across varying sample sizes (from $n = 10$ to $10^4$). The rightmost column of each series of heatmaps is the ground truth. We present a plot of the Mean Squared Error (MSE) for each scenario as the accuracy metric for the density estimation, bootstrapped ten times for confidence interval.

treatment-covariate term is constant in scenario 1 and linear in scenario 2. In scenario 1, the treatment assignment does not affect the outcome. The linear term in scenario 2 implies that the optimal treatment decision is always the extreme point in the feasible region. In both scenarios, we let $p = 8$.

Scenario 3 simulates an observational setting in which the treatment is influenced by the covariates. We simulate the following covariates: $X_4 \sim \text{Unif}(-1, 1)$ and $X_1 = \sqrt{|X_4|} + \text{Unif}(-1, 1)$, $X_2 = 0.5 \times X_1 + \text{Unif}(-0.5, 0.5)$, $X_3 = 0.3 \times X_1 + 0.3 \times X_2 + \text{Unif}(-0.4, 0.4)$ . For the treatment, we have $\tau = \sin(X_2 X_3) + \text{Unif}(-0.6, 0.6)$, which are influenced by covariates $X_2$ and $X_3$. Similarly, we simulate $\beta \sim \text{Unif}(-1, 1)^4$ and $\xi \sim \text{Unif}(\mathbb{Q}_1)$. Lastly, we have $\bar{Y} = X^T\beta - 0.5\log(|X^T\xi - \tau|)$ which indicates that there exists a unique optimal treatment given the covariates.

For scenario 4, we simulate our covariates as a mixture of categorical and continuous variables in an RCTs (Randomized Control Trials) setting. Let $n_c, n_b$ denote the number of continuous and binary variables. We simulate $X_c \sim$

$\mathcal{N}(\mu_c, \Sigma_c)$ with $\mu_c \in \mathbb{R}^{n_c}, \Sigma_c \in \mathbb{S}_+^{n_c}$. For binary variables, we simulate $X_b[1] \sim \text{Bernoulli}(p_1)$ with $p_1 \sim Unif(0, 1)$. For $i = 2, ..., n_b$, we simulate $X_b[i] = \text{Bernoulli}(p_i)$ with $p_i = \frac{1}{1+\exp(-(a_i \cdot X_b[i-1]+b_i))}$. The coefficients $a_i$ and $b_i$ control the dependency on the previous random variable. We generate $\tau \sim Unif(-1, 1)$ independent of the covariates for the treatment to resemble an RCT scenario. We set $n_c = 12, n_b = 8, a_i = 0.5, b_i = -0.25$. The treatment-covariates interaction term induces some multimodal locally optimal treatment decisions, which can be seen in the heatmap of scenario 4 in Fig. 2.

We estimated the model for each scenario using the constrained optimization framework in (11). We used the Tree-based Parzen Estimators (TPE) optimization algorithm [14] within a defined search space for $\xi \in \Theta$. Our findings are illustrated in Fig. 2.

In scenario 1, the convergence towards the ground truth is achieved swiftly, even with a limited number of training samples. Notably, with $n = 10^4$, the estimated heatmap

accurately reflects the treatment's non-influence, maintaining a consistent scale along the y-axis. A similar swift convergence is observed in scenario 2. Scenario 3 requires a larger sample size to align with the ground truth. Starting from $n = 10^3$, the estimation progressively discerns the unique optimal treatment, as captured by the interaction term. Scenario 4 introduces increased dimensionality, more complex data types, and a multimodal interaction term. The estimation necessitates a larger number of samples to achieve empirical convergence and discern this multimodality. For all scenarios, the MSE and the variances drop with increasing sample sizes, and the first three scenarios achieve significantly lower values than scenario 4 due to lower dimensionality and problem complexity. Moreover, our proposed estimation works empirically for different data types as well as intervention settings.

## IV. APPLICATION TO AN ANTICOAGULANT STUDY

We applied our framework to an anticoagulant study with a curated dataset of over 6000 patients from the International Warfarin Pharmacogenomics Consortium (IWPC) [15]. Warfarin is a treatment for blood clots that can lead to thromboembolism. Personalized warfarin dosing is valuable because of its narrow therapeutic range and diverse individual responses [16].

The cohort selection process is shown in Fig. 3. The selected cohort is divided into train and test sets with a ratio of 9:1. We extracted 25 covariates with both clinical and pharmacogenetic variables, including physical attributes, medical conditions, concomitant medications, genotype status of functional warfarin genetic polymorphisms, therapeutic INR (International Normalized Ratio), etc. In practice, pharmacogenetic-guided dosing can improve dosing effectiveness in highly sensitive responders versus that of patients who received fixed-dose [17]. Different from studies that considered a fixed target INR [12], we constructed a binary outcome that indicates whether a patient's INR falls within their personalized target range. This individualization is important since for deep vein thrombosis and atrial fibrillation, the target is 2.0–3.0; for high-risk heart valve patients, it's above 3.0; a 1.5–2.0 range is recommended for some heart valve patients to reduce bleeding risks [18].

In this study, our objective diverges from refining the accuracy of warfarin dosage predictions, as explored in previous research [19], [20], [21]. Instead, our model focuses on providing a quantifiable interpretation of the influence of baseline covariates through the learned coefficients while accounting for their interactions with the treatment.

### A. Distillation of Soft Labels

Note that in this case, only binary labels are available for training, which causes the left-hand side of Eq. (3), $\bar{Y}_i = \log \frac{Pr(Y_i=1|X_i,\tau_i)}{Pr(Y_i=0|X_i,\tau_i)}$, to be ill-defined. To address this issue, we train an intermediary model and then use the prediction of each sample (i.e., soft labels) in lieu of the binary ground truth to construct $\bar{Y}_i$. We refer to this intermediary model as the "expert model" and our model as the "student model"
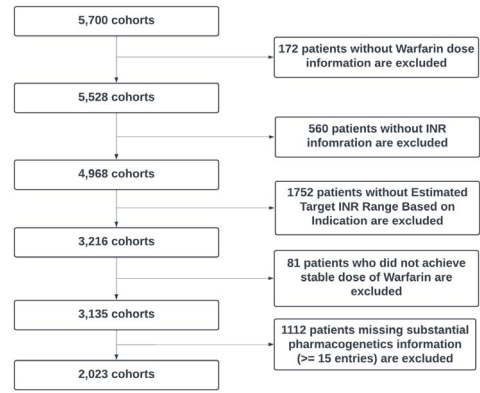


Fig. 3. Cohort selection in the IPWC dataset.

following [22]. In our implementation, we used XGBoost as the expert model.

For data preprocessing, we normalized all continuous variables and one-hot encoded all categorical variables. Since the labels are highly imbalanced, we augmented the data with SMOTE resampling. The efficacy of the expert model is summarized in the top-left panel in Fig. 4. The expert model attaches a soft label to every patient as the training input for the student model.
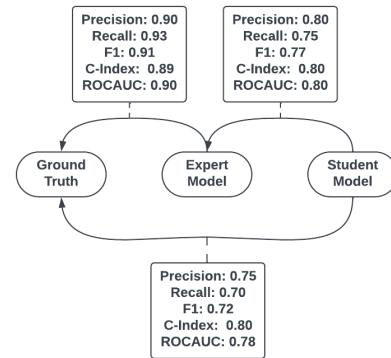


Fig. 4. Predictive performance comparison. The arrow direction indicates the comparison's truth benchmark.

### B. Student Model Training

We aligned the train, validation, and test split with the expert model and estimated the student model. To mitigate the computational load, we used the Epanechnikov kernel with compact support, which is defined as $K(t) = \frac{3}{4}(1 - t^2)\mathbb{1}_{\{|t|\leq 1\}}$. Fig. 5 and 7 show the estimated coefficients with bootstrapped confidence intervals with $k = 30$.

As a classification model, the student model's performance is constrained by the expert model because the student model does not have direct access to the ground truth and is only trained by the expert model. Using a more sophisticated and well-specified expert model may improve the accuracy.

In our model, $X^T \beta$ provides a diagnostic score of the patient. This term can be analogized to the linear predictor in
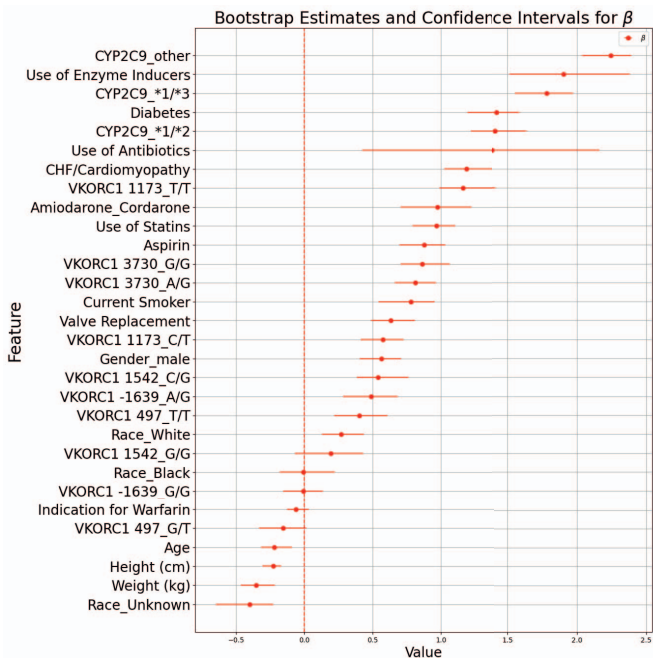
Fig. 5. Values of $\beta$.

logistic regression. We compare the linear score estimation between our model and other linear models that offer similar insights, including Logistic Regression (LR) and Linear Discriminant Analysis (LDA), each with a variant in which the treatment variable acts as a baseline covariate.

The comparison in Fig. 6 shows a consensus in the sign of many parameters across various models, with our estimations in red standing out. Our model shows caution with genetic polymorphism features, pushing many coefficients towards zero, unlike other linear models. It takes a bolder stance on medical history and condition features, resulting in larger coefficient estimates. However, this analysis does not fully assess the clinical significance of these differences, necessitating further evaluation by domain experts to understand the clinical impact of our findings.

*C. Interpretation of Estimates*

For $X^T\beta$, Fig. 5 indicates that larger values for "age", "height", and "weight" have a negative association with pretreatment and may adversely affect the patient's ability to achieve the target INR range. This effect is attributed to higher values for these variables generally corresponding to greater body sizes. Additionally, specific genotypes, such as variants of CYP2C9, are linked to diminished enzyme activity. Individuals carrying these alleles typically metabolize warfarin more slowly, thereby facing a higher risk of bleeding at conventional doses. Consequently, they often necessitate a lower warfarin dose to attain the target range.

For the analysis of the second score $X^T\xi$, it suffices to examine function $g$ at a fixed value of $X^T\beta$. As illustrated in Fig. 8, a higher warfarin dose is beneficial for patients with scores from -1.0 to 0.5. Conversely, for those with scores



Fig. 6. Linear diagnosis score comparison.

above 1.0, a significantly lower dosage is recommended. Referring back to Fig. 7, "age" has a large coefficient, likely due to older patients' increased sensitivity to warfarin and escalated risk of bleeding complications, necessitating lower doses for the same therapeutic effect. Likewise, the substantial negative coefficient associated with "weight" suggests that higher body weight typically requires greater warfarin doses to regulate their INR levels. Additionally, the negative sign for "valve replacement" reflects that patients with valve replacements often need higher warfarin doses. This requirement is due to the elevated risk of clot formation in these patients, warranting a higher target INR.

## V. CONCLUSION

In this study, we introduce a modified partially linear model that captures both the linear effects of covariates on outcomes and the interaction between covariates and treatment. Although our model and the estimation approach show promise, it is limited by the lack of a theoretical underpinning for parameter estimation: identifying sufficient conditions under which our approach achieves statistical consistency is an important future direction. Other future directions include
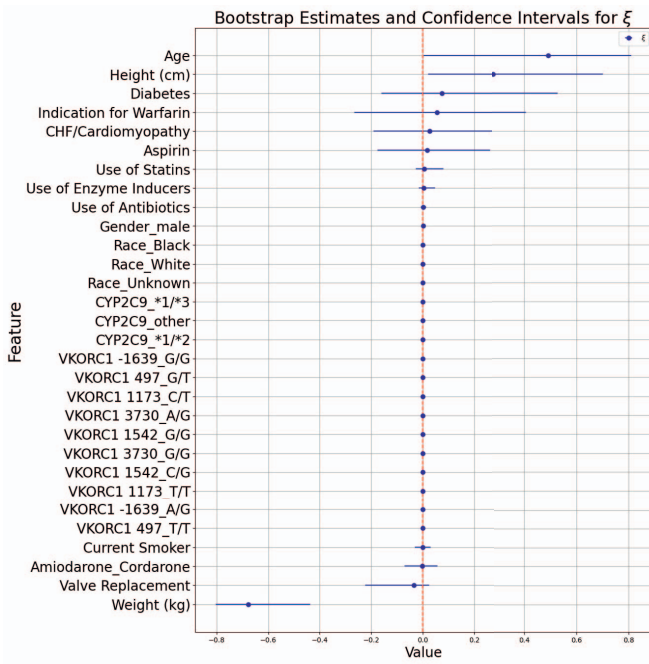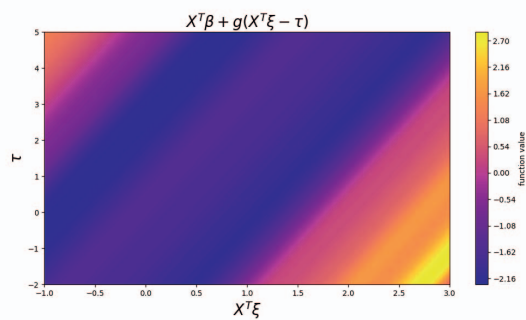
Fig. 7. Values of $\xi$.



Fig. 8. Estimated density function at fixed $X^T\beta = 0.5$.

incorporating domain knowledge for better feature selection, using longitudinal data for temporal insights, and enhancing the model's classifier effectiveness through techniques such as adversarial sample identification [23]. Addressing these limitations can bolster the model's potential as a valuable tool in data analysis and optimization of treatments, particularly in the healthcare space.

## REFERENCES

[1] D. Ruppert, M. P. Wand, and R. J. Carroll, *Semiparametric regression*. Cambridge University Press, 2003.

[2] P. M. Robinson, "Root-n-consistent semiparametric regression," *Econometrica*, pp. 931–954, 1988.

[3] M. Hristache, A. Juditsky, and V. Spokoiny, "Direct estimation of the index coefficient in a single-index model," *Annals of Statistics*, pp. 595–623, 2001.

[4] X. d'Haultfoeuille, C. Gaillac, and A. Maurel, "Partially linear models under data combination," National Bureau of Economic Research, Tech. Rep., 2022.

[5] Y. Wang, Y. Wu, M. H. Jacobson, M. Lee, P. Jin, L. Trasande, and M. Liu, "A family of partial-linear single-index models for analyzing complex environmental exposures with continuous, categorical, time-to-event, and longitudinal health outcomes," *Environmental Health*, vol. 19, pp. 1–16, 2020.

[6] H. Park, E. Petkova, T. Tarpey, and R. T. Ogden, "A constrained single-index regression for estimating interactions between a treatment and covariates," *Biometrics*, vol. 77, no. 2, pp. 506–518, 2021.

[7] H. Park, T. Tarpey, M. Liu, K. Goldfeld, Y. Wu, D. Wu, Y. Li, J. Zhang, D. Ganguly, Y. Ray *et al.*, "Development and validation of a treatment benefit index to identify hospitalized patients with covid-19 who may benefit from convalescent plasma," *JAMA network open*, vol. 5, no. 1, pp. e2 147 375–e2 147 375, 2022.

[8] G. Chen, D. Zeng, and M. R. Kosorok, "Personalized dose finding using outcome weighted learning," *Journal of the American Statistical Association*, vol. 111, no. 516, pp. 1509–1521, 2016.

[9] N. Kallus and A. Zhou, "Policy evaluation and optimization with continuous treatments," in *International conference on artificial intelligence and statistics*. PMLR, 2018, pp. 1243–1251.

[10] P. Jin, W. Lu, Y. Chen, and M. Liu, "Change-plane analysis for subgroup detection with a continuous treatment," *Biometrics*, vol. 79, no. 3, pp. 1920–1933, 2023.

[11] E. B. Laber and Y.-Q. Zhao, "Tree-based methods for individualized treatment regimes," *Biometrika*, vol. 102, no. 3, pp. 501–514, 2015.

[12] H. Park, E. Petkova, T. Tarpey, and R. T. Ogden, "A single-index model with a surface-link for optimizing individualized dose rules," *Journal of Computational and Graphical Statistics*, vol. 31, no. 2, pp. 553–562, 2022.

[13] R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 58, no. 1, pp. 267–288, 1996.

[14] S. Watanabe, "Tree-structured parzen estimator: Understanding its algorithm components and their roles for better empirical performance," *arXiv preprint arXiv:2304.11127*, 2023.

[15] M. Whirl-Carrillo, R. Huddart, L. Gong, K. Sangkuhl, C. F. Thorn, R. Whaley, and T. E. Klein, "An evidence-based framework for evaluating pharmacogenomics knowledge for personalized medicine," *Clinical Pharmacology & Therapeutics*, vol. 110, no. 3, pp. 563–572, 2021.

[16] S. Kim, A. E. Gaweda, D. Wu, L. Li, S. N. Rai, and M. E. Brier, "Simplified warfarin dose–response pharmacodynamic models," *Biomedical Engineering: Applications, Basis and Communications*, vol. 27, no. 01, p. 1550001, 2015.

[17] N. Dietz, C. Ruff, R. P. Giugliano, M. F. Mercuri, and E. M. Antman, "Pharmacogenetic-guided and clinical warfarin dosing algorithm assessments with bleeding outcomes risk-stratified by genetic and covariate subgroups," *International Journal of Cardiology*, vol. 317, pp. 159–166, 2020.

[18] J. D. Puskas, M. Gerdisch, D. Nichols, L. Fermin, B. Rhenman, D. Kapoor, J. Copeland, R. Quinn, G. C. Hughes, H. Azar *et al.*, "Anticoagulation and antiplatelet strategies after on-x mechanical aortic valve replacement," *Journal of the American College of Cardiology*, vol. 71, no. 24, pp. 2717–2726, 2018.

[19] R. J. Chen, J. J. Wang, D. F. Williamson, T. Y. Chen, J. Lipkova, M. Y. Lu, S. Sahai, and F. Mahmood, "Algorithmic fairness in artificial intelligence for medicine and healthcare," *Nature biomedical engineering*, vol. 7, no. 6, pp. 719–742, 2023.

[20] Z. Wang, J. Poon, J. Yang, and S. Poon, "Warfarin dose estimation on high-dimensional and incomplete data," 2021.

[21] H. Lee, H. J. Kim, H. W. Chang, D. J. Kim, J. Mo, and J.-E. Kim, "Development of a system to support warfarin dose decisions using deep neural networks," *Scientific Reports*, vol. 11, no. 1, p. 14745, 2021.

[22] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," 2015.

[23] B. Heo, M. Lee, S. Yun, and J. Y. Choi, "Knowledge distillation with adversarial samples supporting decision boundary," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 33, no. 01, 2019, pp. 3771–3778.