

# Data-Centric OPE: Evaluating Off-Policy Evaluation Without an Environment

Anonymous Authors<sup>1</sup>

## Abstract

Evaluating the value of a hypothetical target policy with only a logged (offline) dataset is important but challenging. On the one hand, it brings opportunities for safe policy improvement under high stake scenarios like clinical guidelines. On the other hand, such opportunities raise a need for precise off-policy evaluation (OPE). While previous work on OPE focused on improving the algorithm in value estimation, in this work, we emphasize the importance of the offline dataset, hence putting forward a data-centric framework for OPE. We propose DataCOPE, a [data-centric framework for evaluating off-policy evaluation](#), that answers the questions of whether and to what extent we can evaluate a hypothetical policy given a dataset. DataCOPE (1) forecasts the performance of general OPE algorithms without access to the environment, which is especially useful before real-world deployment where evaluating OPE is impossible; (2) discovers the subgroup in the dataset where OPE can be inaccurate; (3) [permits environment-free evaluations or comparisons of dataset or data-collection strategies for OPE problems](#). Our empirical analysis of DataCOPE in the logged contextual bandit settings using healthcare datasets confirms its ability to evaluate both machine-learning and human-based policies, such as clinical guidelines.

## 1. Introduction

Introducing novel policies and guidelines in high-stakes settings such as healthcare and criminal justice comes with great potential harm, and should be backed by appropriate data and evidence before enacting (Woolf et al., 1999; Suresh & Guttag, 2021).

The challenge is that it can be hard to actually know when

<sup>1</sup>Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.

your data is sufficient, and without the ability to run a policy to see what *actually* happens, we must resort to estimating its effect given this logged *offline dataset* of previously seen observations and actions. Generally, this is known as **off-policy evaluation** (OPE) (Precup, 2000; Beygelzimer & Langford, 2009; Dudík et al., 2011) or the **off-policy prediction** (OPP) when an *instance-wise return prediction is emphasized* (Zhang et al., 2022; Taufiq et al., 2022), using data collected previously under one policy to predict the performance of *another* policy. While this terminology emerged in the reinforcement learning (RL) literature, it’s rooted in a causal problem and is often considered under the guise of treatment effect estimation when intervention occurs only in one time step (Powers et al., 2018). The causal nature highlights why the offline RL problem is hard - there is a limit to what can be said about a policy from *observational* data alone (Pearl, 2009).

OPE has been extensively studied both methodologically and empirically (Strehl et al., 2010; Jiang & Li, 2016) - however, this has mostly been from the perspective of improving estimators. As such, we challenge the underlying assumption that the dataset is suitable for all potential OPE methods, a neglected focus so far. While we still consider improving the estimator to an extent, we take a *data-centric* approach, looking at the problem:

Given a **dataset** and **OPE method**:

To what extent can we evaluate a **policy**?

And can we use this to see how we might improve?

Throughout the paper, we will return to the concrete example of introducing a new clinical guideline (*cf. policy*), such as the introduction of the Model for End-Stage Liver Disease (MELD) scoring in liver transplant allocation (Habib et al., 2006). These are highly impactful, committee-made decisions, for which we really must know the answers to the questions of will it be effective? and if so who will it actually benefit?

**What is needed from the community?** We require a method that can *evaluate whether an OPE method will be reliable and identify for what contexts/policies an evaluation is uncertain*. For clarity, consider the following desiderata:

1. Method-Agnostic (Data-Centric): it should be robust to change of underlying OPE methods, being able to cap-

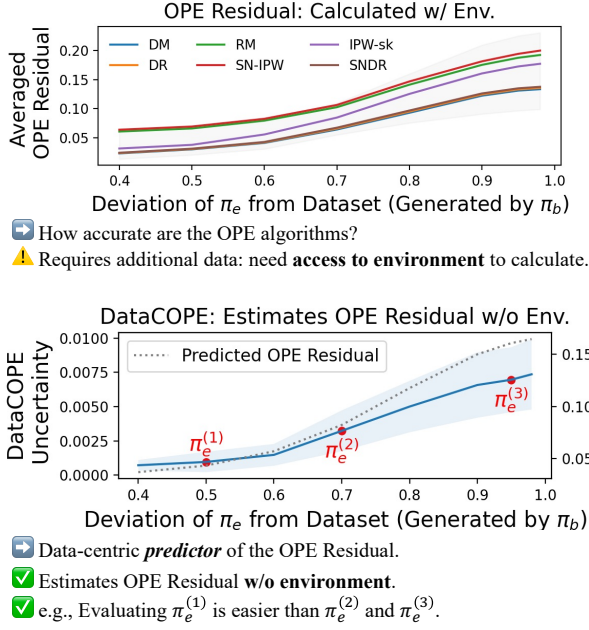


Figure 1. (1) Upper Plot: The difficulty of OPE problems is not uniform and increases as the mismatch between the behavior policy and the target policy becomes larger. This observation is a key motivation behind our work. (2) Lower Plot: evaluating OPE methods typically requires additional data, which is often impractical to obtain. It is valuable to have a proxy indicator that can predict the difficulty of OPE problems. Our proposed method, DataCOPE, can serve as a proxy and distinguish between easy and hard OPE problems. It can inform users whether a given dataset can accurately evaluate a target policy from a data-centric perspective.

ture the intrinsic difficulty of OPE problems regardless of the methods being selected.

2. **Decomposable Evaluation:** for hard-to-evaluate policies, it should be able to identify the reason (examples) that causes the difficulty in evaluation. Specifically, it will be useful if hard-to-evaluate examples can be discovered.
3. **Enables Dataset Comparison:** It should be able to identify which dataset, or data collection strategy, is more appropriate in evaluating a certain target policy. i.e., to evaluate the target policy with higher confidence and lower error.

Fulfilling this, in this work we propose Data-Centric Off-Policy Evaluation (DataCOPE), a framework that evaluates the inherent difficulty of OPE problems, and is able to predict the general performance of OPE algorithms by decomposing the estimation uncertainty in their models. DataCOPE detects dataset-target policy mismatch and thus compares data collection strategies for more accurate OPE without an environment. Figure 2 illustrates how our work develops.

Contributions of our paper is threefold:

- Methodologically, our research diverges from previous model-centric studies on Off-Policy Evaluation (OPE), which have primarily concentrated on developing algorithms. Instead, our investigation places a significant emphasis on the crucial role of data, particularly with regards to the target policy, in OPE problems. Thus, our study marks an initial attempt at implementing the principles of Data-Centric AI that is prevailing in supervised learning to address OPE issues.
- Practically, we introduce DataCOPE as an evaluation proxy for OPE problems. Traditional evaluation of OPE algorithms requires access to the true target policy value or live environment, but DataCOPE can serve as a proxy to predict whether OPE algorithms perform well in the absence of an environment.
- Empirically, we have demonstrated that DataCOPE (1) serves as an effective evaluation proxy for OPE, (2) provides a detailed performance prediction on instance-wise value estimation, and (3) can be applied to real-world datasets, such as evaluating clinical guidelines.

## 2. Preliminaries

We focus on a **logged contextual bandit** setting (Joachims et al., 2018). In particular, we consider **contexts**  $x \in \mathcal{X}$ , **actions**  $a \in \mathcal{A} := \{1, 2, \dots, k\}$ , and **rewards**  $r^a \in \mathbb{R}^+$  generated by the stochastic reward generation process taking action  $a$  given context  $x$ . In this environment, one can act according to a **policy**  $\pi \in \Pi := \Delta(\mathcal{A})^{\mathcal{X}}$ , the **value** of which is given by:

$$V(\pi) = \mathbb{E}_{x \sim p(X), a \sim \pi(x)} [r^a], \quad (1)$$

the expected reward obtained by executing the given policy. A **contextual bandit** problem then involves finding the solution to:

$$\arg \max_{\pi \in \Pi} V(\pi), \quad (2)$$

considered the **optimal policy**  $\pi^*$ . Typically this is solved via repeated interaction with the environment, which is not available when we move to the **logged** setting. In this case, we are unable to interact with the environment but alternatively have access to a **dataset**  $D = \{(x_i, a_i, r_i)\}_{i=1}^N$  of context, action, reward tuples. These have been generated via a **behavior policy**  $\pi_b \in \Pi$  that we *do not* observe but has previously interacted with the environment as follows:

1. An examples  $x_i \sim p(X)$  is drawn.
2. The behavior policy  $\pi_b$  selects action  $a_i \sim \pi_b(x_i)$ .
3. A stochastic reward  $r^a$  is observed.
4.  $(x_i, a_i, r_i)$  is added to  $D$ .

This makes the optimization task in (2) difficult as in the dataset  $A \sim \pi_b$ , which will introduce significant bias if we attempt to Monte Carlo estimate the expectation directly

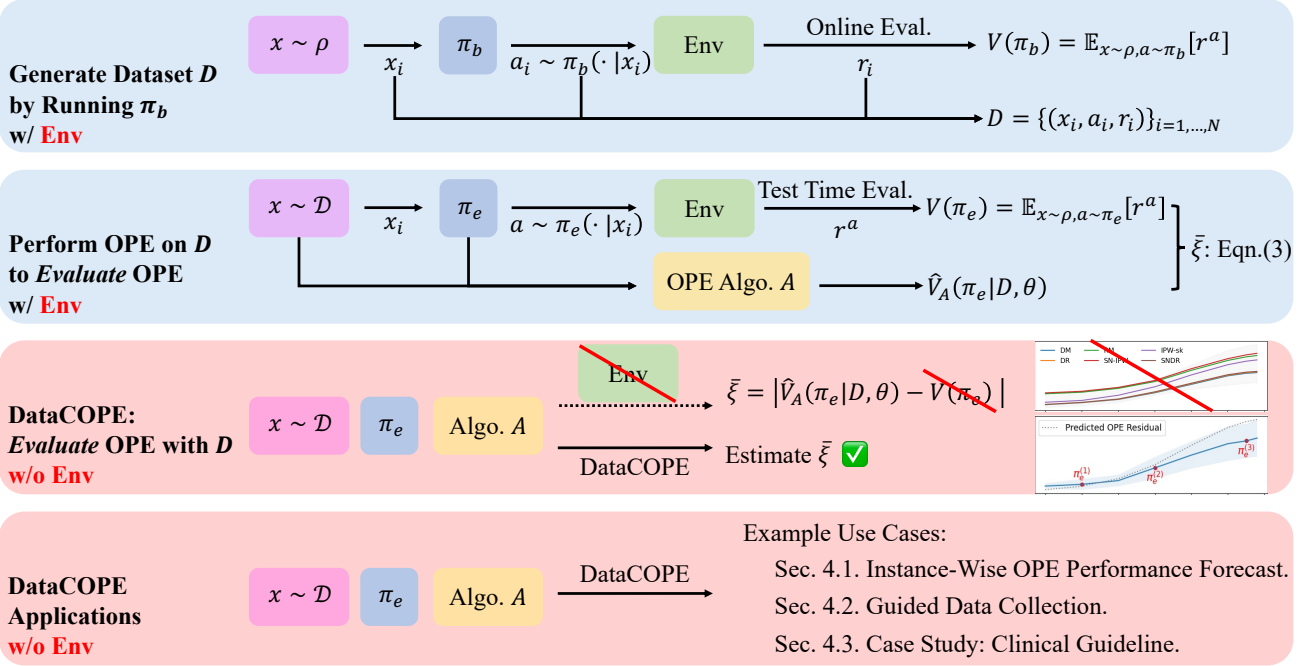


Figure 2. Road map of DataCOPE. **1-st row:** illustrates the offline dataset collection process. In the context of healthcare, it corresponds to treatment records abiding by an existing guideline. **2-nd row:** (Sec.2) The collected dataset  $\mathcal{D}$  is then used for OPE. For an OPE algorithm, Equation (3) calculates the test-time residual (error) between the value estimation result from an algorithm and the true value. **3-rd row:** (Sec.3) DataCOPE is proposed as a proxy for such an evaluation residual, that can work in test-time as an OPE performance indicator without access to the true value  $V(\pi_e)$  and the environment. Only the DataCOPE part in Fig 1 is accessible in practice. **4-th row:** (Sec.4) DataCOPE can be applied to various use cases, which is demonstrated with extensive empirical studies. **Notions are explained in detail in Section 2.**

using our given samples. Thus, a large part of the logged contextual bandit problem revolves around the accurate estimation of (1), a task referred to as **Off-Policy Evaluation** (OPE), as we wish to *evaluate* a policy using data collected *off-policy* (i.e., with a behavioral policy). Many algorithms have been proposed for such estimation, the learning objective of which is normally to minimize the mean-square error (*OPE residual*) between real and estimated values. Considering a value function estimator  $\hat{V}(\pi|D, \theta)$ , built with dataset  $D$  and parameterised by  $\theta$ , the objective is to minimise:

$$\bar{\xi} := \text{MSE}(\hat{V}) = \mathbb{E} \left[ \left( V(\pi) - \hat{V}(\pi|D, \theta) \right)^2 \right]. \quad (3)$$

**Plug-in Estimation of the Value Function.** An important sub-class of OPE algorithms, based on the **Direct Method** (DM) (Beygelzimer & Langford, 2009) revolves around constructing a directly parameterized estimator of the reward, a function  $\hat{q}_\theta$  taking contexts and actions and returning the predicted reward. This is typically learned using supervised learning based on the dataset  $D$ , such that:

$$\hat{q}_\theta = \arg \min_{\hat{q}_\theta} \mathbb{E}_{(x,a,r) \sim D} \left[ (r - \hat{q}_\theta(x, a))^2 \right]. \quad (4)$$

Armed with  $\hat{q}_\theta$ , DM estimates the predicted reward of taking actions according to the policy and plugs them into the value

function calculation, producing an estimator:

$$\hat{V}_{\text{DM}}(\pi|D, \theta) = \mathbb{E}_{x \sim D, a \sim \pi} [\hat{q}_\theta(x, a)]. \quad (5)$$

While this class of methods is by far one of the most popular (Saito et al., 2020; Fu et al., 2021) (and the one we shall build our method around), there are a number of alternatives. For example, in Inverse Probability Weighting (IPW) (Precup, 2000; Strehl et al., 2010), a weighted sum of behavior rewards is used in estimating the value of a target policy. The Doubly Robust (DR) method (Dudík et al., 2011; Jiang & Li, 2016; Su et al., 2020) improves the DM and IPW by leveraging the strength of both. Shrinkage techniques (Su et al., 2020), self-normalization (Swaminathan & Joachims, 2015), and switch method (Wang et al., 2017) further address the variance issue of those estimators. The recent advance of distribution correction estimation (DICE) family (Nachum et al., 2019; Zhang et al., 2020a;b) achieve promising performance and are unified as regularized Lagrangians of the same linear program (Yang et al., 2020). Differently, in this work, we focus on the *importance of the dataset* used for those OPE estimators, which is why we emphasize the specific dependence on  $D$  in Equation (5).

**Data-Centric AI and Data Quality.** Previous data-centric works (e.g. Data Maps (Swayamdipta et al., 2020)

and Data-IQ (Seedat et al., 2022)) evaluate the data in classification settings with the goal of characterizing examples in a dataset into easy, hard and ambiguous based on analyzing the training dynamics of individual examples. Such methods assume access to the prediction probability for the ground-truth class to compute uncertainty measures and are largely focused on curating a high-quality training dataset. However, Off-Policy Evaluation has three distinct differences: (1) we are focused on test time evaluation, where (2) the ground-truth label in value prediction is not available, therefore, the prediction confidence used in prior works is unavailable and (3) these methods are unable to tackle regression tasks, which are always needed in OPE.

**Benchmarking OPE Algorithms.** As with the problem of evaluating policies, actually evaluating the quality of OPE algorithms is similarly difficult. In recent years, the community has proposed various large-scale datasets for the purpose of benchmarking OPE algorithms, including Open-Bandit (Saito et al., 2020) and DOPE (Fu et al., 2021) benchmark OPE in the commercial recommendation and robotics settings. (Voloshin et al., 2019) empirically studies many current methods, and stress tests OPE algorithms in the RL setting. We note an important difference in health-care that there are always human-based clinical guidelines, rather than machine-learning policies.

Moreover, a central problem remains - given a new task for which OPE is required, there is no reliable way to estimate the quality of any value estimation - which brings us to our contribution.

Extended discussions on related work is elaborated in Appendix A.

### 3. Evaluating OPE through a Data-Centric Perspective

Given a dataset, target policies are not created equal for OPE. Consider an offline dataset generated by some behavior policy  $\pi_b$ , then intuitively the more similar a target policy  $\pi_e$  is to the behavior policy the easier it should be to evaluate the performance of  $\pi_e$  using the dataset generated by  $\pi_b$ . When the decisions made by  $\pi_e$  and  $\pi_b$  are similar, outcomes of the actions from  $\pi_e$  should be well represented in the dataset, while the unsupported actions' values will be challenging to predict. However, without explicit access to the behavior policy  $\pi_b$ , measuring its distance to the target policy is a highly non-trivial task.

#### 3.1. Data-Centric Difficulty Forecasting

A critical objective of this work is to provide a forecast of how accurate an OPE problem can be solved from a data-centric perspective: we emphasize the importance of *data*, rather than the algorithms. Such a perspective permits a

---

#### Algorithm 1 DataCOPE

---

**Input**

Logged Dataset  $D = \{x_i, a_i, r_i\}$ , target policy  $\pi_e$ .

**Output**

$h$  that Forecasts OPE Residual

# 1. Build Value Estimator

Optimize  $\hat{q}_\theta$  with Equation (7) for non-binary reward estimation and classifiers with logits outputs otherwise.

# 2. Uncertainty Decomposition

Decompose the uncertainties in estimating  $\hat{q}_\theta$  with Equation (8).

# 3. (Optional) Calibration

Quantify OPE residual according to Equation (9).

**Return**

Difficulty in estimating different examples

---

hierarchical analysis of the problem:

1. First and foremost, we want to provide an overall description of the *inherent difficulty* of the OPE problem. Given the offline dataset  $D$ , and target policy  $\pi_e$ , we formally write the difficulty of the OPE problem as

$$h(D, \pi_e) \propto \mathbb{E}_\theta \text{MSE}(\hat{V}_A(\pi_e|D, \theta)), \forall \hat{V}_A \quad (6)$$

where  $A$  is an arbitrary OPE algorithm and  $\theta$  denotes its instantiation.

2. Then, a method should permit a case-by-case analysis of the OPE problem, and give fine-grained explanations of the difficulties identified above.
3. In doing so, it should identify the sources of OPE difficulty. While in some cases the difficulty originates from inherent stochasticity, in some other cases it requires more diverse samples to match the target policy for an accurate OPE.

#### 3.2. Direct Method Uncertainty Decomposition

A central part of our idea for quantifying the difficulty of OPE is to quantify decomposed uncertainty and use that to build a model that is capable of predicting the OPE residual. As is common in data-centric approaches, to satisfy the previous desiderata we propose to separate the aleatoric and epistemic parts of the prediction (Seedat et al., 2022), which in our case is the value estimator. This has many benefits, especially for informing future data collection since the only relevant part here is the epistemic uncertainty. This section is organized as follows: first, we propose our tailored DM making use of a distributional reward estimator; we then introduce the uncertainty decomposition algorithm used for breaking down this estimator; finally, we present a practical method for predicting OPE difficulty using these components.



**DM with Distributional Reward Estimators** To achieve such a separation, instead of using a regular regression model that only learns an expected value, we leverage a probabilistic network to *capture the distributional information* in value estimation. When the reward model is binary (i.e., only success with +1 or failure with 0), a network with a standard softmax output can be considered to output the logits of a Bernoulli distribution. For continuous reward models, normal regression models are insufficient and so we adopt the DM with a mixture density network (MDN) (Bishop, 1994) as a reward estimator.

Specifically, an MDN predicts a set of parameters used in a mixture of Gaussian predictive distribution so as to maximize the log-likelihood of observed data. This likelihood with  $K$  Gaussian is given as

$$\mathcal{L} = \sum_{i=1}^{i=N} \mathcal{L}(r_i | x_i, a_i, \theta) = \sum_{i=1}^{i=N} \sum_{k=1}^K w_k(x_i, a_i, \theta) \times \phi(r_i | \mu_k(x_i, a_i, \theta), \sigma_k(x_i, a_i, \theta)), \quad (7)$$

where  $\theta$  denotes the parameters of networks with three branches of outputs  $\{w_k, \mu_k, \sigma_k\}_{k=1, \dots, K}$ ,  $\phi$  is the probability density function of normal distribution, and  $\sum_{i=1}^K w_i = 1$  is the normalized weight. In this case, the reward prediction, given context  $x$  and action  $a$ , is a random variable, we denote it with  $\hat{R}_x^a$ .

**Uncertainty Decomposition** As we are interested in the predictive random variable  $\hat{R}_x^a$ , we model the uncertainty of with value predictive variance:  $v(x, a) = \mathbb{V}_{\hat{R}^a | X=x, A=a}(\hat{R}^a | X=x, A=a)$ . According to the law of total variance, we can write the total uncertainty as:

$$v(x, a) = \mathbb{V}_{\Theta} \left[ \mathbb{E}_{\hat{R}^a | X=x, A=a} \left( \hat{R}^a | X=x, A=a, \Theta \right) \right] + \mathbb{E}_{\Theta} \left[ \mathbb{V}_{\hat{R}^a | X=x, A=a} \left( \hat{R}^a | X=x, A=a, \Theta \right) \right], \quad (8)$$

where  $\Theta$  is a random variable that has an empirical distribution over the set of parameters with different instantiations, i.e., ensemble. In this way, the overall uncertainty is split into two components: *epistemic uncertainty*

$$v_{\text{ep}} = \mathbb{V}_{\Theta} \left[ \mathbb{E}_{\hat{R}^a | X=x, A=a} \left( \hat{R}^a | X=x, A=a, \Theta \right) \right],$$

and *aleatoric uncertainty*

$$v_{\text{al}} = \mathbb{E}_{\Theta} \left[ \mathbb{V}_{\hat{R}^a | X=x, A=a} \left( \hat{R}^a | X=x, A=a, \Theta \right) \right].$$

Regarding the epistemic component  $v_{\text{ep}}$ , the variance originates from the model's oscillation due to a lack of training data. In this case, collecting more data could help to reduce such variance. On the other hand, in the aleatoric

component  $v_{\text{al}}$ , the variance originates from the inherent uncertainty in the data, which results in the inability to make specific predictions in some cases. In this case, collecting more informative features rather than more data could help to reduce such variance.

**Residual Prediction through Uncertainty** Now that we are equipped with a decomposition of the uncertainty, we are able to compare the evaluation confidence of target policies: naturally, policies that induce  $(x, a)$  pairs with higher epistemic and aleatoric uncertainties are harder to evaluate. While the former can be alleviated by collecting more examples in the training data, the latter is an inherent problem of the task. Moreover, we introduce a practical method to quantitatively estimate the accuracy of OPE algorithms with the help of training data: Given the uncertainty decomposition of  $(x, a)$  pairs, we are able to build a linear regression model  $h$  that takes  $v_{\text{al}}, v_{\text{ep}}$  as the inputs and outputs the OPE residual. We call this model  $h$  the calibrated hardness predictor, and fit it with a group of held-out training data from  $D$ :

$$h = \arg \min_h \mathbb{E}_{(x,a) \sim D} (\bar{\xi} - h(v_{\text{ep}}(x, a), v_{\text{al}}(x, a)))^2, \quad (9)$$

where  $\bar{\xi}$  is defined in Eqn.(3). In this way, such a model  $h$  is able to work as a hardness indicator function of the OPE problem. Pseudocode is provided in Algorithm 1. [Implementation details for the calibration step can be found in Appendix D](#)

## 4. Experiments

Recall the three desiderata (and hence goals of our method) that we established in the Introduction: we want to be able to 1) identify at an instance-wise level difficulty in evaluation, 2) [be able to compare datasets for evaluating a certain policy by comparing coverage](#), and 3) robustly evaluate arbitrary OPE algorithms. In experiments, we evaluate DataCOPE's ability to fulfill them point-by-point in Section 4.1, 4.2, and Appendix B limited by the space, alongside ablations that consider how our ability is affected by factors such as policy complexity and bias, as well as data coverage. Section 4.3 further verifies those abilities on a real-world clinical dataset — We demonstrate how DataCOPE can be applied to evaluate clinical guidelines.

**Synthetic Dataset Generation** Following standard procedures in the OPE literature (Chu et al., 2011; Li et al., 2012; Agrawal & Goyal, 2013), we adapt supervised learning datasets into example logged bandit datasets where we know the true underlying generative process. In particular, we use two medical tabular datasets, namely Breast Cancer and Diabetes (Dua & Graff, 2017) — representing both classification and regression tasks respectively

Table 1. DataCOPE is able to predict the instance-wise evaluation error of various OPE algorithms. (Reported numbers: correlation for DataCOPE rows, and for the ablation studies we report the performance difference compared with DataCOPE, **higher is better**)

Dataset	Ablations	DM	DR	RM	SNIPW	IPWsk	SNDR
Breast Cancer	DataCOPE	$0.708 \pm 0.006$	$0.899 \pm 0.001$	$0.936 \pm 0.000$	$0.936 \pm 0.000$	$0.936 \pm 0.000$	$0.901 \pm 0.001$
	w/o $v_{al}$	$-0.136 \pm 0.004$	$-0.360 \pm 0.012$	$-0.450 \pm 0.016$	$-0.450 \pm 0.016$	$-0.450 \pm 0.016$	$-0.364 \pm 0.012$
	w/o $v_{ep}$	$-0.129 \pm 0.008$	$-0.044 \pm 0.003$	$-0.022 \pm 0.002$	$-0.022 \pm 0.002$	$-0.022 \pm 0.002$	$-0.043 \pm 0.003$
	w/o Decomposition	$-0.129 \pm 0.007$	$-0.044 \pm 0.003$	$-0.022 \pm 0.002$	$-0.022 \pm 0.002$	$-0.022 \pm 0.002$	$-0.043 \pm 0.003$
Diabetes	DataCOPE	$0.743 \pm 0.017$	$0.743 \pm 0.020$	$0.745 \pm 0.020$	$0.745 \pm 0.020$	$0.744 \pm 0.020$	$0.744 \pm 0.019$
	w/o $v_{al}$	$-0.015 \pm 0.009$	$-0.014 \pm 0.008$	$-0.013 \pm 0.008$	$-0.013 \pm 0.008$	$-0.013 \pm 0.008$	$-0.013 \pm 0.008$
	w/o $v_{ep}$	$-0.357 \pm 0.022$	$-0.365 \pm 0.018$	$-0.371 \pm 0.020$	$-0.371 \pm 0.020$	$-0.370 \pm 0.020$	$-0.370 \pm 0.020$
	w/o Decomposition	$-0.357 \pm 0.022$	$-0.364 \pm 0.018$	$-0.370 \pm 0.020$	$-0.370 \pm 0.020$	$-0.370 \pm 0.020$	$-0.369 \pm 0.020$

— in order to validate that DataCOPE is effective and powerful. Further validation on other common UCI datasets is provided in Appendix E to verify the generalization of DataCOPE.

We use linear models for the behavior policy  $\pi_b$ , which is trained on the dataset examples  $\{(x_i, y_i)\}_{i=1}^N$ . Given context  $x_i$  with corresponding labels  $y_i$  we learn a policy to generate actions by minimizing the expected negative log-likelihood

$$\mathbb{E}_{x,y} [\text{NLL}(\pi_b(x_i), y_i)], \quad (10)$$

under a predictive Gaussian/Bernoulli distribution for regression/classification respectively.

We then evaluate  $(x_i, a_i)$  with the help of ground-truth labels  $y_i$ . For classification tasks, the reward of action  $a$  is given as  $r_i^a = \mathbf{1}(a_i = y_i)$ , where  $\mathbf{1}$  is the indicator function; while for regression tasks, the reward of action  $a_i$  is given by the coefficient of determination of the prediction. i.e.,  $r_i^a = 1 - \frac{u}{v}$ , where  $u$  is the residual sum of squares  $u = \|y - a\|_2^2$  and  $v$  is the total sum of squares  $v = \|y - \bar{y}\|_2^2$ ,  $\bar{y}$  denotes the mean of  $y$ . In this way, we can generate  $\{(x_i, a_i, r_i^a)\}_{i=1}^N$  tuples as our dataset containing  $N$  examples.

**Target Policy** We also generate target policies using neural network models trained on the *same* training data  $\{(x_i, y_i)\}_{i=1}^N$ , but with injected noise, that will depend on the particular experiment, as  $\pi_e$ . The ground-truth performance of policy  $\pi_e$  is then given by  $r_j^b = \mathbf{1}(b_j = y_j)$ , where  $b_j = \pi_e(x_j)$  - importantly this is available to us and thus allows for the evaluation of the quality of the OPE.

In particular, the off-policy evaluation problem is by definition estimating  $\mathbb{E}[r_j^b]$ , for  $j = 1, \dots, N$ . As we have access to  $y_j, j = 1, \dots, N$ , we can quantitatively evaluate our method and compare results with the ground-true policy values.

**General Experiment Settings** We conducted all our experiments using 8 random seeds and reported the mean and standard deviation of the results. The observed performance differences with 8 trials were found to be statistically significant. For the ensemble-based uncertainty quantification, we

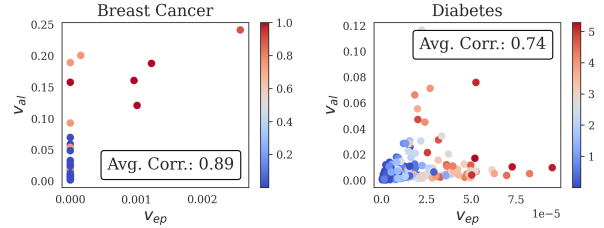


Figure 3. DataCOPE decomposes the uncertainty and provides an instance-wise prediction of the estimation error (Shallow color indicates larger residual value). We find in our experiments that such a prediction is highly correlated with the instance-wise OPE residual: examples with high uncertainty always have high-value estimation errors. The result holds with all datasets and various OPE algorithms.

determined that using 100 models was computationally feasible for tabular datasets in the healthcare domain. However, we also found that using 10 models produced satisfactory results.

To maintain consistency with prior literature, we differentiated between the behavior policy and target policy in two distinct ways. The first approach involved injecting random noise into the labels or regression target during the generation of behavior policies, thereby ensuring that the behavior trajectories (i.e., dataset) did not completely conform to the target policy behavior. The second, slightly more advanced method involved biasing the policy through dataset sculpture, as described in Section 4.2. The aim of both approaches was to evaluate the performance of DataCOPE on various datasets with differing degrees of mismatch between the behavior and target policies. As shown in the upper plot of Figure 1, we observed that an increase in mismatch resulted in a higher OPE error for all algorithms.

#### 4.1. Using DataCOPE for Instance-Wise Difficulty Indication

In this section, we zoom in on the individuals who will be acted on by the policies and use DataCOPE to predict instance-wise OPE difficulty. We provide an analysis at the algorithm level, asking how well will an OPE method do *on average* in Appendix B.

**Experimental Setup** For each dataset, we run different OPE algorithms and are able to calculate the value estimation residual (error) over every datum, individually comparing this to the output of DataCOPE, which has decomposed the uncertainties into aleatoric and epistemic components. The correlation between the OPE residuals and the uncertainties is then calculated and reported, noting that a *high* correlation implies that the DataCOPE uncertainties have *predictive power* when it comes to determining if the OPE is accurate.

**Results** As shown in Table 1, DataCOPE effectively predicts the OPE residual, but we can see the effect on an individual level more clearly in the scatter plot of Figure 3 which qualitatively shows the relation between two uncertainties and the averaged OPE performance in terms of residual over 6 algorithms. Normally, the examples whose value can not be precisely estimated are located at the upper right in the scatter plots, meaning they either have high aleatoric or epistemic uncertainty.

Our ablation studies also serve to demonstrate how the decomposition is important: In *Diabetes*, the epistemic uncertainty is more important than the aleatoric component in predicting the OPE residual, indicating that such difficulty can potentially be alleviated by collecting more data. While in *Breast Cancer*, we find aleatoric uncertainty dominates the performance prediction. This is important to be aware of since the aleatoric dominance in *Breast Cancer* suggests there is limited opportunity in the future to be able to reduce this uncertainty by collecting additional data - which is not the case at all for *Diabetes*. In both cases, the decomposition is vitally important as it manifests the uncertainty components even at a different magnitude.

**Take-away:** *DataCOPE is able to predict instance-wise difficulty in OPE for both discrete and continuous tasks. The uncertainty decomposition step in DataCOPE is important in isolating the uncertainty components' effect even when they are of different magnitudes.*

#### 4.2. Matching Dataset and Target Policy with DataCOPE

So far we have shown that DataCOPE can tell if our OPE estimates will be any good, and for which individuals this will apply. We now move to show how the uncertainty components discovered by DataCOPE can be informative in comparing datasets for evaluating a certain target policy.

We consider the scenario of the clinical setting where new policies need to be evaluated before verifying those policies in clinical trials. We show that DataCOPE can act as a performance indicator and be informative in the pursuit of efficient and accurate OPE.

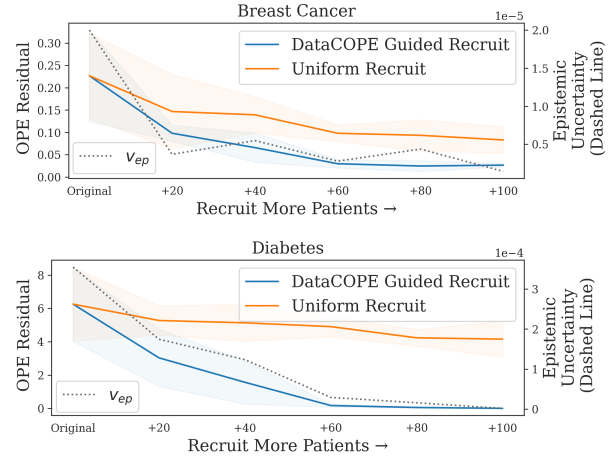


Figure 4. Adding examples that have high epistemic uncertainty in the learning of  $\pi_b$  minimizes epistemic uncertainty and clearly improves OPE performance. While uniformly sampling new examples only improves OPE a little. DataCOPE identifies which dataset or data collection strategy is better for evaluating a given target policy.

**Experimental Setup** We consider the setting where different data collection strategies are available, such that we can choose from different data collection strategies for OPE. e.g., a medical guideline is still not optimal and exploration can be done to try to improve it.

Here, we synthetically generate a biased logged dataset by running biased and highly deterministic behavior policies  $\pi_b$  for collecting those data. Specifically, we explore the *coverage* of the behavioral policy by thresholding examples by their quantile values on certain features. We remove examples with the top 60% high values at the feature of the *worst concave point* in the *Breast Cancer* dataset and remove examples with top 60% high values of the feature *log of serum triglycerides level*<sup>1</sup> during behavior policy generation. In the meantime, the target policy  $\pi_e$  is still generated with the full dataset, and hence there exists a behavior bias on the training examples out of the threshold.

We demonstrate the effectiveness of DataCOPE by experimenting with the two most straightforward data collection strategies: uniform sampling (e.g., uniformly recruiting patients in clinical trials) and uncertainty-guided sampling — according to the principle of *optimism in the face of uncertainty* (Kearns & Singh, 2002; Brafman & Tennenholtz, 2002) (i.e. picking those with the highest uncertainty).

**Results** In Figure 4, we report the changes in OPE residuals as newly recruited patients are added. Therefore, the aligned x-axis indicates the size of the dataset is the same — yet their quality can be different — due to different data

<sup>1</sup>these are selected as the most influential dimension for linear models to manifest the difference

collection strategies and discrepancies from the behavior policy. We report the epistemic uncertainty in OPE when recruitment is guided by DataCOPE. Curves are get by averaging over 6 benchmark algorithms. For both datasets, we find DataCOPE is able to identify the better data collection strategy according to the epistemic uncertainty.

**Take-away:** *DataCOPE is able to identify the target policy  $\pi_e$ 's deviation from the behavior policy. The epistemic uncertainty component of DataCOPE informs the inherent difficulty of OPE problem in general as well as in sub-groups. The quality of the datasets for a specific target policy evaluation task can be compared according to DataCOPE.*

### 4.3. Case Study: The Introduction of MELD

We apply DataCOPE to a real-world healthcare guideline and examine its potential to benefit high-stakes OPE problems, with a specific focus on evaluating organ transplant allocation policy.

**Background** In the medical field, the official guidance on organ transplantation and which potential recipients are offered organs when they become available has evolved multiple times over the past few decades (Starzl et al., 1982; Adam et al., 2012; Hüyük et al., 2022). This setting can be seen as a contextual bandit problem where every arriving patient and organ feature instance is considered the context, while the policy makes an allocation decision as the action, with the corresponding survival time for the patient perceived as an internal reward. We examine data from the Organ Procurement & Transplant Network (OPTN) (Leppke et al., 2013), which includes information on patients registered for liver transplants from 1995 to 2020. More details on data processing can be found in Appendix D.

We specifically focus on the deployment of the prevailing allocation policy — MELD (Bernardi et al., 2011), and study the OPE problem of MELD before its deployment time, 2002, to show that DataCOPE:

1. Has a strong correlation with the OPE residual;
2. Predicts and explains when and why additional collected data improves OPE;
3. Identifies vulnerable sub-groups of patients that are more likely to suffer from in-accurate OPE estimation, hence potentially having high risk in medical practice.

This is important to consider whether there actually was sufficient evidence at point of deployment to justify if the guidance would be - or if there were some (sub-)groups who might be left worse off.

**Experiment Setting** We consider the behavior policy  $\pi_b$  to be that which was deployed before 2002. In order to approximate the process of selecting the most in-need patient,

Table 2. DataCOPE strongly correlates with the OPE residual on the Organ Transplant dataset.

DataCOPE	w/o $v_{ep}$	w/o $v_{al}$	w/o Decomposition
$0.712 \pm 0.069$	$0.654 \pm 0.044$	$0.309 \pm 0.198$	$0.654 \pm 0.044$

for each patient who receives an organ at some time point, we add 9 other contemporary patients not being allocated to the organ as reference examples.  $\pi_b$  then identifies the selected patient out of the 10 patients — with a discrete action, and the corresponding reward is given by the survival time of the selected patient after receiving the donated liver.

This is then used to estimate the value of the MELD policy, the  $\pi_e$  in the context of this OPE problem, that is deployed after 2002. Similar to the training dataset, we collect test-time data using the transplant records between 2003 and 2005. As MELD was applied as the allocation policy during this period, we have the real value of  $\pi_e$  — the average survival time of the patients who receive organs during this period. We apply the direct method as a demonstrative OPE solver in experiments considering it does not require a parameterized policy.

**Results: Correlation.** We calculate the Pearson R between uncertainty components calculated by DataCOPE and the OPE residual reported in Table 2 shows DataCOPE is highly correlated with OPE performance. And our ablation studies demonstrate again the importance of uncertainty decomposition.

**Results: Evaluating MELD Evaluation over time.** By utilizing DataCOPE, we are able to examine snapshots of off-policy datasets taken at different points in time, analyze the uncertainties, and compare them with the residuals resulting from OPE, as illustrated in Figure 5.

1. Before the year 2000, the Institute of Medicine allocation recommendations (*IoM Alloc. Rec.*) were used as the behavior policy. Collecting more data during this period can improve the evaluation of the MELD policy. By comparing the results provided by DataCOPE before and at 2000, we can see that the epistemic uncertainty does not change significantly, while the aleatoric uncertainty decreases. This suggests that the OPE task should become easier, which is verified by the reduction in OPE residuals at 2000.

2. After the year 2000, the *OPTN Final Rule* was implemented as the allocation policy. This change in policy affected the pattern of collected data and subsequently, the performance of OPE for MELD. During the period when OPTN was in operation (2000-2002), the aleatoric component of uncertainty decreased while the epistemic component increased. This reminds us to be cautious when using data from this period to evaluate MELD. In fact, the OPE residual using data from this period increases.



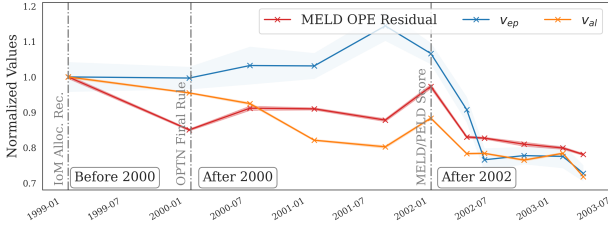


Figure 5. Evaluating MELD over time. DataCOPE can be applied to monitor the OPE performance without the real policy value.

3. After 2002, the *MELD* guideline was implemented as the allocation policy, replacing *OPTN*. As a result, the data collected since 2002 is unbiased. Both uncertainty components decreased significantly and so did the OPE residual when looking back in hindsight.

To summarize, our study on the Organ Transplant dataset supports our findings in Section 4.2 regarding data collection. Specifically, it shows that in order to achieve more accurate OPE results, efforts should be made to reduce the epistemic uncertainty as identified by DataCOPE through targeted recruitment of new patients.

#### Results: Vulnerable Sub-Groups Identification with $v_{ep}$ .

In this section, we use DataCOPE to investigate the sub-group with high epistemic uncertainty in OPE, as improving the performance of this sub-group is possible through collecting more data. We examine the examples with the highest epistemic uncertainty in OPE, and found those selected patients in general have a larger *Weight Difference*.

To manifest how the performance of OPE evolves over such a sub-group, we compute the averaged OPE residual on the sub-group and compare it with the overall performance. To be specific, we choose to use the data snapshots from the years 2000 (*IoM*), 2002 (*OPTN*), and 2003 (*MELD*) respectively to contrast the effect induced by data collected by different behavior policies.

Results are presented in Table 3. In addition to the epistemic uncertainty and OPE residuals of both sub-group and in population, we additionally calculate the ratio to highlight the vulnerability of the sub-group. We find the inaccuracy of OPE in this sub-group remains clearly higher than average until at least 2003, a year or so after *MELD* was introduced. DataCOPE successfully identifies and monitors such vulnerability with the epistemic uncertainty component, and would have been able to highlight this uncertainty to clinicians implementing the policy, allowing them to take more care in decisions for this sub-group.

Table 3. Results on vulnerable sub-group identification. DataCOPE discovers the vulnerable group in OPE problem of *MELD*, and can be used to forecast the performance.

	Data	Sub-Group	Overall	Ratio
$v_{ep}(10^{-5})$	2000 ( <i>IoM</i> )	$2.261 \pm 0.124$	$1.252 \pm 0.050$	$\times 1.808$
	2002 ( <i>OPTN</i> )	$2.238 \pm 0.146$	$1.358 \pm 0.054$	$\times 1.648$
	2003 ( <i>MELD</i> )	$1.443 \pm 0.081$	$1.151 \pm 0.048$	$\times 1.254$
$\xi(10^{-2})$	2000 ( <i>IoM</i> )	$8.869 \pm 0.116$	$6.625 \pm 0.022$	$\times 1.339$
	2002 ( <i>OPTN</i> )	$7.885 \pm 0.125$	$6.838 \pm 0.023$	$\times 1.153$
	2003 ( <i>MELD</i> )	$6.351 \pm 0.043$	$6.098 \pm 0.022$	$\times 1.041$

## 5. Conclusion

In this work, we propose DataCOPE to tackle the problem of off-policy evaluation (OPE) which is widely applicable in high-stakes real-world tasks like healthcare. DataCOPE serves as an effective proxy for the true OPE residual without the need for direct access to the environment. While previous work has focused on learning estimators from offline datasets without considering the challenges inherent in the mismatch between such datasets and the target policy being evaluated, we propose to re-think the problem of OPE from a data-centric perspective and demonstrate the benefits of such a perspective with the proposed method DataCOPE through empirical studies in both synthetic and real-world healthcare datasets.

## References

- Adam, R., Karam, V., Delvart, V., O'Grady, J., Mirza, D., Klempnauer, J., Castaing, D., Neuhaus, P., Jamieson, N., Salizzoni, M., et al. Evolution of indications and results of liver transplantation in europe. a report from the european liver transplant registry (eltr). *Journal of hepatology*, 57(3):675–688, 2012.
- Agrawal, S. and Goyal, N. Thompson sampling for contextual bandits with linear payoffs. In *International conference on machine learning*, pp. 127–135. PMLR, 2013.
- Bernardi, M., Gitto, S., and Biselli, M. The meld score in patients awaiting liver transplant: strengths and weaknesses. *Journal of hepatology*, 54(6):1297–1306, 2011.
- Beygelzimer, A. and Langford, J. The offset tree for learning with partial labels. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 129–138, 2009.
- Bishop, C. M. Mixture density networks. 1994.
- Brafman, R. I. and Tennenholtz, M. R-max-a general polynomial time algorithm for near-optimal reinforcement learning. *Journal of Machine Learning Research*, 3(Oct): 213–231, 2002.
- Chu, W., Li, L., Reyzin, L., and Schapire, R. Contextual bandits with linear payoff functions. In *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*, pp. 208–214. JMLR Workshop and Conference Proceedings, 2011.
- Dua, D. and Graff, C. UCI machine learning repository, 2017. URL <http://archive.ics.uci.edu/ml>.
- Dudík, M., Langford, J., and Li, L. Doubly robust policy evaluation and learning. *arXiv preprint arXiv:1103.4601*, 2011.
- Fu, J., Norouzi, M., Nachum, O., Tucker, G., Wang, Z., Novikov, A., Yang, M., Zhang, M. R., Chen, Y., Kumar, A., et al. Benchmarks for deep off-policy evaluation. *arXiv preprint arXiv:2103.16596*, 2021.
- Habib, S., Berk, B., Chang, C.-C. H., Demetris, A. J., Fontes, P., Dvorchik, I., Egthesad, B., Marcos, A., and Shakil, A. O. Meld and prediction of post-liver transplantation survival. *Liver transplantation*, 12(3):440–447, 2006.
- Houlsby, N., Huszár, F., Ghahramani, Z., and Lengyel, M. Bayesian active learning for classification and preference learning. *arXiv preprint arXiv:1112.5745*, 2011.
- Hüyük, A., Jarrett, D., and van der Schaar, M. Inverse contextual bandits: Learning how behavior evolves over time. In *International Conference on Machine Learning*, pp. 9506–9524. PMLR, 2022.
- Jiang, N. and Li, L. Doubly robust off-policy value evaluation for reinforcement learning. In *International Conference on Machine Learning*, pp. 652–661. PMLR, 2016.
- Joachims, T., Swaminathan, A., and De Rijke, M. Deep learning with logged bandit feedback. In *International Conference on Learning Representations*, 2018.
- Kearns, M. and Singh, S. Near-optimal reinforcement learning in polynomial time. *Machine learning*, 49:209–232, 2002.
- Kirsch, A., Van Amersfoort, J., and Gal, Y. Batchbald: Efficient and diverse batch acquisition for deep bayesian active learning. *Advances in neural information processing systems*, 32, 2019.
- Leppke, S., Leighton, T., Zaun, D., Chen, S.-C., Skeans, M., Israni, A. K., Snyder, J. J., and Kasiske, B. L. Scientific registry of transplant recipients: collecting, analyzing, and reporting data on transplantation in the united states. *Transplantation Reviews*, 27(2):50–56, 2013.
- Lewis, D. D. A sequential algorithm for training text classifiers: Corrigendum and additional data. In *Acm Sigir Forum*, volume 29, pp. 13–19. ACM New York, NY, USA, 1995.
- Li, L., Chu, W., Langford, J., and Schapire, R. Contextual-bandit approach to personalized news article recommendation, January 19 2012. US Patent App. 12/836,188.
- Musmann, S. and Liang, P. On the relationship between data efficiency and error for uncertainty sampling. In *International Conference on Machine Learning*, pp. 3674–3682. PMLR, 2018.
- Nachum, O., Chow, Y., Dai, B., and Li, L. Dualdice: Behavior-agnostic estimation of discounted stationary distribution corrections. *Advances in Neural Information Processing Systems*, 32, 2019.
- Nguyen, V.-L., Shaker, M. H., and Hüllermeier, E. How to measure uncertainty in uncertainty sampling for active learning. *Machine Learning*, 111(1):89–122, 2022.
- Pearl, J. *Causality*. Cambridge university press, 2009.
- Powers, S., Qian, J., Jung, K., Schuler, A., Shah, N. H., Hastie, T., and Tibshirani, R. Some methods for heterogeneous treatment effect estimation in high dimensions. *Statistics in medicine*, 37(11):1767–1787, 2018.
- Precup, D. Eligibility traces for off-policy policy evaluation. *Computer Science Department Faculty Publication Series*, pp. 80, 2000.

- Saito, Y., Aihara, S., Matsutani, M., and Narita, Y. Open bandit dataset and pipeline: Towards realistic and reproducible off-policy evaluation. *arXiv preprint arXiv:2008.07146*, 2020.
- Sebastiani, P. and Wynn, H. P. Maximum entropy sampling and optimal bayesian experimental design. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 62(1):145–157, 2000.
- Seedat, N., Crabbé, J., Bica, I., and van der Schaar, M. Dataiq: Characterizing subgroups with heterogeneous outcomes in tabular data. *arXiv preprint arXiv:2210.13043*, 2022.
- Starzl, T. E., Iwatsuki, S., Van Thiel, D. H., Gartner, J. C., Zitelli, B. J., Malatack, J. J., Schade, R. R., Shaw Jr, B. W., Hakala, T. R., Rosenthal, J. T., et al. Evolution of liver transplantation. *Hepatology (Baltimore, Md.)*, 2(5): 614, 1982.
- Strehl, A., Langford, J., Li, L., and Kakade, S. M. Learning from logged implicit exploration data. *Advances in neural information processing systems*, 23, 2010.
- Su, Y., Dimakopoulou, M., Krishnamurthy, A., and Dudík, M. Doubly robust off-policy evaluation with shrinkage. In *International Conference on Machine Learning*, pp. 9167–9176. PMLR, 2020.
- Suresh, H. and Guttag, J. A framework for understanding sources of harm throughout the machine learning life cycle. In *Equity and access in algorithms, mechanisms, and optimization*, pp. 1–9. 2021.
- Swaminathan, A. and Joachims, T. The self-normalized estimator for counterfactual learning. *advances in neural information processing systems*, 28, 2015.
- Swayamdipta, S., Schwartz, R., Lourie, N., Wang, Y., Hajishirzi, H., Smith, N. A., and Choi, Y. Dataset cartography: Mapping and diagnosing datasets with training dynamics. *arXiv preprint arXiv:2009.10795*, 2020.
- Taufiq, M. F., Ton, J.-F., Cornish, R., Teh, Y. W., and Doucet, A. Conformal off-policy prediction in contextual bandits. *arXiv preprint arXiv:2206.04405*, 2022.
- Thomas, P., Theodorou, G., and Ghavamzadeh, M. High-confidence off-policy evaluation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 29, 2015.
- Tucker, A. D. and Joachims, T. Variance-optimal augmentation logging for counterfactual evaluation in contextual bandits. *arXiv preprint arXiv:2202.01721*, 2022.
- Udagawa, T., Kiyohara, H., Narita, Y., Saito, Y., and Tateno, K. Policy-adaptive estimator selection for off-policy evaluation. *arXiv preprint arXiv:2211.13904*, 2022.
- Voloshin, C., Le, H. M., Jiang, N., and Yue, Y. Empirical study of off-policy policy evaluation for reinforcement learning. *arXiv preprint arXiv:1911.06854*, 2019.
- Wan, R., Kveton, B., and Song, R. Safe exploration for efficient policy evaluation and comparison. In *International Conference on Machine Learning*, pp. 22491–22511. PMLR, 2022.
- Wang, Y.-X., Agarwal, A., and Dudík, M. Optimal and adaptive off-policy evaluation in contextual bandits. In *International Conference on Machine Learning*, pp. 3589–3597. PMLR, 2017.
- Woolf, S. H., Grol, R., Hutchinson, A., Eccles, M., and Grimshaw, J. Potential benefits, limitations, and harms of clinical guidelines. *Bmj*, 318(7182):527–530, 1999.
- Yang, M., Nachum, O., Dai, B., Li, L., and Schuurmans, D. Off-policy evaluation via the regularized lagrangian. *Advances in Neural Information Processing Systems*, 33: 6551–6561, 2020.
- Zhang, R., Dai, B., Li, L., and Schuurmans, D. Gendice: Generalized offline estimation of stationary values. *arXiv preprint arXiv:2002.09072*, 2020a.
- Zhang, S., Liu, B., and Whiteson, S. Gradientdice: Rethinking generalized offline estimation of stationary values. In *International Conference on Machine Learning*, pp. 11194–11203. PMLR, 2020b.
- Zhang, Y., Shi, C., and Luo, S. Conformal off-policy prediction. *arXiv preprint arXiv:2206.06711*, 2022.

## A. Extended Related Works

### A.1. High-Level Difference: DataCOPE is a Framework for Evaluating OPE, rather than an OPE Estimator.

Although DataCOPE is situated within the general field of OPE, it differs from typical OPE papers that propose methods for solving existing OPE problems. Instead, DataCOPE functions as a type of Meta-OPE that “COPE with the Data” — it evaluates general OPE algorithms by forecasting whether the conditional expected return can be estimated accurately.

It is important to note that not all OPE problems are equal, as some can be intrinsically difficult, while others may be much easier. In our work, we focus on evaluating OPE problems (i.e., a dataset-target policy pair) rather than proposing or evaluating OPE algorithms. We present DataCOPE as a practical method and a first step toward quantifying the difficulty of OPE problems. This difficulty is an inherent quantity that relates to the mismatch between the behavior dataset and the target policy. We are motivated to introduce DataCOPE as a data-centric solution that captures the properties of the dataset and target policy.

### A.2. OPE Estimators

We provide a further discussion of the OPE literature in this section. Specifically, we connect and differentiate DataCOPE with (Tucker & Joachims, 2022; Wan et al., 2022; Jiang & Li, 2016; Thomas et al., 2015; Taufiq et al., 2022; Zhang et al., 2022; Udagawa et al., 2022).

One of the fundamental differences between our work and the literature above is that DataCOPE challenges the OPE problem itself, rather than aiming to optimize for a better OPE estimator. Our perspective is to question whether evaluating a specific target policy is reasonable, given the current dataset. This data-centric perspective is a significant contribution to the field. Instead of developing OPE algorithms without considering the data quality, we propose to evaluate whether a particular target policy can be appropriately evaluated using the available dataset. Not all target policies can be evaluated with the same level of confidence, and some may be harder to evaluate than others. DataCOPE aims to identify the easy problems and subsets from the hard ones, rather than focusing on developing OPE algorithms that perform well on some problems but not others.

To elaborate on the connections and differences between our work and the above literature in more detail:

We emphasize the importance of the data-centric perspective of OPE throughout the paper, and our contribution is not a method for data collection that treats other approaches as baselines under specific settings (budget, safety, etc.). Instead, we provide insight into the importance of a match between the dataset and the target policy. This perspective is supported by the results of (Tucker & Joachims, 2022; Wan et al., 2022), as they implicitly reveal the importance of data quality, which we focus on explicitly.

Additionally, the settings and assumptions in these works are different. Tucker & Joachims (2022) considers the problem of counterfactual predictions that is often the case in advertising. Wan et al. (2022) explores efficient data collection under safety constraints. In our work, data collection is not an option with freedom in the high-stake clinical setting. We apply DataCOPE as a demonstrative example to show how data quality affects OPE performance and use it as a proxy indicator. Our work complements this line of literature by quantifying the quality of the dataset for a given policy in terms of uncertainty.

Jiang & Li (2016) extends the doubly robust method to the sequential setting. Its counterpart in the contextual bandit setting, DR, is already used as a backbone method in verifying the universal predictive power of DataCOPE.

The work by (Thomas et al., 2015) estimates a lower bound for OPE but relies on the assumption that the behavior policy adequately covers the action space, which can not hold in high-stakes applications like clinical guidelines. In contrast, our work focuses on the problem of evaluating a target policy with an imperfect dataset, and aims to identify which policies can be accurately evaluated on which subset of the data. Thus, our approach offers a different perspective on the problem of OPE.

The works (Taufiq et al., 2022; Zhang et al., 2022) use conformal prediction techniques to improve the accuracy and reliability of OPE estimates, but they still belong to the OPE algorithm class (algorithm-centric). In contrast, DataCOPE is not an OPE algorithm but rather an evaluation method that provides a data-centric perspective on the feasibility of OPE for a given target policy and dataset. It evaluates the quality of the data and the assumptions underlying OPE, which can inform the development of OPE algorithms as well as the interpretation of their results.



Udagawa et al. (2022) is a concurrent work that examines the issue of estimator selection. In contrast, our work places greater emphasis on the notion that for certain ill-defined problems, such estimator selection may be of little help since all estimators would perform similarly poorly in comparison to when they are matched with a more suitable dataset-target policy. This highlights the fundamental difference between a data-centric approach and an algorithm-centric approach to tackling OPE problems.

### A.3. Exploration

While exploration aims to minimize long-term regret, our work is focused on minimizing the error in off-policy policy value estimation. To illustrate this, consider a clinical trial where the purpose of introducing new treatments (i.e., exploration) is to improve long-term feedback through exploration. In contrast, our work focuses on answering the question of whether, given a past or future dataset of patient information, we can minimize the estimation error of an introduced new treatment. Therefore, we believe that our work can provide valuable insights into improving the accuracy of policy value estimation and complement the exploration strategies in the literature.

### A.4. Active Learning and Active Data Collection

DataCOPE deals with a fundamentally different task from active learning (Lewis, 1995; Sebastiani & Wynn, 2000; Musmann & Liang, 2018; Houlisby et al., 2011; Kirsch et al., 2019; Nguyen et al., 2022). The goal of active learning-based uncertainty sampling is to improve the model by selecting which samples to best label next (model-centric task). This is different from DataCOPE which is focused on an already given dataset, with the goal of characterizing the types of samples and their effect on OPE. In addition, the cost functions in acquiring new samples are always very clearly defined in active learning, yet our setting does not set explicitly trade-offs.

Our proposed metric could be extended for the purposes of data acquisition. The active learning literature has explored two main paradigms, Maximum Entropy Sampling (MES) or Uncertainty Sampling, and Bayesian Active Learning by Disagreement (BALD), which aim to query samples with the maximum predictive entropy or maximize the mutual information between the observation and the model parameters, respectively. These approaches reflect model uncertainty and are linked to our notion of epistemic uncertainty, albeit measured differently. Our work could inspire future research on how to leverage our proposed metric for more principled data acquisition in OPE settings.

Different from previous active learning works or OPE data collection methods like (Tucker & Joachims, 2022), DataCOPE is used for evaluating datasets at the same size yet with different samples to demonstrate the importance of the match between dataset and target policy in OPE problems. Those different samples are collected by different strategies: the baseline one is collected through uniform sampling, while the other one is collected according to the epistemic uncertainty.

In our experiments, we are not actively recruiting new data points that have the highest epistemic uncertainty, nor can we be able to add new features to the task to decrease aleatoric uncertainty. What we intended to emphasize in this section is that DataCOPE is a hindsight descriptive tool that compares the quality of a dataset given a certain target policy.

DataCOPE is distinct from prior active learning literature or OPE data collection methods such as (Tucker & Joachims, 2022) in that DataCOPE assesses datasets of the same size but with varying samples to emphasize the significance of matching the dataset with the target policy in OPE problems. Different sampling strategies are employed to collect these diverse samples in Section 4.2: the baseline approach employs uniform sampling, while the other method leverages epistemic uncertainty to collect samples. The takeaway message is that: datasets are not equal in evaluating the same target policy. A lower epistemic uncertainty leads to a generally improved performance among a batch of OPE algorithms. There is no doubt that advanced data-collecting algorithms like (Tucker & Joachims, 2022) can potentially be more efficient, though might not be practical in high-stake scenarios.

To sum up the difference: in our experiments, we do not actively recruit new data points that have the highest epistemic uncertainty, nor can we add new features to the task to reduce aleatoric uncertainty. The purpose of this section is to emphasize that DataCOPE is a retrospective descriptive tool that compares the quality of a dataset given a specific target policy.

## B. Using DataCOPE for Evaluating OPE Algorithms

We quantitatively demonstrate that the uncertainty components from DataCOPE can be used to estimate the OPE residuals with calibration across a range of OPE methods and therefore could be broadly used as a performance indicator when the ground truth is unknown.

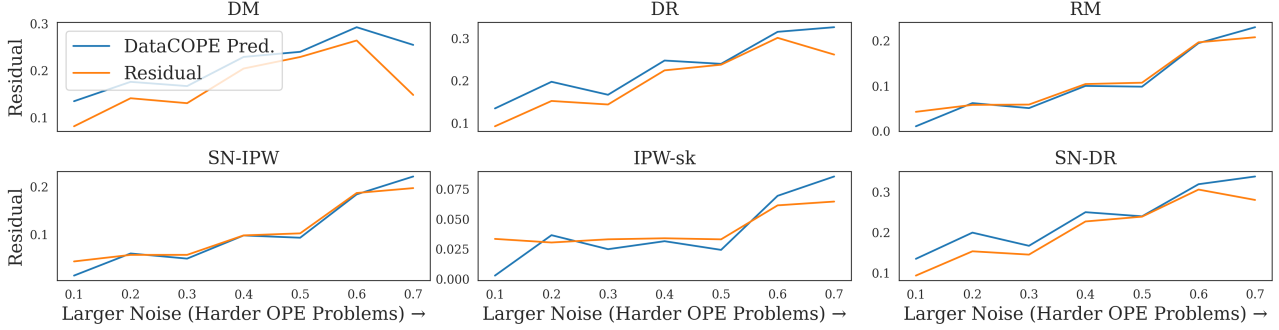


Figure 6. DataCOPE works as an accurate proxy of the OPE residual. DataCOPE is able to predict OPE residuals of different algorithms with calibration. Dataset: Breast Cancer.

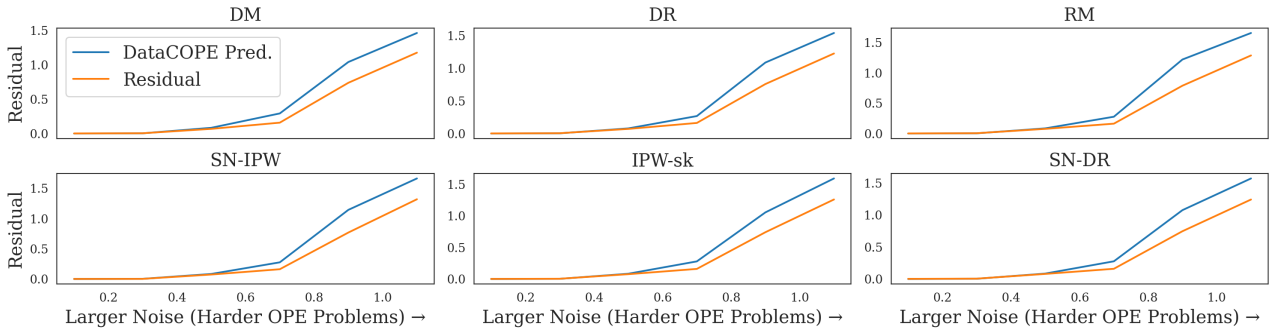


Figure 7. DataCOPE works as an accurate proxy of the OPE residual. DataCOPE is able to predict OPE residuals of different algorithms with calibration. Dataset: Diabetes.

**Experiment Settings** For each dataset, we run different OPE algorithms and are able to calculate the value estimation residual (error), comparing this to the output of DataCOPE, which has decomposed the uncertainties into aleatoric and epistemic components. The correlation between the OPE residuals and the uncertainties is then calculated and reported, noting that a strong correlation implies that the DataCOPE uncertainties have predictive power when it comes to determining if the OPE is accurate.

To test the impact of mismatched behavior and target policies, in the `Diabetes` dataset, we consider the impact of injecting different levels of noise during policy generation. With higher training noise, the policies will be forced to rely on more general features and employ simpler heuristics - in this way the relative *complexity* of the target policy compared to the behavior policy *increases* allowing the difficulty of OPE given the mismatched behavior policy and target policy to be manipulated and explored.

Specifically, we test across a range of values of complexity and add Gaussian noise with variance  $[0.1, 0.3, 0.5, 0.7, 0.9, 1.1]$  during behavior policy learning. We then perform calibration to predict the OPE residuals using DataCOPE according to Equation (9): we fit model  $h$  on the held-out training data, and dub such a model the calibrated OPE residual predictor, before again aiming to predict OPE residuals using our uncertainty decomposition.

**Results** Raw numerical results are reported in Table 1, which quantitatively shows the correlation between uncertainties obtained from DataCOPE and value estimation residuals of different OPE algorithms. DataCOPE in general performs well across all datasets and algorithms, with a strong correlation with the residuals.

**Take-away:** *Calibrated uncertainties decomposed by DataCOPE can be used to accurately predict the performance across a range of OPE algorithms.*

On the impact of complexity, Figure 6 and Figure 7 visualizes our experiment results, plotting both the prediction of DataCOPE, and the true OPE residual for multiple OPE methods as the complexity differential increases. DataCOPE is able to predict OPE residuals of various algorithms with high accuracy. In general, the deviations of predicted values from real residuals get higher as the complexity increases. Moreover, we empirically find DataCOPE tends to overestimate the OPE residual, lending its use as a performance lower-bound, being most suitable for high-stake application scenarios like healthcare.

**Take-away:** *As complexity differential between policies increases, the OPE performance decreases - an effect DataCOPE is able to predict and track well.*

## C. Implementation Details

### C.1. Code

Our code is open-sourced anonymously at

<https://anonymous.4open.science/r/Data-Centric-OPE-C974>.

### C.2. Hardware and Training Time

We experiment on a machine with 2 TITAN X GPUs and 32 Intel(R) E5-2640 CPUs. In general, the computational expense of DataCOPE is cheap. With our PyTorch-based implementation, decomposing the uncertainty components with DataCOPE using 100 neural network ensembles takes approximately 10 minutes to run. OPE algorithms take 2 to 10 minutes to run depending on their complexity and can be accelerated by using GPUs.

### C.3. OPE Algorithms

Our implementation of OPE algorithms is based on the implementation of (Saito et al., 2020), open-sourced at <https://github.com/st-tech/zr-obp/tree/master/obp/oep>.

Different from (Saito et al., 2020), we do not assume the access of the behavior policy  $\pi_b$ , hence we build estimators of  $\pi_b$  based on the dataset by training classification or regression models with neural networks if a parameterized behavior policy is needed. e.g., in the case of IPW.

### C.4. Hyper-Parameters for Neural Approximators

The hyper-parameters for neural-network-based approximators we used for MDN are reported in Table 4.

Table 4. Hyper-parameters of neural network approximators

Hyper-param	Choice
Network Layer	3
Activation Function	ReLU
Hidden Units	64
Training Epoch	100
Optimizer	Adam
Learning Rate	1e-3

### C.5. Hyper-Parameters for DataCOPE

In DataCOPE, there are two hyper-parameters to be specified: the number of ensemble models used and the number of Gaussians in the MDN model in regression tasks. DataCOPE is robust to both hyper-parameters in our empirical study.

In our experiments, we find using 100 ensemble models is computationally affordable on tabular data. As it takes less than

10 minutes to run. In the meantime, we find using 10 ensemble models can provide decent performance.

The main objective of introducing MDN is to be able to quantitatively decompose the uncertainty in regression settings. Therefore, while the number of Gaussian mixtures used should depend on the inherent structure of data distribution, whether or not capturing the exact distribution is of less importance than being able to capture the uncertainty components. The number of Gaussians we use in our experiment is set to be 3 for all datasets. We empirically find using 3, 5, 10 mixture of Gaussians does not clearly affect the performance.

## D. Experiment Details

### D.1. Calibration

In practice the OPE residual  $\bar{\xi}$  is always infeasible, such a difficulty leads us to introduce a practical substitute that estimates  $\bar{\xi}$  without relying on that infeasible ground-truth policy value. We elaborate on the calibration step in this section.

We leveraged a cross-validation type method in the calibration step. To conform with typical procedures in supervised learning, we hold out a portion of the training data, where we have knowledge of their accurate instance-wise value (i.e., the return of the context-action pairs). By doing this, we are able to break down the uncertainties of the instance-wise value prediction. Subsequently, we use the reserved ground-truth value to calculate the OPE residuals of each OPE algorithm as the objective, and uncertainties as the input of a regression model. Our implementation employs the bootstrap sampling method.

It is also worth mentioning that such a calibration step is optional, and DataCOPE can evaluate OPE problems without calibration:

In fact, our contribution and novelty of DataCOPE do not rely on such a calibration step. This is demonstrated in our experiments: The difficulty of OPE problems is correlated with our decomposed uncertainties. Such a decomposition, without calibration, is useful in that (1) it permits a comparison among datasets to answer the question of “which dataset is most appropriate in evaluating a certain given policy” (Section 4.2) (2) it permits a hindsight interpretation of clinical guideline deployment and vulnerable group identification. (Section 4.3)

In the revision, we have updated our manuscript to highlight the application of DataCOPE in practice without calibration.

### D.2. Synthetic Dataset

**Policy Generation** Following standard procedures in the OPE literature (Chu et al., 2011; Li et al., 2012; Agrawal & Goyal, 2013), we adapt supervised learning datasets into example logged bandit datasets where we know the true underlying generative process.

We use linear models for the behavior policy  $\pi_b$ , which is trained on the dataset examples  $\{(x_i, y_i)\}_{i=1}^N$ . Given context  $x_i$  with corresponding labels  $y_i$  we learn a policy to generate actions by minimizing the expected negative log-likelihood

$$\mathbb{E}_{x,y}[\text{NLL}(\pi_b(x_i), y_i)], \quad (11)$$

under a predictive Gaussian/Bernoulli distribution for regression/classification respectively.

We then evaluate  $(x_i, a_i)$  with the help of ground-truth labels  $y_i$ . For classification tasks, the reward of action  $a$  is given as  $r_i^a = \mathbf{1}(a_i = y_i)$ , where  $\mathbf{1}$  is the indicator function; while for regression tasks, the reward of action  $a_i$  is given by the coefficient of determination of the prediction. i.e.,  $r_i^a = 1 - \frac{u}{v}$ , where  $u$  is the residual sum of squares  $u = \|y - a\|_2^2$  and  $v$  is the total sum of squares  $v = \|y - \bar{y}\|_2^2$ ,  $\bar{y}$  denotes the mean of  $y$ . In this way, we can generate  $\{(x_i, a_i, r_i^a)\}_{i=1}^N$  tuples as our dataset containing  $N$  examples.

**Target Policy** We also generate target policies using neural network models trained on the *same* training data  $\{(x_i, y_i)\}_{i=1}^N$ , but with injected noise (that will depend on the particular experiment, as  $\pi_e$ ). The ground-truth performance of policy  $\pi_e$  is then given by  $r_j^b = \mathbf{1}(b_j = y_j)$ , where  $b_j = \pi_e(x_j)$  - importantly this is available to us and thus allows for the evaluation of the quality of the OPE.

In particular, the off-policy evaluation problem is by definition estimating  $\mathbb{E}[r_j^b]$ , for  $j = 1, \dots, N$ . As we have access to  $y_j, j = 1, \dots, N$ , we can quantitatively evaluate our method and compare results with the ground-truth policy values.



### D.3. Organ Transplant

**Data Description and Logged Contextual Bandit Formalism.** We examine data from the Organ Procurement & Transplant Network (OPTN) (Leppke et al., 2013), which includes information on patients registered for liver transplants from 1995 to 2020. Our focus is on the policy of matching organs that become available to patients who are waiting for a transplant.

For each decision, the set of potential patients (action space) includes those on the waitlist at the time the organ becomes available, and the information considered for each patient (context) includes both the organ and the patient’s characteristics.

**Data Preparation.** The OPTN dataset includes 308,912 patients who were either waiting for or had received a liver transplant. We eliminated patients who hadn’t received a transplant, were under 18 or had a donor under 18, or had missing data for certain variables, leaving us with 31,045 patients. Consequently, we consider 8 features in total: *ABO Mismatch*, *Age*, *Creatinine*, *Dialysis*, *INR*, *Life Support*, *Bilirubin*, and *Weight Difference*.

The focus of our experiments is on patients who received organs prior to 2005, allowing us to divide the data into distinct phases based on evolving allocation guidelines. Specifically, 512 patients were allocated organs under the Institute of Medicine recommendations before 2000, 1969 patients under the OPTN Final Rule between 2000 and 2002, and 2515 patients after the implementation of MELD between February 26th, 2002 and May 1st, 2003. Those data serve as the training set, with additional data containing 3914 patients collected between May 1st, 2003 and January 1st, 2005 used to evaluate OPE for MELD and DataCOPE’s estimation of MELD.

We streamline the data following (Hüyük et al., 2022) and select patients out of their contemporary patients from the transplant waitlist for an organ to be allocated to. Hence formalizing the dataset as a logged contextual bandit task. For every coming organ, the allocation policy selects one from the waiting patients to allocate the organ and receives the patients’ survival time as the reward.

### E. Additional Experiments

In order to validate that DataCOPE is effective and powerful, we further validate DataCOPE on other UCI datasets. We experiment with synthetic contextual bandit dataset generated from both classification tasks *Digits*, *Wine* and regression task *Boston*. Results are presented in Table 5. On all datasets, we find DataCOPE is able to act as a well-performing proxy of the OPE residual.

Table 5. DataCOPE is able to predict the difficulty of various OPE algorithms under various settings. The correlation between DataCOPE’s two uncertainties and OPE residuals under different settings is high. The ablation studies show that uncertainty decomposition is important in predicting OPE performance. (Reported numbers: correlation for DataCOPE rows, and for the ablation studies we report the performance difference compared with DataCOPE, **higher is better**)

Dataset	Ablations	DM	DR	RM	SNIPW	IPWsk	SNDR
Wine	DataCOPE	0.801	0.873	0.915	0.918	0.829	0.895
	w/o $v_{al}$	-0.004	-0.018	-0.238	-0.193	-0.176	-0.045
	w/o $v_{ep}$	-0.049	-0.027	-0.037	-0.02	-0.019	-0.008
	w/o Decomposition	-0.015	-0.003	-0.09	-0.062	-0.057	-0.0
Digits	DataCOPE	0.912	0.907	0.963	0.996	0.899	0.957
	w/o $v_{al}$	-0.224	-0.22	-0.004	-0.028	-0.099	-0.126
	w/o $v_{ep}$	-0.3	-0.295	-0.0	-0.063	-0.154	-0.188
	w/o Decomposition	-0.241	-0.238	-0.002	-0.035	-0.111	-0.14
Boston	DataCOPE	0.805	0.787	0.786	0.779	0.776	0.787
	w/o $v_{al}$	-0.024	-0.062	-0.06	-0.057	-0.045	-0.068
	w/o $v_{ep}$	-0.071	-0.105	-0.117	-0.107	-0.094	-0.107
	w/o Decomposition	-0.07	-0.104	-0.117	-0.106	-0.094	-0.106