003 004

010 011

012

013

014

015

016

017

018

019

021

# CONFORMAL PREDICTION FOR DOSE-RESPONSE MODELS WITH CONTINUOUS TREATMENTS

Anonymous authors

Paper under double-blind review

## ABSTRACT

Understanding the dose-response relation between a continuous treatment and the outcome for an individual can greatly drive decision-making, particularly in areas like personalized drug dosing and personalized healthcare interventions. Point estimates are often insufficient in these high-risk environments, highlighting the need for uncertainty quantification to support informed decisions. Conformal prediction, a distribution-free and model-agnostic method for uncertainty quantification, has seen limited application in continuous treatments or dose-response models. To address this gap, we propose a novel methodology that frames the causal dose-response problem as a covariate shift, leveraging weighted conformal prediction. By incorporating propensity estimation, conformal predictive systems, and likelihood ratios, we present a practical solution for generating prediction intervals for dose-response models. Additionally, our method approximates local coverage for every treatment value by applying kernel functions as weights in weighted conformal prediction. Finally, we use a new synthetic benchmark dataset to demonstrate the significance of covariate shift assumptions in achieving robust prediction intervals for dose-response models.

025 026 027

028

### 1 INTRODUCTION

029 How can we determine the optimal dose for a patient to ensure the best therapeutic outcome? What is the impact of discounts in an online store on sales? What impact does  $CO_2$  concentration have 031 on local climates? At the core of each of these questions lies a shared causal idea: understanding 032 the dose-response relation under continuous treatments to inform decision-making. In many cases, 033 these decisions bear significant consequences, where relying solely on point estimates may be insuffi-034 cient (Feuerriegel et al., 2024). Particularly in high-stakes situations, augmenting predictions with uncertainty quantification (UQ) can significantly improve decision-making processes (Feuerriegel et al., 2024). For instance, while the estimated causal effect of a continuous treatment may appear positive, prediction intervals could suggest a largely negative outcome for a specific individual. Such 037 insights are crucial for deciding interventions. To tackle this, conformal prediction (CP) offers a robust solution for UQ, being both distribution-free and model-agnostic, with formal coverage guarantees (Vovk et al., 2022). 040

In this work, we seek to extend CP to UQ in dose-response models, aiming to aid decision-makers
 with more informed estimates to tackle such questions. We introduce a novel approach for deriving
 prediction intervals in the continuous treatment setting using weighted conformal prediction by
 combining propensity estimation with weighted conformal predictive systems. Furthermore, with the
 aid of a novel synthetic benchmark, we show how viewing the problem as a covariate shift approach
 provides coverage across all treatment values to help create more individualized dose-response curves.

047 048

## 2 BACKGROUND

In this paper we expand upon the potential outcomes framework introduced in Rubin (2005), otherwise known as the Rubin framework to accommodate continuous treatments. Consider a continuous treatment variable  $T \in [t_L, t_U]$  with a lower bound  $t_L$  and upper bound  $t_U$ , observed covariates X, and potential outcomes  $Y(t) \in \mathbb{R}$  representing the outcome that would be observed under treatment level t. The Conditional Average Dose-Response Function (CADRF) is defined as  $\nu(x, t) =$  054 E[Y(t)|X = x], the expected value over the Individual Dose-Response Functions (IDRF) for all 055 individuals with observed X. Similar to Conditional Average Treatment Effects (CATE), to estimate 056 the CADRF we make the following standard assumptions (Rubin, 2005; Hirano & Imbens, 2004): • Unconfoundedness:  $Y(t) \perp T \mid X, \forall t \in T$ . This assumption states that, conditional on 058 the observed covariates, the treatment assignment is independent of the potential outcomes. In other words, there are no unobserved confounders that influence both the treatment 060 assignment and the outcome. 061 • Overlap or positivity:  $0 < P(T = t | X = x) < 1, \forall t \in T \text{ with } x \in X.$  The overlap 062 assumption ensures that for every covariate value x, there is a positive probability of 063 receiving any treatment level. This is crucial for estimating treatment effects across the 064 entire range of treatment levels. 065 • Consistency: Y = Y(t) with probability 1. This assumption links the observed outcomes to 066 the potential outcomes, stating that the observed outcome is equal to the potential outcome 067 corresponding to the treatment received. 068 069 Ouantifying the IDRF requires observing the Y(t) for all possible treatment values. These treatment values are all counterfactuals and thus impossible to observe as we only can observe Y for a single 071 treatment value t at a time. Furthermore for estimating the CADRF, likewise with CATE estimation, 072 the distribution of the treatment assignment can bias the estimation (Hirano & Imbens, 2004). This distribution of the treatment assignment is called the propensity distribution, which was initially 073 defined for binary treatments. Hirano & Imbens (2004) introduced the generalized propensity score 074 (GPS) for continuous treatments that aims to unbias the CATE estimation for continuous treatments. 075 The GPS is defined as  $\pi(t_i|x) = f_{T|X}(T = t_i|X = x)$ , which is the evaluation of  $T = t_i$  on 076 the conditional probability density function T|X (Hirano & Imbens, 2004). If the treatment is 077 independent of X, i.e. there are no confounders that influence treatment assignment, then  $f_{T|X}$  is equal for all possible X. Furthermore, the treatment assignment is considered uniformly assigned 079 between lower  $t_L$  and upper  $t_U$  possible treatment if  $f_{T|X}$  represents the density function of the uniform distribution between  $t_L$  and  $t_U$ . The GPS can then be used to mimic the randomly assigned 081 treatment to estimate the unbiased CADRF (Wu et al., 2024). 082 The simplest method to estimate the CADRF is using an S-learner where a single learner is fit on 083 both the covariates X and the treatment T to estimate Y. This approach provides a CADRF for each 084 specific sample by keeping the covariates X constant and changing T to all different treatment values. 085 However, if the treatment in the data is not uniformly assigned then the epistemic error can increase

for specific treatment values  $t_i$  and X = x in low overlap regions or where  $\pi(t_i|x)$  becomes very small. Consequently inferring  $T = t_i$  in these regions would yield unreliable model estimates which should be communicated to ensure correct usage of a CADRF model. The estimated  $\widehat{IDRF}$  can also be seen as follows:  $\widehat{IDRF} = \nu(x, t) + \epsilon_{a,IDRF}(x, t) + \epsilon_{e,IDRF}(x, t)$ .

The eleatoric uncertainty is symbolized by  $\epsilon_{a,IDRF}(x,t) = \nu(x,t) + \epsilon_{a,IDRF}(x,t) + \epsilon_{e,IDRF}(x,t)$ . The aleatoric uncertainty is symbolized by  $\epsilon_{a,IDRF}(x,t)$  created by the inherent variability between individuals having the same covariates.  $\epsilon_{e,IDRF}(x,t)$  symbolises the epistemic uncertainty coming from model specification and finite samples. Estimating both uncertainties creates the opportunity to estimate the ranges of the  $\widehat{IDRF}$ :

**Problem Definition** To accurately estimate the  $\widehat{IDRF}$  for all possible treatment values we require correctly estimating both uncertainties for all treatment values equally, or more formally; for a specific significance level  $\alpha$ , lower treatment bound  $t_L$ , upper treatment bound  $t_U$ , and covariates X, we require prediction intervals C(t, X) such that

$$\mathbb{P}(Y(t) \in C(X, t)) \ge 1 - \alpha, \quad \forall t \in [t_L, t_U]$$
(1)

This requirement necessitates prediction intervals that guarantee coverage for each possible treatment value individually.

103 104 105

100 101

095

# 3 RELATED WORK

107 Our proposed solution combines three different domains: propensity score methods, conformal prediction, and treatment effect or dose-response modelling.

108 Propensity score methods, introduced by Rosenbaum & Rubin (1983), have become widespread 109 in causal inference, especially in observational studies. These methods aim to balance confounders 110 across treatment groups, reducing bias in treatment effect estimates. Hirano & Imbens (2004) gen-111 eralized this propensity score to continuous instead of binary treatments, introducing the generalized 112 propensity score and building the foundation for causal inference with continuous exposures. Wu et al. (2024) used the generalized propensity score for matching continuous treatments to debias 113 the treatment assignment and more accurately estimate the average dose-response curve for all 114 treatment values. Other approaches adapt machine learning techniques to dose-response modelling. 115 For instance, Athey et al. (2019) developed generalized random forests for heterogeneous treatment 116 effect estimation, adaptable to continuous treatments. 117

To provide UQ, this work adapts conformal prediction. Conformal prediction is a model-agnostic 118 method introduced by Vovk et al. (2022) that constructs prediction intervals with guaranteed finite-119 sample coverage under distribution-free assumptions. Conformal prediction uses conformity scores 120 to assess uncertainty. Various improvements, such as the adaptive version by Romano et al. (2019), 121 have increased the flexibility and applicability to even heteroscedastic settings. Additionally, Lei et al. 122 (2018) and Papadopoulos et al. (2002) introduced split conformal prediction, significantly improving 123 computational efficiency. For scenarios involving covariate or distribution shifts, Tibshirani et al. 124 (2019) introduced weighted conformal prediction to ensure coverage under mismatched training and 125 testing data distributions, with additional work by Gibbs & Candes (2021; 2024) and Barber et al. 126 (2023). By reweighting the calibration samples similar to weighted conformal prediction, Guan 127 (2023) introduced localized conformal prediction where the prediction intervals are determined by 128 calibration samples localized around the test sample. Vovk et al. (2019) also introduced conformal 129 predictive systems (CPS); an extension of full conformal prediction that allows extracting predictive distributions instead of prediction intervals. More recently, Jonkers et al. (2024a) combined previous 130 concepts, introducing weighted conformal predictive systems to also account for covariate shifts. 131

132 In causal inference, conformal prediction has mainly been applied to binary treatments. For instance, 133 Lei & Candès (2021) were among the first to apply conformal prediction to treatment effects estima-134 tion in randomized experiments and confounded or observational data. Jonkers et al. (2024b) and Alaa 135 & Ahmad (2024) extended this approach to the potential outcomes framework, providing uncertainty to quantify individual treatment effects. However, the use of conformal prediction in continuous 136 137 treatment settings remains largely unexplored. Schröder et al. (2024) proposed a conformal prediction framework for prediction intervals of treatment effects for continuous treatment interventions. How-138 ever, their approach mainly covers single-treatment interventions and is computationally intensive, 139 requiring optimization per confidence level, treatment, and sample where they provide prediction 140 intervals for a single treatment value. For a more in-depth analysis of Schröder et al. (2024), see 141 Appendix E. 142

Our goal is to achieve predictive coverage across the entire range of the treatment variable in estimating the dose-response curve. To our knowledge, no existing UQ methods offer conformal prediction guarantees for dose-response models with continuous treatments. To address this gap, we propose a novel methodology that seeks to provide this coverage by integrating weighted conformal prediction with propensity score weighting thereby guaranteeing coverage for any treatment value in continuous treatment dose-response models.

149

## 4 Method

150 151 152

153

## 4.1 INTRODUCTION TO CONFORMAL PREDICTION

Before delving into our proposed method, we provide a formal introduction to conformal prediction (Jonkers et al., 2024a; Tibshirani et al., 2019). Conformal prediction offers a powerful method for constructing prediction intervals with guaranteed finite-sample coverage under distribution-free assumptions (Vovk et al., 2022). The key insight of conformal prediction lies in its use of a nonconformity measure to quantify the degree to which a new observation differs from previously observed data.

160 Let us consider a regression problem with the training data being n independent and identically 161 distributed (i.i.d.) data pairs  $Z_1 = (X_1, y_1), ..., Z_n = (X_n, y_n)$ , where  $X_i \in \mathbb{R}^d$  represents a vector of d features and  $y_i \in \mathbb{R}$  the corresponding label. Consider  $Z_{n+1} = (X_{n+1}, y_{n+1})$  a new exchangeable point being the test observation to evaluate and provide prediction intervals. Conformal prediction aims to construct a prediction interval  $\hat{C}(X_{n+1})$  such that

$$\mathbb{P}\{y_{n+1} \in \hat{C}(X_{n+1})\} \ge 1 - \alpha \tag{2}$$

for a pre-specified significance level  $\alpha \in (0, 1)$  where the probability is calculated over the points  $Z_i, i = 1, ..., n$ .

To achieve this, we first define a nonconformity measure  $S((X, y), Z_{1:n})$  that quantifies how different the pair (X, y) is from a multiset  $Z_{1:n} = \{Z_1, ..., Z_n\}$  of data points. The lower the nonconformity measure, the more the pair conforms to the multiset  $Z_{1:n}$ . The most commonly used nonconformity measure is the absolute error  $S((X, y), Z_{1:n}) = |y - \hat{\mu}(X)|$  with  $\hat{\mu}$  an estimator fitted on  $Z_{1:n}$ .

Next, for each possible value  $y \in \mathbb{R}$  that  $y_{n+1}$  could be, we compute the nonconformity scores:

$$R_i^y := S((X_i, y_i), \{(X_1, y_1), \dots, (X_{i-1}, y_{i-1}), (X_{i+1}, y_{i+1}), \dots, (X_n, y_n), (X_{n+1}, y)\}), i = 1, \dots, n$$
(3)

176 177 178

179

181 182 183

185

186

187

174 175

166

 $R_{n+1}^y := S((X_{n+1}, y), \{(X_1, y_1), ..., (X_n, y_n)\})$ (4)

Finally, we construct the prediction interval containing all y where (Jonkers et al., 2024a)

$$\hat{C}(X_{n+1}) = \left\{ y \in \mathbb{R} : \frac{\#\{i=1,\dots,n+1: R_i^y \ge R_{n+1}^y\}}{n+1} \ge 1-\alpha \right\}$$
(5)

Tibshirani et al. (2019) presented conformal prediction slightly differently by using quantile functions instead, which will be more convenient for weighted conformal prediction later on. Tibshirani et al. (2019) defines the  $1 - \alpha$  quantile function as follows, where  $F_R(y)$  represents the distribution of nonconformity scores  $R_i^y$  consisting of a sum of point masses  $\delta_a$  with mass at a where  $R^y \sim$  $F_R(y)$  (Tibshirani et al., 2019).  $F_R(y)$  can then be used to calculate probabilities:

$$Quantile(1-\alpha; F_R(y)) = inf\{R_i^y : \mathbb{P}\{R^y \le R_i^y\} \ge 1-\alpha\}$$
(6)

$$F_R(y) = \frac{1}{n+1} \sum_{i=1}^n \delta_{R_i^y} + \frac{1}{n+1} \delta_{\infty}$$
(7)

Finally, we construct the prediction interval containing all y where

$$\hat{C}(X_{n+1}) = \{ y \in \mathbb{R} : R_{n+1}^y \le Quantile\left(1 - \alpha; F_R(y)\right) \}$$
(8)

This procedure guarantees that  $P(y_{n+1} \in \hat{C}(X_{n+1})) \ge 1 - \alpha$  for any exchangeable distribution of the data and any choice of nonconformity measure (Tibshirani et al., 2019).

### 4.1.1 INDUCTIVE CONFORMAL PREDICTION

204 The previously mentioned conformal prediction approach is computationally heavy as it requires fitting  $n \cdot \#\{\mathbb{R}\} + 1$  estimators  $\hat{\mu}$ . Inductive or split conformal prediction (ICP), introduced 205 by Papadopoulos et al. (2002), tackles this computation issue by splitting the training sequence 206  $Z_{1:n} = \{Z_1, ..., Z_n\}$  into two sets: the proper training set  $Z_{1:m} = \{Z_1, ..., Z_m\}$  and the calibration 207 set  $Z_{m+1:n} = \{Z_{m+1}, ..., Z_n\}$ . A single regression model  $\hat{\mu}$  is fit on the proper training set while the 208 nonconformity scores (e.g.,  $R_i = |y_i - \hat{\mu}(X_i)|, i = m + 1, ..., n$ ) are generated from the calibration 209 set. These scores are sorted in descending order denoted as  $R_1^*, ..., R_{n-m}^*$ . Then, for a new sample 210 with features  $X_{n+1}$ , a point prediction is made  $\hat{y}_{n+1} = \hat{\mu}(X_{n+1})$ . Finally, given a target coverage of 211  $1-\alpha$ , the prediction interval becomes 212

214

$$C(X_{n+1}) = [\hat{y}_{n+1} - R_s^*, \hat{y}_{n+1} + R_s^*]$$
(9)

where  $s = \lfloor \alpha(n - m + 1) \rfloor$  represents the  $1 - \alpha$  quantile of the ordered nonconformity set with size n - m (Jonkers et al., 2024a).

196 197

201 202

203

# 4.1.2 WEIGHTED CONFORMAL PREDICTION

218 Evaluating and requiring coverage guarantees for the dose-response model at all possible treatment values changes the test distribution compared to the training distribution. In the training data, all 219 treatment values are sampled according to their (conditional) training distribution, which can be 220 determined by other variables in the case of confounding. However, every treatment value is possible 221 in testing, and thus, every treatment sample can be sampled. This mimics sampling a new test 222 sample with the treatment value from a uniform distribution, which can be vastly different from the 223 treatment distribution in the training data. Standard conformal prediction only guarantees coverage 224 if the joint distribution of the new sample  $Z_{n+1}$  and  $Z_{1:n}$  remains the same under permutations, 225 which is called the exchangeability assumption (Vovk et al., 2022; Tibshirani et al., 2019). This issue 226 is called covariate shift; The features  $X_{n+1}$  come from a different distribution compared to  $X_{1:n}$ , 227 while the relation between X and y remains the same. More formally:  $X_i \sim P_X$ , i = 1, ..., n and 228  $X_{n+1} \sim P_X$  where  $P_X \neq P_X$  while  $y_i \sim P_{Y|X}$ , i = 1, ..., n.

Weighted conformal prediction provides a solution to tackle this issue (Tibshirani et al., 2019). However, their main assumption is that the likelihood ratio between the training  $P_X$  and the test covariate distribution  $\tilde{P}_X$  is known, defined as

$$w(x) = \frac{d\dot{P}(x)}{dP(x)} \tag{10}$$

The rationale is that they reweight the distribution of nonconformity scores  $F_R(y)$  to make the nonconformity scores more exchangeable with the test population by using the following weights in equation 7 (Tibshirani et al., 2019):

$$p_i^w(X_{n+1}) = \frac{w(X_i)}{\sum_{j=1}^n w(X_j) + w(X_{n+1})} \qquad p_{n+1}^w(X_{n+1}) = \frac{w(X_{n+1})}{\sum_{j=1}^n w(X_j) + w(X_{n+1})}$$
(11)

243 244

251

252

260

265 266 267

269

233 234 235

236

237

238 239

$$F_R(y) = \sum_{i=1}^n p_i^w(X_{n+1})\delta_{R_i^y} + p_{n+1}^w(X_{n+1})\delta_\infty$$
(12)

Consequently, these weights adjust the distribution of nonconformity scores to give more weight to nonconformity scores that are more likely in the test set and vice versa while in standard conformal prediction, every  $R_i$  has equal weight. Also, note that the weights  $p^w(x)$  are normalized, cancelling out any constant terms resulting in w(x) being proportional to  $w(x) \propto \frac{d\tilde{P}(x)}{dP(x)}$ . An extension to split weighted conformal prediction can be done similarly as in section 4.1.1 (Tibshirani et al., 2019).

#### 4.1.3 CONFORMAL PREDICTIVE SYSTEMS

In some cases, providing a prediction interval often does not suffice and a complete predictive distribution is required. The extension proposed by Vovk et al. (2019) produces a predictive distribution by arranging p-values, created using specific conformity measures, into a probability distribution function. A requirement to create a Conformal Predictive System (CPS) is to use a specific type of conformity measures <sup>1</sup> which include monotonic measures. Then, given the training data  $Z_{1:n}$ and observed test sample  $X_{n+1}$ , we define an example of this specific conformity measure S and conformity scores  $R_i^y$  similar as in equations 3 and 4:

$$S((X,y), Z_{1:n}) = y - \hat{\mu}(X)$$
(13)

With  $\hat{\mu}$  an estimator fitted on the training set  $Z_{1:n}$ .  $R_i^y$  and  $R_{n+1}^y$  are then similarly defined as in equation 3 for a CPS. Then, as defined in Vovk et al. (2022) we can define a predictive distribution Q for value y, using a distribution of nonconformity scores  $F_R(y)$  of y to calculate  $\mathbb{P}$ , similarly to the quantile function in equation 6 as follows:

$$Q_R(y,\phi) = \mathbb{P}_{F_R(y)}\{R^y < R_{n+1}^y\} + \phi \cdot \mathbb{P}_{F_R(y)}\{R^y = R_{n+1}^y\}$$
(14)

Where  $\phi$  is a random number sampled from a uniform distribution between 0 and 1 to ensure a smooth predictive distribution. Using the same approach as section 4.1.2, these conformal predictive

<sup>&</sup>lt;sup>1</sup>For the specific definition see Vovk et al. (2020)

systems can be expanded to weighted conformal predictive systems by adjusting  $F_R(y)$  to account for the covariate shift (Jonkers et al., 2024a).

Additionally, conformal predictive systems also suffer from computational issues, therefore Vovk et al.
 (2020) introduced split conformal predictive systems to tackle the same issues in a way analogous to section 4.1.1.

# 4.2 PROPOSED METHODOLOGY: PROPENSITY WEIGHTED CONFORMAL PREDICTION

278 Taking into account the background knowledge of conformal prediction, we first need to formally 279 define the target distribution to tackle our problem definition. A CADRF model  $\hat{\nu}(X,T)$  is trained on triples (X, T, Y) with X d-dimensional observed covariates  $X \in \mathbb{R}^d \sim P_X$  and continuous treatment variables  $T \in [t_L, t_U] \sim P_{T|X}$  to predict responses  $Y \in \mathbb{R} \sim P_{Y|T,X}$ .  $P_X$  represents the covariate 281 distribution,  $P_{T|X}$  represents the observational conditional treatment distribution given confounders 282 X, and  $P_{Y|T,X}$  represents the outcome distribution.  $P_{T|X} = P_T$  if there are no confounders for 283 T. A CADRF model will be used to query the dose-response for all  $T \in [t_L, t_U]$ , creating an 284 interventional distribution  $P_T$ . As every treatment value t is equally likely in this guery we can define 285  $\tilde{P}_T = \tilde{P}_{T|X} = Uniform(t_L, t_U).$ 286

To attain marginal coverage across the interventional test set for a CADRF we can use weighted conformal prediction (Tibshirani et al., 2019). This requires defining the weights w for  $X_i$  and treatment value t using equation 11, which we will call the global (g) propensity (p) weights  $w_{g,p}$ :

$$w_{g,p}(X_i, T_i) = \frac{d\tilde{P}_{X,T}(X_i, T_i)}{dP_{X,T}(X_i, T_i)} = \frac{dP_{T|X}(X_i, T_i)dP_X(X_i)}{dP_{T|X}(X_i, T_i)dP_X(X_i)} = \frac{dP_{T|X}(X_i, T_i)dP_X(X_i)}{dP_{T|X}(X_i, T_i)dP_X(X_i)} = \frac{d\tilde{P}_{T|X}(X_i, T_i)}{dP_{T|X}(X_i, T_i)} = \frac{f_{U(t_L, t_U)}(T_i)}{\pi(T_i|X_i)} = \frac{\frac{1_{[t_L, t_U]}(T_i)}{t_U - t_L}}{\pi(T_i|X_i)} \propto \frac{\mathbbm{1}_{[t_L, t_U]}(T_i)}{\pi(T_i|X_i)}$$
(15)

294 295 296

309 310

312

287

289

with  $\mathbb{1}_{[t_L, t_U]}(T_i)$  the indicator function for  $T_i \in [t_L, t_U]$ .

P

297 We For simplicity, we assume that there is no distribution shift for X and thus  $\tilde{P}_X(X_i) = P_X(X_i)$ 298 (The covariate shift approach for X is detailed in Appendix D.1). Additionally,  $f_{U(t_L,t_U)}$  is the 299 probability density function for the uniform distribution. We also define the propensity function 300  $\pi(T_i|X_i)$  as the probability density function for  $P_{T|X}(T_i)$  as specified in Section 2. To gener-301 ate the prediction intervals at treatment value t for a new sample  $X_{n+1}$  the weights change to  $w_{g,p}(X_{n+1},t) = \frac{1}{\pi(t|X_{n+1})}$ . According to the weighted exchangeability defined in (Tibshirani et al., 302 303 2019), this guarantees marginal coverage over the interventional distribution, for all  $T \in [t_L, t_U]$ , 304 and  $X \sim P_X$ . Tibshirani et al. (2019) also suggested a method to attain local coverage around a 305 predetermined target point  $x_0$  using weighted conformal prediction. Consequently, this can provide 306 varying prediction intervals for different values of  $x_0$  providing another heteroscedastic approach. The proposed weights, which we call the local (l) weights  $w_l$ , utilize kernel functions with bandwidth 307 parameter h: 308

$$w_l^{x_0}(X_i) \propto K\left(\frac{X_i - x_0}{h}\right)$$
 (16)

311 These weights then guarantee

$$\mathcal{D}_{x_0}\{Y_{n+1} \in \hat{C}(X_{n+1}; x_0)\} \ge 1 - \alpha$$
 (17)

This assures coverage *around*  $x_0$ , but  $x_0$  must be determined beforehand. Additionally, if a new  $x_0$ must be evaluated, a new calibration procedure must be performed which should be considered when applying it to general regression use cases. However, for this work, the target interventional treatment distribution is known in advance and can all be computed before deployment. Consequently, for a target treatment value t we can define  $w_l^t(T_i) \propto K(\frac{T_i - t}{b})$  instead.

The local weights guarantee coverage where  $d\tilde{P}_T(T_i)/dP_T(T_i) \propto K(\frac{T_i-t}{h})$ . To adjust the local weights for a CADRF model we need to be aware of the covariate shift introduced by evaluating the interventional distribution and thus must combine  $w_{g,p}$  with  $w_{local}$  to achieve weighted exchangeability. These new weights are defined as  $w_{l,p}$  for target treatment t:

$$w_{l,p}^t(X_i, T_i) \propto \frac{\mathbb{1}_{[t_L, t_U]}(T_i) K\left(\frac{T_i - t}{h}\right)}{\pi(T_i | X_i)} \tag{18}$$

To generate the prediction intervals for target treatment t for a new sample  $X_{n+1}$  the weights are then  $w_{l,p}^t(X_{n+1},t) = \frac{\mathbb{1}_{[t_L,t_U]}(T_i)K((t-t/h))}{\pi(t|X_i)} = \frac{\mathbb{1}_{[t_L,t_U]}(T_i)}{\pi(t|X_i)}$ , which is equal to  $w_{g,p}^t(X_{n+1},t)$ . By using these weights in a weighted conformal prediction framework, we provide a solution to the problem definition in Section 2. Theoretical coverage results of our approach are shown and discussed in Appendix A.

330 331

332

5 EXPERIMENTS

333 334 5.1 Synthetic Data

We evaluate the proposed approach on synthetic data as evaluating the true individual dose-response curve requires knowing the counterfactuals which is not feasible in real-world data.

We used three experimental setups using synthetic data, each having different scenarios that change specific parameters. Setup 1 is inspired by Wu et al. (2024) and Setup 2 follows the experimental setup of Schröder et al. (2024). Both Setup 1 and 2 are clarified in Appendix B. Setup 3 is novel, proposed by us, which mimics a situation where, for every scenario, two different possible dose-response functions are possible that each depends on the covariates, resulting in heavy confounding and thus limited overlap.

For each scenario (over the different setups), 5000 samples were generated using 50 different random 344 seeds resulting in 50 datasets for each scenario. These datasets were split into 25% test (1250), 25% 345 calibration (1250), and 50% training (2500) samples. For each scenario, two different  $\alpha$  (significance 346 values) were evaluated (i.e., 0.1 and 0.05 for a confidence of 90% and 95% resp.). Each sample in the 347 test set is evaluated using 40 treatment values  $t_0$  at equal intervals between the 2% and 98% training 348 treatment value quantile to include varying treatment overlap regions and to mimic the uniform 349 treatment sampling. In the results, the coverage of all treatment values and all samples in the test set 350 are aggregated to a single mean coverage for each experiment, resulting in 50 mean coverage results 351 for every method and scenario.

352 353

354

359 360 5.1.1 SETUP 3

Setup 3 is a new experimental setup proposed in this work to underline the importance of compensating
 for confounding in UQ for CADRF. The covariates are independently sampled from a normal
 distribution. The treatment *T* is confounded by two variables, determining the mean of the treatment
 assignment distribution:

$$X_1, X_2, X_3 \sim \text{Normal}(0, 5)$$
  $T \sim \text{Normal}(X_2 + 0.1 \cdot X_1, 4)$ 

The two scenarios have slightly different outcome distributions, as shown in Table 1. The idea is the same for both scenarios; The individual dose-response function is truly conditional and thus equal treatment values between different individuals or samples do not necessarily translate to each other. In total, there are four different possible dose-response functions depending on the covariates. Furthermore, there is heavy confounding resulting in limited samples where  $T - X_2$  yields high values that in turn create large outcome values. This creates an opportunity for high epistemic uncertainty and limited overlap. For scenario two, the aleatoric uncertainty is also heteroscedastic based on  $X_3$  forcing solutions to look beyond the treatment value to quantify uncertainty.

Scenario	Outcome Distribution
1	$Y \sim sign(X_3) \cdot (2(T - X_2))^2 + 33T \cdot sign(X_1) + Normal(0, 2)$
2	$Y \sim sign(X_3) \cdot (2(T - X_2))^2 + 33T \cdot sign(X_1) \\ + \frac{(sign(X_3) + 1)}{2} \cdot \text{Normal}(0, 30) + \text{Normal}(0, 2)$

375 376 377

369 370

372 373 374

Table 1: The outcome distributions for setup 3

# 378 5.2 IMPLEMENTATION 379

380 In the case of synthetic data, the true propensity distribution, also known as the oracle distribution, is 381 available. However, in real-world applications, the true propensity distribution is mostly unknown. 382 As a result, any method that relies on propensity is evaluated using both the oracle propensity distribution and an estimated propensity distribution in the experiments, denoted as "Oracle" and "Propensity" in the results respectively. The latter can be approximated by estimated distribution in 384 this work is obtained using the Conformal Prediction System (CPS), leveraging conformal prediction, 385 specifically CPS though other propensity estimators could also be used. Do note that CPS quantifies 386 total uncertainty and thus also includes the epistemic uncertainty while ideally only the aleatoric un-387 certainty is included. Additionally, this propensity distribution estimate is not completely guaranteed 388 to be equal to the true conditional propensity distribution, which we theoretically need to get complete 389 finite sample guarantees of validity. Although, in practice, this can still be a valid approximation. A 390 learner is trained on the covariates X to predict the treatment assignment T, deemed the propensity 391 learner. Subsequently, a CPS is calibrated for this learner using the calibration set as it is more 392 practical to extract an empirical density distribution compared to standard conformal prediction. This 393 Since CPS produces an empirical density distribution being a sum of Dirac delta distribution similar to  $F_R$ , thus we require the use of kernel density estimation (KDE) to extract is applied to derive a 394 continuous propensity density function for a treatment value t, given covariates  $X_i$ . Do note that KDE 395 interpolates the density and depending on the KDE parameters may introduce additional epistemic 396 error, which is a drawback of estimating the propensity in this manner. The implementation for 397 the propensity estimation is shown and a computational discussion for Global and Local Propensity 398 WCP is presented in Appendix C.1 and our propensity estimation in Appendix C.2. 399

For the evaluation, several baseline methods were tested and compared, including Gaussian Process, 400 CatBoost with Uncertainty (Duan et al., 2019), Standard Conformal Prediction, and Locally Weighted 401 Conformal Prediction (WCP Local), using weights  $w_l$ ). For the proposed propensity methods we 402 included both variations, using their respective weights: Global Propensity-Weighted Conformal 403 Prediction (WCP Global Oracle and WCP Global Propensity ), using  $w_{a,p}$  and Local Propensity-404 Weighted Conformal Prediction (WCP Local Oracle and WCP Local Propensity, using  $w_{l,p}$ ). The 405 Gaussian Process was included in the comparison due to its widespread use for UQ in regression 406 problems assuming a normal error distribution (Fiedler et al., 2021). All other approaches were based 407 on the CatBoost model used a CatBoost model for the base CADRF learner, chosen for its strong 408 out-of-the-box performance (Dorogush et al., 2018). As a result, the "CatBoost with Uncertainty" 409 method was incorporated as a baseline for comparison of UQ.

410 The propensity learner employed in the propensity-weighted approaches was a 411 CatboostRegressor with 4000 iterations and default hyperparameters. Similarly, the 412 CADRF models were based on CatBoost a CatBoost model with 5000 iterations and default 413 hyperparameters. The CatBoost with Uncertainty approach used the same underlying CatBoost 414 model as the CADRF-other methods to ensure consistency. For the locally weighted conformal 415 approaches, a Gaussian kernel (Theodoridis, 2015) was employed to represent local coverage. The bandwidth parameter for the kernel was set as  $h = 2 \cdot (0.2 \cdot \sigma_{\hat{\pi}})^2$ , where  $\sigma_{\hat{\pi}}$  denotes the standard 416 deviation of the estimated propensity distribution. 417

- 418 419
  - 5.3 Results
- 420

Figure 1 presents the coverage bar plots across all methods for Setup 3 Scenario 1 on the test set. Evaluations on all other More evaluations and CADRF RMSE on all setups and scenarios can be found in Appendix F. The bar plots in Figure 1 clearly illustrate the impact of covariate shift in the treatment on coverage guarantees for methods that did not account for this shift. All propensity-weighting methods assumed uniform treatment sampling during evaluation, mimicking the interpretation of a dose-response curve for decision-making for all treatment values, keeping their coverage guarantees.

428 As can be seen in Figure 1, the global propensity-weighting method shows a high variance in coverage 429 across different experiments. This variance arises due to the calibration process, which considers all 430 possible treatment values between  $t_L$  and  $t_U$ , including those with minimal or no overlap. Depending 431 on the calibration and test set split, certain samples may receive a significantly large likelihood ratio, 436 thereby assigning considerable weight to those values according to Equation 12. This inflates the size



Figure 1: Barplot of the mean coverage calculated over 40 treatment values in 50 experiments for setup 3 scenario 1. Black dotted line is the ideal coverage.

446

447 448

449 of the prediction intervals, leading to conservative estimates. The oracle estimates are also notably more conservative, as they tend to provide narrower propensity distributions. This increases the 450 frequency of large likelihood ratios when compared to the estimated propensity distribution, where the 451 epistemic uncertainty of the propensity learner is also taken into account by the CPS procedure. On 452 the contrary, for a new sample, the local propensity method uses calibration samples with treatment 453 values close to the predefined value  $t_0$  and weighting the propensities as well. Our presented approach 454 uses more comparable calibration samples rather than the entire dataset, resulting in more conditional 455 prediction intervals, provided there are enough calibration samples. Our method thus combines 456 the strengths of both the local and the propensity weighting techniques. These trends are further 457 supported in Figure 2, which shows the prediction intervals for all weighting methods alongside the 458 treatment assignment distribution for a specific test observation. This example highlights the necessity 459 of the uniform treatment sampling assumption for the evaluation of dose-response curves, as both the 460 local weighting method and standard conformal prediction produce inaccurate prediction intervals in regions with low treatment overlap. In these regions, there is insufficient data to support predictions 461 for the model, making these predictions unreliable. Consequently, propensity-weighted methods 462 produce much larger prediction intervals in these areas to compensate for this lack of data support. If 463 there is almost no support or extremely low propensity values, then the propensity-weighted methods 464 provide intervals with an infinite width to show that there is no support in these regions. It is important 465 to note, however, that these intervals may be overly conservative if the model has indeed generalized 466 effectively in such regions. The only way to validate this is through additional data collection in these 467 areas to confirm the model's performance. 468

Note that Schröder et al. (2024) also introduced a conformal prediction method to provide prediction intervals in the continuous treatment setting. However, we did not include a direct comparison in this study due to the high computational complexity of their approach, which would require several years to complete the same experiments we executed in a matter of hours. For a more detailed comparison, including a discussion of the difference in assumptions and methodologies, see Appendix E.

Implementing local propensity weighting in practice is less straightforward as it involves calibrating 474 for a set of predefined treatment values and either storing these models for later use during inference 475 or performing this action in parallel. This has the advantage that it allows conditional prediction 476 intervals to be calculated more quickly during inference. However, a drawback is that evaluating a 477 treatment value not included in the predefined set requires recalibration, and must be considered for 478 inference. Still, this approach is particularly useful in fields like drug dosing, where treatment ranges 479 are often predefined and personalized CADRF is highly relevant or where inference of new treatment 480 values is not time-critical. Additionally, an important factor to consider is the effective sample size 481  $\hat{n}$  in local propensity weighting (Tibshirani et al., 2019; Jonkers et al., 2024a). Reweighting  $F_R(y)$ 482 can significantly reduce the effective sample size, which increases variability in empirical coverage 483 compared to standard conformal prediction. This issue is especially pronounced in regions with low treatment overlap, where the effective sample size can become extremely small. However, as 484 prediction intervals with infinite length are possible using weighted conformal prediction, these 485 infinite intervals additionally provide information to the user where the model cannot be trusted



516 dose-response models can be used to reduce the interval widths and provide even more informative 517 prediction intervals.

## 6 CONCLUSION

525 526 527

524

In this work, we have introduced a novel approach to weighted conformal prediction for UQ in doseresponse models, utilizing propensity estimation and kernel functions as weights for the likelihood ratio. Alongside a newly proposed synthetic dataset, our approach highlights the necessity of compensating for the covariate shift in the treatment assignment when evaluating dose-response models across all possible treatment values. This is achieved by assuming uniform treatment sampling during testing, similar to methods used in discrete treatment effect estimation. Additionally, by leveraging conformal predictive systems to estimate propensity distributions, we offer a practical solution to implement UQ in continuous dose-response estimation for various practical use cases.

Our contribution not only adds to the field of dose-response modelling but also facilitates delivering
 reliable, individualized dose-response functions. Our approach has the potential to aid decision making for personalized dosing in fields such as marketing, policy-making, and healthcare. With this
 UQ for continuous treatments, we are one step closer to achieving truly personalized interventions
 that optimize outcomes for individuals.

540 541	References
5/12	Ahmed M Alaa and Zaid Ahmad. Conformal Meta-learners for Predictive Inference of Individual
543	Treatment Effects. In Advances in Neural Information Processing Systems. Curran Associates,
544	Inc., 2024.
545	Susan Athay Julia Tibebirani and Stafan Wagar Generalized random forests. The Annals of Statistics
546	Susan Atney, June Hosmirani, and Stefan Wager. Generalized fandoin forests. <i>The Annals of Statistics</i> , 47(2):1148–1178. April 2010, ISSN 0000-5364, 2168-8066, doi: 10.1214/18. AOS1700. Publisher:
547	Institute of Mathematical Statistics
548	Institute of Multifulitual Statistics.
549	Rina Foygel Barber, Emmanuel J. Candès, Aaditya Ramdas, and Ryan J. Tibshi-
550	rani. Conformal prediction beyond exchangeability. The Annals of Statistics, 51
551	(2):816–845, April 2023. ISSN 0090-5364, 2168-8966. doi: 10.1214/23-AOS2276.
552	URL https://projecteuclid.org/journals/annals-of-statistics/
553	volume-51/issue-2/Conformal-prediction-beyond-exchangeability/
554	10.1214/23-A0522/6.1011. Publisher: Institute of Mathematical Statistics.
555	Anna Veronika Dorogush, Vasily Ershov, and Andrey Gulin. CatBoost: gradient boosting with
556	categorical features support. arXiv:1810.11363 [cs, stat], October 2018. URL http://arxiv.
557	org/abs/1810.11363. arXiv: 1810.11363.
558	Tony Duan Anand Avati Daisy Vi Ding Khanh K Thai Sanjay Dean Andrew V Ng and Alajardra
559	Schuler NGBoost: Natural Gradient Roosting for Probabilistic Prediction October 2010 UDI
560	https://arxiv.org/abs/1910.03225v4
561	neeps.,, arniv.org, aso, 1910.0022001.
562	Stefan Feuerriegel, Dennis Frauen, Valentyn Melnychuk, Jonas Schweisthal, Konstantin Hess, Alicia
563	Curth, Stefan Bauer, Niki Kilbertus, Isaac S. Kohane, and Mihaela van der Schaar. Causal
564	machine learning for predicting treatment outcomes. <i>Nature Medicine</i> , 30(4):958–968, April 2024.
565	ISSN 1540-170X. doi: 10.1038/841591-024-02902-1. UKL https://www.nature.com/
566	arcretes/s4r59r=024=02902=1. I donisher. Mature I donishing Group.
567	Christian Fiedler, Carsten W. Scherer, and Sebastian Trimpe. Practical and Rigorous Uncertainty
568	Bounds for Gaussian Process Regression. Proceedings of the AAAI Conference on Artificial
569	Intelligence, 35(8):7439–7447, May 2021. ISSN 2374-3468. doi: 10.1609/aaai.v35i8.16912. URL
570	https://ojs.aaai.org/index.php/AAAI/article/view/16912. Number: 8.
572	Rina Foygel Barber, Emmanuel J Candès, Aaditya Ramdas, and Ryan J Tibshirani. The limits
573	of distribution-free conditional predictive inference. Information and Inference: A Journal of
574	the IMA, 10(2):455-482, June 2021. ISSN 2049-8772. doi: 10.1093/imaiai/iaaa017. URL
575	https://doi.org/10.1093/imaiai/iaaa017.
576	Isaac Gibbs and Emmanuel Candes. Adaptive Conformal Inference Under Distribution Shift. In
577	Advances in Neural Information Processing Systems, volume 34, pp. 1660–1672. Curran As-
578	sociates, Inc., 2021. URL https://proceedings.neurips.cc/paper/2021/hash/
579	0d441de75945e5acbc865406fc9a2559-Abstract.html.
580	Isaga Cibbs and Emmanuel Condes Conformal Information for Online Prediction with Arbitrary
581	Distribution Shifts Journal of Machine Learning Research 2024 LIPL https://jmln.org/
582	papers/volume25/22-1218/22-1218 pdf
583	papers, voramero, 22 1210/22 1210.pdr.
584	Leying Guan. Localized conformal prediction: a generalized inference framework for conformal
585	prediction. <i>Biometrika</i> , 110(1):33–50, March 2023. ISSN 1464-3510. doi: 10.1093/biomet/asac040.
586	UKL https://doi.org/10.1093/biomet/asac040.
587	Keisuke Hirano and Guido W. Imbens. The Propensity Score with Continuous Treatments. In
588	Applied Bayesian Modeling and Causal Inference from Incomplete-Data Perspectives, pp. 73–84.
589	John Wiley & Sons, Ltd, 2004. ISBN 978-0-470-09045-9. doi: 10.1002/0470090456.ch7. URL
590	https://onlinelibrary.wiley.com/doi/abs/10.1002/0470090456.ch7.
591	Jaf Jankara Glann Van Wallandaal Lua Duahataan and Safa Van Uaaska. Conformal Dradiction
592 502	Systems Under Covariate Shift April 2024a URL http://arviv.org/abs/2404_15018
222	arXiv:2404.15018 [cs, stat].

- 594 Jef Jonkers, Jarne Verhaeghe, Glenn Van Wallendael, Luc Duchateau, and Sofie Van Hoecke. Confor-595 mal Convolution and Monte Carlo Meta-learners for Predictive Inference of Individual Treatment 596 Effects, June 2024b. URL http://arxiv.org/abs/2402.04906. arXiv:2402.04906 [cs, 597 stat].
- 598 Jing Lei, Max G'Sell, Alessandro Rinaldo, Ryan J. Tibshirani, and Larry Wasserman. Distribution-Free Predictive Inference for Regression. Journal of the American Statistical Association, 113 600 (523):1094–1111, July 2018. ISSN 0162-1459. doi: 10.1080/01621459.2017.1307116. URL 601 https://doi.org/10.1080/01621459.2017.1307116. Publisher: Taylor & Francis 602 \_eprint: https://doi.org/10.1080/01621459.2017.1307116. 603
- Lihua Lei and Emmanuel J. Candès. Conformal Inference of Counterfactuals and Individual Treatment 604 Effects. Journal of the Royal Statistical Society Series B: Statistical Methodology, 83(5):911–938, 605 November 2021. ISSN 1369-7412. doi: 10.1111/rssb.12445. URL https://doi.org/10. 606 1111/rssb.12445. 607
- 608 Harris Papadopoulos, Kostas Proedrou, Volodya Vovk, and Alex Gammerman. Inductive Confidence 609 Machines for Regression. In Tapio Elomaa, Heikki Mannila, and Hannu Toivonen (eds.), Machine Learning: ECML 2002, pp. 345–356, Berlin, Heidelberg, 2002. Springer. ISBN 978-3-540-36755-610 0. doi: 10.1007/3-540-36755-1\_29. 611
- 612 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-613 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and 614 E. Duchesnay. Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, 615 12:2825-2830, 2011.
- Yaniv Romano, Evan Patterson, and Emmanuel Candes. Conformalized Quantile Regres-617 sion. In Advances in Neural Information Processing Systems, volume 32. Curran Asso-618 ciates, Inc., 2019. URL https://proceedings.neurips.cc/paper\_files/paper/ 619 2019/hash/5103c3584b063c431bd1268e9b5e76fb-Abstract.html. 620

624

625

626

627

628

629

630 631

- 621 Paul R. Rosenbaum and Donald B. Rubin. The central role of the propensity score in observational 622 studies for causal effects. Biometrika, 70(1):41-55, April 1983. ISSN 0006-3444. doi: 10.1093/ biomet/70.1.41. URL https://doi.org/10.1093/biomet/70.1.41. 623
  - Donald B Rubin. Causal Inference Using Potential Outcomes. Journal of the American Statistical Association, 100(469):322–331, March 2005. ISSN 0162-1459. doi: 10.1198/016214504000001880. URL https://doi.org/10.1198/016214504000001880.
  - Maresa Schröder, Dennis Frauen, Jonas Schweisthal, Konstantin Heß, Valentyn Melnychuk, and Stefan Feuerriegel. Conformal Prediction for Causal Effects of Continuous Treatments, July 2024. URL http://arxiv.org/abs/2407.03094. arXiv:2407.03094 [cs, stat].
- Sergios Theodoridis. Chapter 11 Learning in Reproducing Kernel Hilbert Spaces. In Sergios 632 Theodoridis (ed.), Machine Learning, pp. 509-583. Academic Press, Oxford, January 2015. ISBN 978-0-12-801522-3. doi: 10.1016/B978-0-12-801522-3.00011-2. URL https://www. 633 sciencedirect.com/science/article/pii/B9780128015223000112. 634
- 635 Ryan J Tibshirani, Rina Foygel Barber, Emmanuel Candes, and Aaditya Ramdas. Conformal Pre-636 diction Under Covariate Shift. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d' Alché-Buc, E. Fox, and R. Garnett (eds.), Advances in Neural Information Processing Systems, volume 32. Cur-638 ran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper\_files/ 639 paper/2019/file/8fb21ee7a2207526da55a679f0332de2-Paper.pdf.
- 640 Vladimir Vovk, Jieli Shen, Valery Manokhin, and Min-ge Xie. Nonparametric predictive dis-641 tributions based on conformal prediction. Machine Learning, 108(3):445-474, March 2019. 642 ISSN 1573-0565. doi: 10.1007/s10994-018-5755-8. URL https://doi.org/10.1007/ 643 s10994-018-5755-8. 644
- Vladimir Vovk, Ivan Petej, Ilia Nouretdinov, Valery Manokhin, and Alexander Gammerman. Compu-645 tationally efficient versions of conformal predictive distributions. *Neurocomputing*, 397:292–308, 646 2020. ISSN 0925-2312. doi: https://doi.org/10.1016/j.neucom.2019.10.110. URL https: 647 //www.sciencedirect.com/science/article/pii/S0925231219316042.

 Vladimir Vovk, Alexander Gammerman, and Glenn Shafer. Algorithmic Learning in a Random World. Springer International Publishing, Cham, 2022. ISBN 978-3-031-06648-1 978-3-031-06649-8. doi: 10.1007/978-3-031-06649-8. URL https://link.springer.com/10. 1007/978-3-031-06649-8.

Xiao Wu, Fabrizia Mealli, Marianthi-Anna Kioumourtzoglou, Francesca Dominici, and Danielle Braun. Matching on Generalized Propensity Scores with Continuous Exposures. *Journal of the American Statistical Association*, 119(545):757–772, January 2024. ISSN 0162-1459. doi: 10.1080/01621459.2022.2144737. URL https://doi.org/10.1080/01621459.2022. 2144737. Publisher: Taylor & Francis \_eprint: https://doi.org/10.1080/01621459.2022.2144737.

## A FINITE SAMPLE COVERAGE GUARANTEES

For counterfactual prediction intervals, the ideal goal is to achieve the following general conditional coverage guarantee:

$$\mathbb{P}_{Y \sim P_{Y|T=t,X=x}}(Y(t) \in \hat{C}(x,t)|X=x) \ge 1 - \alpha, \text{ where } t \in [t_{L}, t_{U}]$$
(19)

, which, under the strong ignorability assumption, is equivalent to:

$$\mathbb{P}_{Y \sim P_{Y|T=t,X=x}}(Y \in \hat{C}(x,t) | X = x, T = t) \ge 1 - \alpha.$$
(20)

However, constructing non-trivial prediction intervals with such conditional guarantees is generally impossible without additional modeling assumptions, as shown in Foygel Barber et al. (2021). Even under the relaxed conditional guarantee, where conditioning is only on the treatment value, as in binary treatment settings (Lei & Candès, 2021):

$$\mathbb{P}_{Y \sim P_X \times P_Y|_{T=t,X}} (Y \in \hat{C}(X,t)|_{T=t}) \ge 1 - \alpha,$$
(21)

the problem persists when the treatment variable t is continuous.

## A.1 PROPOSED FRAMEWORK

681To address this challenge, we introduce a distribution shift in the treatment variable by moving from682the generalized propensity distribution to a user-specified interventional distribution,  $T_{n+1} \sim \tilde{P}_{T|X}$ .683We then leverage the weighted conformal prediction (WCP) framework to construct prediction684intervals. This approach allows us to build on prior theoretical coverage results under both oracle685and estimated likelihood functions(Tibshirani et al., 2019; Lei & Candès, 2021).

Table 2 outlines the two interventional distributions utilized in this work: global propensity, local propensity, and  $\delta$ -propensity (Dirac delta). The latter corresponds to a hard intervention. Relaxing the  $\delta$ -propensity to the local propensity enables the construction of non-trivial prediction intervals (see Remark4). Notably, when  $T \in \{0, 1\}$ , our approach under  $\delta$ -propensity aligns with the counterfactual inference framework for binary treatments proposed in Lei & Candès (2021).

<sup>692</sup> Table 2: Translation of general interventional distribution framework to WCP global, local, and  $\delta$ -propensity.

General	Global propensity	Local propensity	<u><i>δ</i>-propensity</u>
$\tilde{P}_{\mathcal{T} \mid \mathcal{X}}$	$Uniform(t_L, t_U)$	$\frac{\mathbb{1}_{[t_L,t_U]}(T)K(\frac{T-t}{h})}{\sqrt{t_L}\mathbb{1}_{[t_L,t_U]}(T)K(\frac{T-t}{h})dT}$	$\delta(T-t)$
$\underline{w}(X,T)_{\sim}$	$\frac{\mathbbm{1}_{[t_L,t_U]}(T)}{\sqrt{\pi}(\mathcal{I} \mathbf{X})} \sim $	$\frac{\mathbb{1}_{[t_L,t_U]}(T)K\left(\frac{T-t}{h}\right)}{\swarrow} \xrightarrow{\pi(\mathcal{I} X)} \xrightarrow{\pi(\mathcal{I} X)}$	$rac{\delta(T-t)}{\pi(T+X)}$
$\underbrace{\hat{w}(X,T)}_{\sim}$	$\frac{\mathbbm{1}_{[t_L,t_U]}(T)}{\sim \hat{\pi}(T \mid X) \sim}$	$\frac{\mathbb{1}_{[t_L,t_U]}(T)K\left(\frac{T-t}{h}\right)}{\sim\!\!\!\!\!\!\sim\!\!\!\!\!\!\!\!\!\!\sim\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!\!$	$\sim \frac{\delta(T-t)}{\hat{\pi}(T \downarrow X)} \sim$

A.2 PROPOSTION: FINITE-SAMPLE GUARANTEES	
<b>Proposition 1</b> (following Tibshirani et al. (2019); Lei & Candès (2021)).	Assume
$(X_i, T_i, Y_i) \stackrel{i.i.d.}{\sim} P_{\mathbf{Y}} \times P_{\mathbf{Y} \mathbf{Y}} \times P_{\mathbf{Y} \mathbf{T} \mathbf{Y}}$ $i = 1$ n: the likelihood ratio $w(X, T) \propto \frac{d\tilde{P}_T}{dr}$	$\frac{ x }{ x }$ · and
(1) = (1) + (1)	
the estimated likelihood ratio $w(X, 1)$ . Using wCP to construct $C(X, 1)$ , the jo finite-sample bounds apply:	ollowing
<b>S1.</b> ( <i>Oracle Likelihood Ratio</i> ) If $\hat{w}(\cdot, \cdot) = w(\cdot, \cdot)$ , <i>i.e. oracle likelihood ratio function;</i>	then,
$1 - \alpha \le \mathbb{P}_{(X,T,Y) \sim P_X \times \tilde{P}_{T X} \times P_{Y T,X}} \{ Y \in \hat{C}(X,T) \}$	(22)
S2 (Einite Sample with Decrylarity Conditions) If $\hat{\omega}(x) = \omega(x)$ , the new equivalent	
Sz. (Finite Sample with Regularity Conditions) If $w(\cdot, \cdot) = w(\cdot, \cdot)$ , the non-condition	<u>yormu y</u>
scores $S_i$ have no ties almost surely; $P_{T X} \times P_X$ is absolutely continuous with	respect
to $P_{T X} \times P_X$ ; and $(\mathbb{E}_{(X,T)\sim P_X \times P_{T X}}   w(X,T)' )^{\overline{r}} \leq M_T < \infty$ where $r > 0$ of denotes the upper bound of the r-th moment of the likelihood ratio; then,	and $M_r$
$1 - \alpha \le \mathbb{P}_{(X,T,Y) \sim P_X \times \tilde{P}_{T X} \times P_{Y T,X}} \{ Y \in \hat{C}(X,T) \} \le 1 - \alpha + cn^{\frac{1}{r-1}}$	(23)
where c is an arbitrary positive constant depending on $M_r$ and r.	
<b>S3.</b> (Estimated Likelihood Ratio) If $\hat{w}(\cdot, \cdot) \neq w(\cdot, \cdot); \Delta_{w} = \frac{1}{2} \mathbb{E}_{(X,T) \approx P_{X \times P_{T \times Y}}} [ \hat{w}(X) ]$	(T) - w(X, T)
$(\mathbb{E}_{(X,T)}, \mathbb{P}_{Y,T}) = [\hat{w}(X,T)^r])^{\frac{1}{r}} \leq M_r \leq \infty$ and further assuming the same assuming	mptions
as in S2.: then.	
$1 - \alpha - \Delta_w \le \mathbb{P}_{(X,T,Y)\sim P_X \times \tilde{P}_{T X} \times P_{Y T Y}} \{Y \in \hat{C}(X,T)\} \le 1 - \alpha + \Delta_w + \alpha$	$cn^{\frac{1}{r-1}}$
	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~
	(24)
<b>Proof</b> We can reformulate our problem as a covariate shift scenario by treating the tr	antmont
variable as part of the covariates, i.e., defining $X^* = [X, T]$ . Under this transformation:	
• The proof for setting S.1 follows directly from Theorem 2 in Tibshirani et al. (2019	<u>9).</u>
• The proof for setting <b>S.2</b> aligns with Proposition 1 in Lei & Candès (2021). their work focuses explicitly on split-weighted conformalized quantile regression (Romano et al., 2019), the argument extends to WCP because it only depends weighted exchangeability of nonconformity scores and the boundedness of the lik ratio function	While (CQR) on the celihood
• Similarly, the proof for setting S.3 follows from Theorem 3 in Lei & Candès (2021	), along
with its corresponding derivation.	
	_
<b>Remark 1.</b> $r$ specifies which moment of the likelihood ratio $w(X,T)$ is being considered.	Larger
r corresponds to stricter regularity conditions on $w(X,T)$ . $M_r$ defines the upper bound on	the r-th
moment of $w(X,T)$ ensuring the likelihood ratio does not arow too large and remains well 1	behaved.
moment of w(21, 1), cashing the inclinious ratio aces not grow too targe and remains well-t	
<b>Remark 2.</b> Note that the term $cn^{\frac{1}{r-1}}$ , represents the upper bound of the expectation of m	aximum
<b>Remark 2.</b> Note that the term $cn^{\frac{1}{r-1}}$ , represents the upper bound of the expectation of <i>m</i> weight (probability), i.e., $\mathbb{E}\left[\max_{i \in [1, \dots, n] \cup \{\infty\}} p_i^w(X_{n+1})\right]$ , which under no covariate shift to $\frac{1}{n+1}$ the upper bound of unweighted conformal prediction.	aximum is equal
<b>Remark 2.</b> Note that the term $cn^{\frac{1}{r-1}}$ , represents the upper bound of the expectation of m weight (probability), i.e., $\mathbb{E}\left[\max_{i \in [1,,n] \cup \{\infty\}} p_i^w(X_{n+1})\right]$ , which under no covariate shift to $\frac{1}{n+1}$ the upper bound of unweighted conformal prediction. <b>Remark 3.</b> The bounding condition assumed in S.2 and S.3 in Pro-	aximum is equal position
<b>Remark 2.</b> Note that the term $cn^{\frac{1}{r-1}}$ , represents the upper bound of the expectation of <i>m</i> weight (probability), i.e., $\mathbb{E}\left[\max_{i\in[1,\dots,n]\cup\{\infty\}}p_i^w(X_{n+1})\right]$ , which under no covariate shift to $\frac{1}{n+1}$ the upper bound of unweighted conformal prediction. <b>Remark 3.</b> The bounding condition assumed in S.2 and S.3 in Pro	aximum is equal position that

756	<b>Demonsh 4</b> For active <b>C</b> 1 the evolution on positivity accumution and havial test $d\tilde{P}_{T X}$
757	<b>Kemark 4.</b> For setting 5.1, the overlap or positivity assumption can be violated, i.e., $\frac{1}{dP_{T X}} = \infty$
758	in terms of the interventional distribution. However, this results in the trivial interval $(-\infty, \infty)$ ,
759	since $w(X_i) = 0, \forall i \in [1,, n]$ and $w(X_{n+1}) = \infty$ resulting in $p_i^w(X_{n+1}) = 0, \forall i \in [1,, n]$ and
760	$p_{n+1}^{\omega} = 1$
761	Remark 5. Since inductive (or split) conformal prediction is a special case of conformal prediction,
762	Proposition 1 also applies to inductive conformal prediction, which we use in our experiments.
763	<b>Remark 6.</b> With an estimated likelihood ratio under weighted COR, our approach also follows the
764	asymptotic double robustness result (see Theorem 1 (Lei & Candès, 2021)).
765	
766	
767	B SYNTHETIC DATA
768	
769	B.1 Setup 1
770	
771	For setup 1, inspired by Wu et al. (2024), six independent covariates are sampled from various
772	distributions representing both continuous and discrete values:
773	$X_{t}, X_{0}, X_{0}, X_{t} \sim Normal(0, 1)$
774	$X_1, X_2, X_3, X_4 \sim \text{Normal}(0, 1)$
775	$\Lambda_5 \sim \text{Uniform}[-2, 2]$ (Integer)
776	$X_6 \sim \text{Uniform}(-3,3)$

The treatment value is confounded by all variables in this setup and thus determined by a treatment function  $T_{\mu}$ . All scenarios share the same treatment function except for scenario 3, where a quadratic term was added. The treatment functions are shown in Table 3. 

Scenario	Treatment function
1, 2, 4, 5, 6, 7, 8	$T_{\mu} = -0.8 + X_1 + 0.1X_2 - 0.1X_3 + 0.2X_4 + 0.1X_5 + 0.1X_6$
3	$T_{\mu} = -0.8 + X_1 + 0.1X_2 - 0.1X_3 + 0.2X_4 + 0.1X_5 + 0.1X_6 + \frac{3}{2}X_3^2$

Table 3: The treatment functions for all scenarios in setup 1.

The true assigned treatment value T is then sampled from a treatment assignment distribution to add randomness and ensure some overlap in the simulated data. This treatment assignment distribution is different for various scenarios to evaluate the differences in the assumed distributions. The various functions are shown in Table 4 

Scenario	Treatment T	Treatment Assignment Distribution
1	$9T_{\mu} + 17$	Normal $(0,5)$
2	$15T_{\mu} + 22$	StudentT(df = 2)
3	$9T_{\mu} + 15$	Normal(0,5)
4	$49 \frac{e^{T_{\mu}}}{1 + e^{T_{\mu}}} - 6$	Normal(0,5)
5	$42\frac{1}{1+e^{T_{\mu}}}+18$	Normal(0,5)
6	$7log( T_{\mu}  + 0.001) + 13$	Normal(0,4)
7	$7T_{\mu} + 16$	Normal(0,1)
8	$7T_{\mu} + 16$	$20 \cdot \text{Beta}(\alpha = 2, \beta = 8)$

Table 4: The propensity functions per scenario for Setup 1

Now, given both the covariates X and the assigned treatment T the outcome function is defined as a random variable sampled from a normal distribution with a variance of 5, with the mean a function dependent on both the treatment and the covariates:  $Y \sim -1 - (2X_1 + 2X_2 + 3X_3^3 - 20X_4 - 2X_5 + 20X_6)$  $-0.1T(1 - X_1 + X_4 + X_5 + X_3^2) + 0.13^2 |T|^3 sin(X_4) + Normal(0,5)$ B.2 SETUP 2 Setup 2 tests the different treatment assignment distributions in the two different scenarios, which is the same experimental setup as proposed by Schröder et al. (2024). The covariates are sampled from a discrete uniform distribution. The treatment is sampled from the treatment assignment distributions shown in Table 5. The outcome function is sampled from a normal distribution with a mean determined by a sinus function based on both X and T:  $X \sim \text{Uniform}[1, 4]$  (Integer)  $Y \sim sin((0.05\pi)(T-X)) + Normal(0, 0.1)$ Treatment Assignment Distribution Scenario  $T \sim p \cdot \text{Uniform}(0, 5X) + (1 - p)\text{Uniform}(5X, 40), p \sim \text{Bernoulli}(0.3)$  $T \sim \text{Normal}(5X, 10)$ Table 5: The propensity functions per scenario for Setup 2 С **PROPENSITY DISTRIBUTION ESTIMATION**ALGORITHM PSEUDOCODE AND COMPUTATIONAL ANALYSIS PROPENSITY-BASED WEIGHTED CONFORMAL PREDICTION PSEUDOCODE C.1 Algorithm 1 presents the fit procedure for both the Local and the Global Propensity WCP, using their respective weights  $w_{l,p}^t$  and  $w_{g,p}^t$  for an array of treatment values we want to evaluate  $t_{eval}$ . The pseudocode is written for any Kernel, although in the experiments, we used the Gaussian kernel as presented in the methodology section. The pseudocode assumes either a pre-fitted propensity estimator  $\hat{\pi}$  or having access to an Oracle estimator. The method used to fit the propensity estimator in this paper is presented in Appendix C.2. Algorithm 2 then presents how the prediction intervals for a significance level  $\alpha$  are generated using both Local and Global Propensity WCP as the implementation is the same for both methods. The get\_interval function is the prediction interval function of the WCP method.

Alg	gorithm 1 Fit and calibrate Local or Global Propensity WCP
1:	<b>Input:</b> Training covariates $X_{tr}$ , calibration covariates $X_{cal}$ , training outcome $u_{tr}$ , calibration
	outcome $y_{cal}$ , training treatment values $T_{tr}$ , calibration treatment values $T_{cal}$ , calibrated
	PropensityEstimator or oracle $\hat{\pi}$ , to evaluate treatments in array $t_{eval}$ , kernel K, CADRF learner
	$\hat{\mu}$
2:	$\widetilde{\text{Fit}}$ CADRF $\hat{\mu}$ on $(X_{tr}, T_{tr})$ to predict $y_{tr}$
3:	Calculate propensities $\pi_{cal} = \hat{\pi}(X_{cal})$
4:	if Global Propensity WCP then
5:	Calculate weights: $w_{a,p} = 1/\pi_{cal}$
6:	Define WCP as Weighted Conformal Prediction with learner $\hat{\mu}$ and weights $w_{a,p}$ or
	$(X_{cal}, T_{cal}, y_{cal})$
7:	Calibrate WCP
8:	else if Local Propensity WCP then
9:	for $t$ in $t_{eval}$ do
10:	Calculate weights: $w_{l,p}^t = K(T_{cal},t)/\pi_{cal}$
11:	Define $WCP_t$ as Weighted Conformal Prediction with learner $\hat{\mu}$ and weights $w_t^t$ or
	$(X \downarrow T \downarrow y \downarrow)$
12.	Calibrate WCP
13.	end for
14·	end if
15:	<b>Output:</b> Calibrated models $\{WCP_t : t \in t_{evol}\}$ for Local Propensity WCP or WCP for Global
	Pronensity WCP
1.	<b>Input:</b> Test sample X $\rightarrow$ calibrated Propensity Estimator or oracle $\hat{\pi}$ k to evaluate treatments
1.	in array $t_{evol}$ , kernel K. CADRF learner $\hat{\mu}_{e}$ calibrated $WCP_{t}$ for all t in $t_{evol}$ , significance $\alpha$
2:	Calculate $\pi_{n+1} = \hat{\pi}(X_{n+1})$
3:	Calculate weights $w = 1/\pi_{cal}$
4:	for $t$ in $t_{eval}$ do
5:	Predict outcome: $\hat{\mu}(X_{n+1}, t)$
6.	Obtain prediction interval: $\hat{C}^t$ = get interval $(WCP_t (X_{r+1}, t), \alpha, w^t)$
7:	end for
8:	<b>Output:</b> Prediction intervals $[\hat{C}^{t_{\text{eval},1}}, \dots, \hat{C}^{t_{\text{eval},k}}]$
C.2	2 PROPENSITY DISTRIBUTION ESTIMATION PSEUDOCODE
Alg	gorithm 3 presents the propensity distribution estimation using Conformal Predictive Systems
(CI	(5). This results in a propensity distribution array $\pi_{arr}$ with the calculated propensity density to
eac	In sample in $X_{cal}$ . $exp$ is the exponential function and $len(X)$ denotes the length of the array $X$
Alg	gorithm 3 Estimating the Propensity Distribution
1:	<b>Input:</b> training Training covariates $X_{tr}$ , calibration covariates $X_{cal}$ , training treatment values
	$T_{tr}$ , calibration treatment values $T_{cal}$ , Kernel Density Estimator $KD$
2:	fit it propensity learner on $Y$ to predict $T$
	$\pi r$ in propensity rearrier on $\Lambda_{tr}$ to predict $T_{tr}$
3:	<b>calibrate CPS on Calibrate CPS using</b> $X_{cal}$ and $T_{cal}$

- 3: calibrate CPS on Calibrate CPS using  $X_{cal}$  and  $T_{cal}$ 4: Define Initialize  $\pi_{arr}$  with as an array of length  $len(X_{cal})$ 5: for i = 1 to  $len(X_{cal})$  do 6: fit KD(Fit KD on CPS( $X_{cal,i}$ )) 7: Set  $\pi_{arr}[i] = exp(KD(T_{cal,i}))$

- 8: end for
- 9: return  $\pi_{arr}$  Output: Propensity array  $\pi_{arr}$

#### 918 C.3 COMPUTATIONAL OVERHEAD 919

920 The computational overhead is greatest for Local Propensity WCP due to the evaluation over multiple treatment values, so we will focus on this version. Let m denote the number of treatment 921 values in the evaluation array  $t_{eval}$ . In this case, the computational overhead compared to standard 922 weighted conformal prediction (WCP) scales linearly with the number of treatment values, i.e., 923  $O(m \cdot WCP)$ , where WCP refers to the cost of standard weighted conformal prediction. In addition, 924 calculating the propensities  $\pi_{cal}$  on the calibration set incurs an additional computational cost, which 925 depends on the size of the calibration set and the chosen propensity estimator. This step can be done 926 once beforehand, so it does not need to be repeated during each evaluation. 927

If the treatment values in  $t_{eval}$  are known and fixed, the calibration for each treatment value can be 928 precomputed and stored, resulting in saved  $WCP_t$  models. This means that, during inference, the 929 computational overhead is reduced to calculating the propensity for a single new sample once and 930 performing m predictions using the CADRF, followed by retrieving the prediction intervals for each 931 treatment value using the pre-calibrated  $WCP_t$ . Thus, the inference overhead is O(m) for a single 932 inference, consisting of a propensity calculation and m predictions and interval retrievals. In the 933 case of a non-static or on-demand  $t_{eval}$ , the overhead is additive as we need O(mWCP) calibrations 934 and directly afterward O(m) for the inference. 935

936 If there is no Oracle propensity estimator, we need to fit the propensity estimator, which, in our case, also involves fitting the Kernel Density Estimator (KDE) for each sample in  $X_{cal}$ , as detailed in 937 Algorithm 3. This introduces an extra layer of computational overhead, which depends on the size 938 of the calibration set and the output of the CPS, which is an empirical distribution of the treatment 939 values for  $x_{cal}$ . The KDE fitting step needs to be performed for each element of  $X_{cal}$ , resulting in a 940 complexity of  $O(\text{len}(X_{\text{cal}}) \cdot \text{KDE})$ , where  $\text{len}(X_{\text{cal}})$  is the number of calibration samples and KDE 941 denotes the cost of fitting the KDE. 942

943 944

945

946 947

948

949 950

951

953

957

#### D **EXTENSIONS AND APPLICATIONS OF WEIGHTED CONFORMAL DOSE-RESPONSE CURVES**

Here, we discuss possible extensions and how the proposed method can be applied to various applications.

## **D.1** EXTENSIONS

952 The paper's current setup assumes no covariate shift in the features X between the training, calibration, and test set, i.e.,  $P_X = P_X$ , to simplify the derivation of the propensity-based 954 weights. However, in real-world applications, covariate shifts are much more common and can 955 hamper the coverage guarantee of conformal prediction, and also thus our proposed method 956 (Tibshirani et al. (2019)). If we assume  $P_X \neq P_X$  in equation 15, we observe that this results in adding a multiplicative term that represents the likelihood term for the covariate shift in X. As such, 958 both  $w_{d,p}^t$  and  $w_{d,p}$  can be easily adjusted to cover a covariate shift in the test set if the covariate shift is known or can be calculated, analogous to Tibshirani et al. (2019), resulting in the following new 959 weights: 960

$$w_{g,p}(X_i, T_i) \propto \frac{\mathbb{1}_{[t_L, t_U]}(T_i)}{\pi(T_i|X_i)} \frac{d\tilde{P}_X}{d\tilde{P}_X}$$
(25)

and

$$w_{l,p}^t(X_i, T_i) \propto \frac{\mathbb{1}_{[t_L, t_U]}(T_i) K\left(\frac{T_i - t}{h}\right)}{\pi(T_i | X_i)} \frac{d\tilde{P}_X}{d\tilde{P}_X}$$
(26)

969 Furthermore, because the method is built using conformal prediction, the whole approach is model-agnostic. As such any possible CADRF model that provides a dose-response curve given 970 features and treatment can be used and thus is not limited to the presented CADRF approach in this 971 paper.

#### 972 D.2 <u>APPLICATIONS</u> 973

974 The classic application is in drug dosing, where the goal is to construct a dose-response curve for every individual to facilitate decision-making when determining an optimal dose for a new patient. 975 In a clinical trial, especially phase 1 and phase 2 where the optimal dose is being determined, the 976 weighted conformal dose-response curve can also act as a tool to analyse the results individually 977 while having an estimate of the uncertainty estimates that is not biased by the treatment assignment 978 distribution. It quantifies uncertainty for individual predictions, compensating for any treatment 979 distribution bias. Furthermore, it highlights areas with insufficient data support with infinite 980 prediction intervals, guiding decisions about whether further trials or treatments are necessary for 981 specific patient subgroups. In the regions where there is support, the model predictions provide the 982 CADRF estimate for this patient and the uncertainty regions show how the outcome would vary. 983

Treatment is not limited to healthcare. Treatment can be generalized as any intervention or action which opens applications in other domains. For example, in predictive maintenance, the model can optimize decisions by estimating the effect of operating pressure on the remaining useful life of equipment like valves. Similarly, in sales, it can help determine the ideal discount for specific clients to maximize the sold units, demonstrating flexibility in various domains.

989 990

D.3 EXPLAINABILITY

The application potential is also not limited to actual treatments and interventions. The method can also be used for the explainability of a model. Suppose we fitted a regression model, regressing  $X = [X_1, ..., X_m]$  on Y. X is observed data; thus, any feature can be confounded or biased. By considering a feature  $X_i$  as a treatment to "intervene" in a model, this method then provides uncertainty quantification on a Ceteris Paribus curve of a model in a similar manner to a dose-response curve<sup>2</sup>. This curve can then give unbiased uncertainty estimates of the "true" outcome for an individual sample if that sample would have had other values for this particular feature.

998 An example is shown in Figure 3 using Local Propensity WCP. This example is generated using 999 the Boston Housing data available native in sklearn (Pedregosa et al., 2011), split into a training and 1000 calibration set using a 75/25 split. A CatBoostRegressor using 300 iterations is fitted on a training 1001 set, and a propensity CatBoostRegressor with the same number of iterations is fitted on the training 1002 set. A CPS is used and calibrated on the calibration set for the propensity distribution estimate, 1003 similar to the experimental setup in this paper. No hyperparameter tuning is applied for simplicity, so note that the epistemic uncertainty could be further reduced. The chosen feature for generating a 1004 ceteris paribus curve is MedInc, the median income, an important variable in predicting the median 1005 house value in this dataset. The figure is for a single data sample where all other variables of this 1006 sample are kept constant except for our "treatment" MedInc. In Figure 3, it is apparent that the 1007 prediction intervals go to infinity for MedInc values below 1 and above 6.5. This indicates that 1008 there is insufficient overlap to evaluate this sample for these values of MedInc, clearly showing 1009 a bias in the data distribution of MedInc, given the other features. Consequently, the predictions 1010 for a sample with these features but with a MedInc of, e.g., 8 cannot be trusted as the model is 1011 simply doing an interpolation in an out-of-bounds region. In the regions with support, i.e., around 1012 1.5 < MedInc < 6.5, we see that the model shows a linear relation with the median house value 1013 with relatively small uncertainty bounds. This analysis can be done for any other regression model 1014 in a likewise manner.

1015 1016

1017

# E COMPARISON TO SCHRÖDER ET AL.

In comparison to the work of Schröder et al. (2024), our approach differs in several key aspects. First, the aim of their work is different from ours. The aim of Schröder et al. (2024) is to provide prediction intervals for the causal effect of treatment interventions where the treatment value is continuous. In our work, the goal is to provide prediction intervals for dose-response models instead of treatment interventions, answering a different causal question. However, adjusting our work to interventions

 <sup>&</sup>lt;sup>1024</sup> <sup>2</sup>A Ceteris Paribus curve visualizes a model's predictions while keeping all features constant except for one explanatory variable. The x-axis represents the explanatory variable, and the y-axis shows the corresponding predictions.



Figure 3: A Ceteris Paribus curve generated with Local Propensity WCP.

is possible; In the case of soft interventions, the target distribution propensity changes and thus 1043 substituting the current uniform distribution in the weights w(x) with the new target propensity 1044 distribution covers the soft intervention case. For hard interventions, this is an evaluation for a single 1045 treatment value which is similar to the local propensity method, but for only that target treatment 1046 value. Secondly, their approach differs in their conformal prediction approach where they want to 1047 provide correct prediction intervals for a single sample, single  $\alpha$  value, and single treatment using 1048 a mathematical solver based on the proposed weighted conformal prediction by Gibbs & Candes 1049 (2021). Thirdly, they frame the propensity or covariate shift differently as either a Dirac distribution 1050 for a hard intervention, or a different propensity distribution in the case of a soft intervention. This 1051 is a direct consequence of their aim to quantify the causal effect of a single intervention, compared to providing a dose-response model in our case which requires a uniform assumption. Fourthly, the 1052 experimental setup of Schröder et al. (2024) does not address the impact of a treatment covariate 1053 shift as shown by Figure 5 and Figure 6 where even standard conformal prediction (CP) achieves 1054 the required empirical coverage. Lastly, we also approach the propensity estimation in cases with 1055 unknown propensity as an uncertainty quantification problem and tackle it with conformal predictive 1056 systems. In the end, our approach offers a different solution on continuous treatment effects through 1057 dose-response modelling. 1058

## F ADDITIONAL RESULTS



Figure 4: Barplot of the mean coverage calculated over 40 treatment values in 50 experiments for setup 3 scenario 2. Black dotted line is the ideal coverage.

1077 1078

1059

1060

1040 1041 1042



Figure 5: Barplot of the mean coverage calculated over 40 treatment values in 50 experiments for setup 2 scenario 1. Black dotted line is the ideal coverage.



Figure 6: Barplot of the mean coverage calculated over 40 treatment values in 50 experiments for setup 2 scenario 2. Black dotted line is the ideal coverage.



Figure 7: Barplot of the mean coverage calculated over 40 treatment values in 50 experiments for setup 1 scenario 1. Black dotted line is the ideal coverage.



Figure 8: Barplot of the mean coverage calculated over 40 treatment values in 50 experiments for setup 1 scenario 2. Black dotted line is the ideal coverage.

1150 1151



Figure 9: Barplot of the mean coverage calculated over 40 treatment values in 50 experiments for setup 1 scenario 3. Black dotted line is the ideal coverage.



Figure 10: Barplot of the mean coverage calculated over 40 treatment values in 50 experiments for setup 1 scenario 4. Black dotted line is the ideal coverage.



Figure 11: Barplot of the mean coverage calculated over 40 treatment values in 50 experiments for setup 1 scenario 5. Black dotted line is the ideal coverage.



Figure 12: Barplot of the mean coverage calculated over 40 treatment values in 50 experiments for setup 1 scenario 6. Black dotted line is the ideal coverage.



Figure 13: Barplot of the mean coverage calculated over 40 treatment values in 50 experiments for setup 1 scenario 7. Black dotted line is the ideal coverage.



Figure 14: Barplot of the mean coverage calculated over 40 treatment values in 50 experiments for setup 1 scenario 8. Black dotted line is the ideal coverage.







Figure 16: Plot of the CADRF RMSE with ± RMSE standard deviation across all repeated
experiments for the considered treatment values for setup 1, scenarios 1 to 4. As All WCP and
CP methods use the same fitted base CatBoost CADRF learner they are represented by "CP and
WCP".



Figure 17: Plot of the CADRF RMSE with  $\pm$  RMSE standard deviation across all repeated experiments for the considered treatment values for setup 1, scenarios 5 to 8. As All WCP and CP methods use the same fitted base CatBoost CADRF learner they are represented by "CP and WCP".