

COMPARING TARGETING STRATEGIES FOR MAXIMIZING SOCIAL WELFARE WITH LIMITED RESOURCES

Anonymous authors

Paper under double-blind review

ABSTRACT

Machine learning is increasingly used to select which individuals receive limited-resource interventions in domains such as human services, education, development, and more. However, it is often not apparent what the right quantity is for models to predict. In particular, policymakers rarely have access to data from a randomized controlled trial (RCT) that would enable accurate estimates of treatment effects – which individuals would benefit more from the intervention. Observational data is more likely to be available, creating a substantial risk of bias in treatment effect estimates. Practitioners instead commonly use a technique termed “risk-based targeting” where the model is just used to predict each individual’s status quo outcome (an easier, non-causal task). Those with higher predicted risk are offered treatment. There is currently almost no empirical evidence to inform which choices lead to the most effect machine learning-informed targeting strategies in social domains. In this work, we use data from 5 real-world RCTs in a variety of domains to empirically assess such choices. We find that risk-based targeting is **typically** inferior to targeting based on even biased estimates of treatment effects. Moreover, these results hold even when the policymaker has strong normative preferences for assisting higher-risk individuals. Our results imply that practitioners may benefit from incorporating even weak evidence about heterogeneous causal effects to inform targeting **in a wider array of settings than current practice**.

1 INTRODUCTION

Policymakers often face the difficulty of allocating a resource-limited intervention with the goal of targeting the intervention towards those who will benefit most from it. Indeed, the causal inference literature documents that any given treatment may not have the same effect on every individual that receives it (Wager & Athey, 2018; Künzel et al., 2019; Varadhan & Seeger, 2013). When there are observable features that correlated with greater benefit from the treatment, such variation can be used for targeting. Heterogeneity of treatment effect (HTE) refers to this nonrandom, explainable variability in the direction and magnitude of treatment effects for individuals within a population. Given this variability, policymakers often face the problem of selecting who to treat when having to assign a particular treatment to a group of people under a fixed budget. Machine learning methods seem to offer the promise of discovering richer forms of heterogeneity, allowing more effective targeting of interventions in practice.

The main challenge is that heterogeneous treatment effects are difficult to learn: doing so requires a potentially large amount of data that is *unconfounded*, i.e., where treatment is assigned in a manner (conditionally) independent of each individual’s outcomes. In an idealized setting, one could conduct a randomized controlled trial (RCT) and experimentally find the subpopulations that benefit most from a treatment. However, this is not always possible for two reasons: 1) Conducting an RCT takes time (potentially years) and resources that the policymaker may not be willing to spend. 2) In some domains, there are ethical objections to experimentation. For example, while the RCT is being conducted, people in genuine need of the treatment could be assigned to the control group and suffer negative outcomes. Policymakers may prefer to prioritize access to treatment

054 via an already-available proxy metric that is believed to align with need if experimenting to gather
055 additional evidence would be seen as unethical.

056 One particularly common proxy is to target according to the *baseline* risk each individual faces, i.e.,
057 their expected outcome in the absence of treatment (as opposed to the treatment effect, which is
058 the difference in outcomes induced by treatment). This strategy has been referred to as *risk-based*
059 *targeting* (Wilder & Welle, 2024). Individuals with poor predicted baseline outcomes may be seen
060 as needing assistance the most. Policymakers may also believe that these individuals also stand the
061 most to benefit since their status-quo prognosis is the worst. Importantly, baseline risks can often
062 be learned using existing administrative data (from before a treatment was introduced) instead of re-
063 quiring a new experiment, making this strategy easily implementable in many practical settings. For
064 all of these reasons, risk-based targeting has seen widespread use by policymakers in a wide range of
065 domains, including targeting humanitarian assistance (Aiken et al., 2021), allocating homelessness
066 services via vulnerability scores (Shinn & Richard, 2022), and the use of “early warning systems”
067 in education (Perdomo et al., 2023).

068 Despite this widespread usage, there is only a limited amount of work which empirically assesses
069 the effectiveness of risk-based targeting: do individuals with the greatest baseline risk actually tend
070 to benefit the most from intervention? The few existing studies speak only to a single, specific
071 application domain each. (Athey et al., 2024) study an RCT where students in a university program
072 were provided a nudge (treatment) as a reminder to renew their financial aid application, concluding
073 that students with intermediate non-renewal risk saw the largest treatment effect. Students with
074 greatest risk of non-renewal, who would be targeted under a risk-based strategy, saw less benefit.
075 (Ascarza, 2018) study a marketing domain and use two field experiments to show that targeting
076 high-risk customers, or customers likely to churn as predicted by a machine learning model, can
077 be ineffective, encouraging practitioners to use RCTs to better inform their decision. However, as
078 discussed above, running a RCT may be infeasible in many settings. The alternative more likely to
079 be available to practitioners is to simply estimate treatment effects using observational data which
080 likely suffers from confounding, potentially leading to biased estimates of treatment effects.

081 How should practitioners navigate this tradeoff between a more easily-learnable label that may not
082 always correlate with benefit from an intervention (baseline risk) and a difficult-to-learn quantity
083 (heterogeneous treatment effects) that captures the impact of the intervention? This corresponds to
084 the choice of the appropriate target for prediction, as opposed to the specific model used to make the
085 prediction. The choice of outcome variable has been observed to exert a disproportionate influence
086 on the impacts of machine learning systems in many settings (Obermeyer et al., 2019; Coster, 2013;
087 Gerdon et al., 2022). In the setting of targeting interventions for causal impact, practitioners have
088 little empirically-grounded guidance. Our goal in this work is to inform the selection of an objective
089 function for machine-learning based targeting of scarce interventions. We make three contributions
090 towards this goal:

091 *First*, we assess the efficacy of risk-based targeting on a wider variety of real RCT datasets encom-
092 passing settings in economics, healthcare and education in contrast to prior studies that generally
093 focus on one dataset. **We find a generally noisy and variable relationship between baseline risk
094 and treatment effects: high-risk individuals seem to benefit more on average in most domains, but
095 with substantial variance in treatment effects which is not explained by baseline risk. Targeting in-
096 stead based on estimated treatment effects produces better results if the policymaker adopts a typical
utilitarian goal of maximizing the expected improvement in outcomes from the intervention.**

097 *Second*, we compare risk-based targeting to targeting policies based on biased estimates of treat-
098 ment effect obtained from confounded data. Such biased estimates are likely when conducting a
099 full-fledged RCT is infeasible and policymakers have to rely on available observational data alone.
100 Accordingly, potentially-biased causal estimates represent the likely alternative to risk-based target-
101 ing in many domains of practical interest. To our knowledge, these two strategies have not been
102 explicitly compared. We find that across even relatively severe levels of confounding, a utilitarian
103 policymaker often prefers targeting according to biased estimates of the treatment effect rather than
104 baseline risk.

105 *Finally*, we analyze the setting where a policymaker has inequality-averse preferences: oftentimes,
106 policymakers may prefer interventions which benefit those who are worse-off to begin with even
107 if they produce less aggregate impact. Such normative goals are one possible justification for risk-

108 based targeting, even if risk-based targeting is less attractive in standard utilitarian terms. We com-
 109 pare the two targeting strategies under popular classes of social welfare functions which capture
 110 inequality-averse preferences. We find that treatment effect based targeting is typically favorable
 111 to risk-based targeting, even in scenarios where policy makers are inequality-averse and the only
 112 available data is confounded to some degree.

114 2 RELATED WORK

116 Measuring heterogeneity in treatment effect and choosing which subpopulations to assign a treat-
 117 ment to has long been an active avenue of research in causal inference literature with a variety of
 118 methods proposed to solve this problem. (Green & Kern, 2012; Hill, 2011; Hill & Su, 2013; Foster
 119 et al., 2011; Wager & Athey, 2018) use forest-based algorithms to identify groups that show hetero-
 120 geneity in treatment effect with other identified groups. (Tian et al., 2014) proposed to measure the
 121 interaction between treatment and covariates by numerically binarizing the treatment and including
 122 the products of this variable with each covariate in a regression model. (Künzel et al., 2019) uses
 123 meta-learners that decompose estimating the CATE into several sub-regression problems that can be
 124 solved with any regression or supervised learning method. The problem of choosing who to treat is
 125 closely related to identifying the heterogeneity in treatment effects. This often involves balancing
 126 policies based solely on estimates of conditional average treatment effect (CATE) with additional
 127 prioritization rules set by the policymaker. (Yadlowsky et al., 2021) proposes rank-weighted average
 128 treatment effect metrics for testing the quality of treatment prioritization rules, providing an example
 129 involving optimal targeting of aspirin to stroke patients.

130 3 PRELIMINARIES

132 Consider a setting where there is a population of individuals who are candidates for a treatment or
 133 intervention. Each individual has a feature vector $X \in \mathbb{R}^d$. Here we are concerned with binary treat-
 134 ments. Following Neyman ((Splawa-Neyman et al., 1990)) and Ruben’s ((Rubin, 1974)) potential
 135 outcomes framework, we use $Y^{(1)}$ to denote the outcome that an individual would experience under
 136 treatment and $Y^{(0)}$ to denote the outcome they would experience if not treated. Their individual
 137 treatment effect, quantifying their benefit from receiving treatment, is $Y^{(1)} - Y^{(0)}$. We assume that
 138 $(X, Y^{(0)}, Y^{(1)})$ are drawn i.i.d. for each individual from some joint distribution. In order to identify
 139 individuals who are likely to benefit, a common strategy is to use individuals’ observed covariates
 140 to predict the expected treatment effect. The conditional average treatment effect (CATE) at $X = x$
 141 is defined as:

$$142 \tau(x) = \mathbb{E} \left[Y^{(1)} - Y^{(0)} \middle| X = x \right]. \quad (1)$$

144 Estimating the CATE in order to target based on treatment effects is a difficult statistical problem.
 145 Suppose we have access to data corresponding to n people, labeled $i = 1, \dots, n$, consisting of fea-
 146 tures X_i , a treatment assignment $W_i \in \{0, 1\}$, and the observed outcome $Y_i = Y_i^{(W_i)}$. Importantly,
 147 for each individual, we can observe only the outcome corresponding to the treatment they were
 148 actually assigned. Accordingly, identifying treatment effects typically requires a a no-unobserved-
 149 confounders assumption (Rosenbaum & Rubin, 1983):

$$150 \{Y^{(0)}, Y^{(1)}\} \perp\!\!\!\perp W \mid X. \quad (2)$$

152 This assumption is most credible in the context of a randomized controlled trial (RCT). In an RCT,
 153 the assignment of treatment, represented by W_i , is assigned independently of the potential outcomes
 154 Y_i (potentially after stratification on covariates X_i). When data is purely observational, practitioners
 155 typically try to select a sufficiently rich set of covariates X such that all potential confounders
 156 between outcomes and treatment assignment are measured. However, ensuring that all confounders
 157 are completely captured is notoriously difficult in practice, creating the likelihood that some bias in
 158 the estimated CATE remains (LaLonde, 1986; Pearl, 2009; Skelly et al., 2012; Milli et al., 2022).

159 As an alternative to targeting on treatment effects, policymakers often decide to treat people who are
 160 more vulnerable or worse-off at present, without attempting to quantify the benefit these individuals
 161 receive from treatment. This is quantified via a ‘baseline risk’; we refer to the resulting allocation
 strategy as ‘risk-based targeting’ (Wilder & Welle, 2024). Baseline risk may sometimes be directly

measured quantity (one of the covariates in X , for example baseline test scores in an educational context). In many settings though, it is estimated using a predictive model that uses the covariates as input. Let b be a function that maps a set of covariates to a baseline risk measurement such that $b(X)$ represents the baseline risk and $b(X_i)$ denotes the baseline risk associated with i . Then this method involves selecting individuals with the highest values of b for treatment, implying that these individuals have the highest 'risk' associated with them, which needs immediate resolution. It is important to note that this strategy requires only data on baseline outcomes prior to program implementation, with no information about the treatment's effect being incorporated in the policymaker's decision.

The example from (Athey et al., 2024) makes the distinction between these two methods clearer. In an experiment where the objective was to "nudge" or remind students in a college program to renew their financial aid applications, targeting based on baseline risk assumes that students predicted to be least likely to renew their aid applications (as determined by a machine learning model) should be prioritized. Meanwhile, the results of the RCT conducted during this experiment show that high values of treatment effect correspond to students with intermediate likelihoods of renewing their aid applications prior to treatment, demonstrating a disparity in the two methods of targeting.

4 METHODS

4.1 OVERVIEW

Our goal is to compare risk-based targeting to treatment effect-based targeting on possibly confounded datasets, with varying degrees of confounding and under different social welfare functions for the policymakers making the treatment assignment policy. To enable this comparison, we use a range of datasets from real-world randomized controlled trials (RCTs). Using RCTs enables us to credibly estimate heterogeneous treatment effects since the no-unobserved-confounders assumption is guaranteed to be satisfied. Then, we simulate each of the targeting policies of interest and compare their effectiveness with respect to the randomization-enabled treatment effect estimates, under varying utility functions for the policymaker. We now detail the methodology used in each step of this process, starting with the datasets used.

4.2 DATASETS

We conduct experiments on a variety of RCTs across different domains as detailed below:

- Targeting the Ultra Poor (TUP) in India ((Banerjee et al., 2021)): This RCT was conducted to study the long-term effects of providing large one-time capital grants to low income-families and observing how family income and overall consumption changed over a period of 7 years. We consider a family's total expenditure as the outcome, which is positively affected by treatment.
- NSW (National Supported Work demonstration) Dataset ((Dehejia & Wahba, 1999; 2002; LaLonde, 1986): This study estimated the impact of the National Supported Work Demonstration, a job training program, on beneficiaries' income in 1978. We consider an individual's income in 1978 as the outcome, which is positively affected by treatment.
- Postoperative Pain Dataset: Patients undergoing operations like tracheal intubations often experience throat pain following treatment (Mchardy & Chung, 1999). This RCT was conducted to test the efficacy of a licorice solution at reducing postoperative sore throat. The outcome we focus on is a patient's throat pain 4 hours after surgery. Here, the effect of the treatment is to reduce the amount of throat pain, hence the treatment effect is negative. In order to maintain consistency with other plots, we present results with the sign for treatment effect reversed.
- Acupuncture Dataset: This RCT aimed to determine the effect of acupuncture therapy on headache severity in patients with chronic headaches. Our outcome variable is headache severity 1 year post-randomization. Here again, the effect of the treatment is to reduce the severity of headaches, hence the treatment effect is negative. In order to maintain consistency with other plots, we present results with the sign for treatment effect reversed.
- Tennessee's Student Teacher Achievement Ratio (STAR) project (Achilles et al., 2008): The Tennessee State Department of Education conducted a comprehensive four-year study

called the Student/Teacher Achievement Ratio (STAR) to examine the effects of class size on student performance. The study design included three different classroom configurations: 1) Small classes with 13-17 students per teacher, 2) Regular classes with 22-25 students per teacher, 3) Regular classes with 22-25 students plus a full-time teacher’s aide. In this paper, we only focus on the first two types of classes mentioned above, so as to maintain consistency with treatment value being binary in other RCTs. We focus on students in kindergarten and a cumulative measure of their scores on various tests as the outcome under consideration.

For each of these datasets, we estimate $E[Y^{(0)}|X]$ using a machine learning model applied to the RCT’s control group and set $b(X) = E[Y^{(0)}|X]$ or $b(X) = -E[Y^{(0)}|X]$ as appropriate (i.e., depending on whether larger outcome values are better or worse). For instance, children with lower baseline test scores in the STAR dataset and patients with high baseline headache severity in the Acupuncture Dataset are considered to be high risk. Additional details about preprocessing for all datasets are included in A.

4.3 ESTIMATING HETEROGENEOUS TREATMENT EFFECTS

We estimate heterogeneous treatment effects on each dataset using a doubly-robust estimator (Kennedy, 2023a). The DR estimator splits the data into separate folds. For each fold, we estimate models for both the expected outcome and the treatment variable (estimating the latter even when the propensity scores are known can increase statistical efficiency Hirano et al. (2000)). Let $\hat{\mu}(X, A)$ be the estimated mean outcome for an individual with covariates X and treatment assignment A , and $\hat{\pi}(X)$ be the estimated propensity score. For each individual in the held-out data for the fold, we estimate their *pseudo-outcomes*, defined as

$$\chi_i(A) = \hat{\mu}(X_i, A) + \frac{1[W_i = A](Y_i - \hat{\mu}(X_i, A))}{A\hat{\pi}(X_i) + (1 - A)(1 - \hat{\pi}(X_i))}.$$

If at least one of $\hat{\mu}$ or $\hat{\pi}$ is correctly specified, $\chi_i(A)$ has expectation (over the random treatment assignment) equal to $Y_i^{(A)}$, which allows us to use it as a proxy for the unobserved outcomes in evaluating counterfactual evaluation policies.

4.4 EXPLORING TREATMENT EFFECT HETEROGENEITY WITH RESPECT TO BASELINE RISK

Our first analysis tests one potential rationale for risk-based targeting strategies: the hypothesis that individuals with greater baseline risk will also tend to have greater treatment effects. We frame this as estimating $E[Y^{(1)} - Y^{(0)}|b(X)]$, a conditional average treatment effect just with respect to value of the risk score b .

We follow the doubly-robust approach to estimating CATEs, where the pseudo-outcome difference $\chi_i(1) - \chi_i(0)$ is regressed on the covariates of interest (Kennedy, 2023a). Because our covariate of interest, b , is one-dimensional, we use a kernel regression method to estimate the CATE as a generic smooth function. Specifically, we employ a Gaussian kernel smoothing method. We sort every individual/household by their baseline risk and center an adaptive Gaussian kernel about each data point. The bandwidth of each kernel is determined adaptively based on the local density of the data, defined as half the range of baseline risk values within a fixed window of 200 data points. This ensures that we are able to estimate greater variation in data-rich regions of the space, while imposing greater smoothness at the extremes where less data is present.

Given the kernel function $K(u) = \exp(-\frac{1}{2}u^2)$, the CATE estimate at $b(X_i)$ is given by:

$$\hat{\tau}(b(X_i)) = \frac{\sum_{j=1}^n K(\frac{b(X_j) - b(X_i)}{\sigma_i}) \hat{\tau}_j}{\sum_{j=1}^n K(\frac{b(X_j) - b(X_i)}{\sigma_i})} \quad (3)$$

where $\hat{\tau}_j$ is the estimated difference in pseudo outcomes, for unit j , $\chi_j(1) - \chi_j(0)$, as determined by the doubly robust estimator, and σ_i is the adaptive bandwidth calculated as:

$$\sigma_i = \frac{1}{2}(b(X_{i+100}) - b(X_{i-99})) \quad (4)$$

for a window of 200 data points centered at i . The confidence intervals are computed using a weighted variance estimate:

$$\text{CI} = \hat{\tau}(b(X_i)) \pm 1.96 \sqrt{\frac{\sum_{j=1}^n K\left(\frac{b(X_j) - b(X_i)}{\sigma_i}\right) (\hat{\tau}_j - \hat{\tau}(b(X_i)))^2}{\left(\sum_{j=1}^n K\left(\frac{b(X_j) - b(X_i)}{\sigma_i}\right)\right)^2}} \quad (5)$$

This approach allows us to capture the heterogeneity in treatment effects across different levels of baseline risk while accounting for the varying density of data points.

4.5 INTRODUCING CONFOUNDING

Our next analysis aims to simulate conditions where we do not have access to perfectly conducted randomized controlled trials for our problem, in order to compare risk-based targeting to a plausible alternative in real world settings: targeting according to observational, and potentially biased, estimates of the CATE.

We introduce varying levels of confounding to the RCTs that we study. We do this by simulating adverse selection into treatment, where units are more likely to be observed if the estimated individual-level treatment effects deviate from the mean. Specifically, we generate the biased ‘‘observational’’ dataset by removing data in a systematic manner. This process, inspired by (Kallus & Zhou, 2021), is controlled by a parameter k giving the fraction of data removed, with higher k corresponding to more biased estimates.

From the treated units, we remove the examples that lie in the top $k\%$ percent when ordered in descending order of $(\chi_i(1) - \chi_i(0))$ (assuming treatment effect is positive) while for the untreated units, we remove the examples that lie in the bottom $k\%$ of examples when ordered in descending order of $(\chi_i(1) - \chi_i(0))$. In simpler terms, for treated units, we remove examples for which the treatment ‘went well’ (most positive), while for untreated units, we remove examples for which the lack of treatment did not go well (least positive). **This can be seen as a strengthening of a typical mechanism for confounding: a typical concern is that individuals are selected for treatment based on unobservable characteristics that are correlated with their potential outcomes, where we simulate selection into treatment based on the actual potential outcomes.**

Our goal is to compare policies learned using this biased data to risk-based targeting. This evaluation is enabled by access to the original, randomized data. Consider a hypothetical policy that assigns treatments $A(X) \in \{0, 1\}$ as a function of individuals’ features X . The mean outcome under policy A , $\mathbb{E}[Y^{(A(X))}]$, can be decomposed as

$$\mathbb{E}[Y^{(A(X))}] = \mathbb{E}[Y^{(0)}] + \Pr(A(X) = 1) \mathbb{E}[Y^{(1)} - Y^{(0)} | A(X) = 1].$$

The term $\mathbb{E}[Y^{(1)} - Y^{(0)} | A(X) = 1]$ represents the treatment effect on the treated population and captures how much the policy improves over no treatment. For policies with the same budget (equal $\Pr(A(X) = 1)$), only this term varies and so we assess allocation policies by their expected treatment-on-the-treated. Following standard doubly-robust estimators for (group) average treatment effects Kennedy (2023b), we empirically estimate this quantity as

$$\frac{1}{\sum_{j=1}^n A(X_j)} \sum_{j=1}^n A(X_j) (\chi_j(1) - \chi_j(0)). \quad (6)$$

4.6 FAMILIES OF WELFARE FUNCTIONS

In order to simulate policymakers with varying preferences for who to treat, we compare risk-based targeting to treatment-effect based targeting on some general utility/social welfare functions that fall under the category of ‘weighted power mean functions’ as described in (Pardeshi et al., 2024). The weighted power mean $M(\mathbf{u}; \mathbf{w}, p)$ for $\mathbf{u} \in \mathbb{R}_+^d$, $\mathbf{w} \in \Delta_{d-1}$, and $p \in \mathbb{R} \cup \{\pm\infty\}$ is defined as:

$$M(\mathbf{u}; \mathbf{w}, p) = \begin{cases} \left(\sum_{i=1}^d w_i u_i^p\right)^{1/p} & p \neq 0 \\ \prod_{i=1}^d u_i^{w_i} & p = 0 \end{cases}$$

We choose this family of welfare functions it contains all function satisfying a standard set of axiomatic properties Pardeshi et al. (2024). The parameter p determines the specific type of welfare function:

- **Utilitarian welfare** ($p = 1$): $M(\mathbf{u}; \mathbf{w}, p = 1) = \sum_{i=1}^d w_i u_i$
- **Nash welfare** ($p = 0$): $M(\mathbf{u}; \mathbf{w}, p = 0) = \prod_{i=1}^d u_i^{w_i}$

We consider utilitarian welfare with two sets of weights. First, the uniform weights $w_i = 1 \forall i \in [n]$. Second, $w_i = \frac{n \cdot e^{\alpha b'(X_i)}}{\sum_{j=1}^n e^{\alpha b'(X_j)}}$ where α is a hyperparameter and $b'(X_i)$ is represents the percentile score of the baseline risk for the i^{th} example, with the example with highest baseline risk having score 1 and the example with lowest baseline risk having score 0. This assigns greater weight to individuals with high values of baseline risk for high values of α , thereby simulating a scenario where a policymaker might value treating these "high risk" individuals for other reasons. α can be interpreted as $2 \log \left(\frac{w_{75}}{w_{25}} \right)$ where w_{75} is the weight given to the 75th percentile example by baseline risk and w_{25} is the weight given to the 25th percentile example by baseline risk.

The Nash social welfare function has commonly been used to achieve a balance between maximizing total welfare(utilitarian) and ensuring equitable distribution(egalitarian) Caragiannis et al. (2019); CHARKHGARD et al. (2022). Egalitarian welfare can sometimes leads to inefficiencies while utilitarian welfare can lead to unjust outcomes: Nash welfare often strikes a useful compromise between these two ends. We consider unweighted Nash welfare ($w_i = 1 \forall i$). In order to avoid the complications of utility when using an unweighted Nash welfare function, we scale up the estimated utilities for each example to a minimum value of 1.

Note that the Nash welfare can be equivalently formulated in log space Caragiannis et al. (2019) as

$$\frac{1}{n} \sum_{i=1}^n \log u_i.$$

When each individual's utility under an allocation policy corresponds to their realized outcome $Y_i^{(A(X_i))}$ (e.g., their income after the interventional period), we compare policies exactly as outlined in Equation 6, but with $Y^{(A(X))}$ replaced by $\log Y^{(A(X))}$, estimated by replicating the same procedure after taking logs of all outcome variables.

5 RESULTS

Figure 1 shows how treatment effect varies as a function of baseline risk for each of the 5 datasets we study, with 95% confidence intervals shaded around the estimated treatment effects. These intervals are pointwise Wald-type confidence intervals (Kennedy et al., 2019) and provide a measure of uncertainty for our smoothed estimates. **The estimated relationship between baseline risk and treatment effect is variable across domains. In most domains, the point estimate shows a general upward trend, indicating that individuals at greater risk benefit more (on average) from treatment. However, in the NSW domain, the point estimate is essentially flat. In addition, the confidence intervals are wide for all domains and there is very little statistically significant evidence in favor of high-risk individuals benefiting more. Wide confidence intervals reflect that there is significant variance in the pseudo-outcomes estimated for different individuals at the same level of baseline risk. That is, there is a great deal of variance in our estimated treatment effects that is not explained by baseline risk. From these results, we form two hypotheses. First, that risk-based targeting should, in most domains, perform better than a random allocation, since the point estimates generally show larger average effects at higher baseline risk. Second, that there is room to improve on risk-based targeting via strategies that leverage some of the substantial variance in treatment effects that is unexplained by baseline risk. The next section provides more statistically precise tests of these hypotheses by comparing the welfare associated with each targeting policy (a single number, which can be quantified more precisely than the entire curve shown in Figure 1).**

Figure 2 shows the comparison between risk-based targeting and treatment effect based targeting for the 5 datasets we study, with varying degrees of confounding and under different social welfare

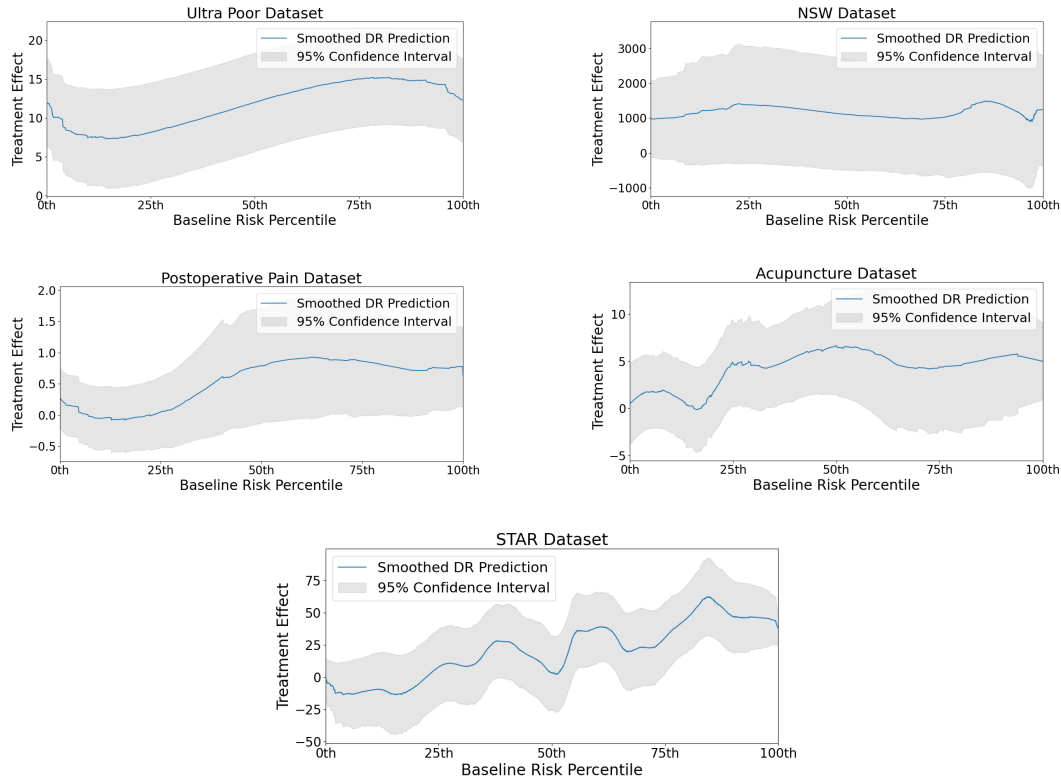


Figure 1: Observing Treatment Effect Heterogeneity across Different Settings by plotting Treatment Effect against Baseline Risk for each of our 5 datasets. We observe a unique trend for each dataset, indicating a lack of a consistent well-defined relation between the two quantities

functions used to represent policymaker preferences. For reference, we also show the performance of a random targeting strategy which allocates the same budget but targets uniformly random individuals. Both targeting methods seem to perform better than random targeting for almost all combinations of dataset/welfare functions we consider. The improved performance of risk-based targeting over random reflects the generally increasing relationship that Figure 1 shows between baseline risk and treatment effects. However, when treatment effects can be accurately estimated ($k = 0$, no confounding), targeting based on treatment effects always produces significantly higher utilitarian welfare (often producing a treatment-on-the-treated effect of three times or more greater than risk-based targeting). This indicates that when a policymaker seeks only to maximize aggregate benefit and can credibly estimate treatment effects, the gains from causal targeting are substantial.

As the level of confounding bias in treatment effect targeting increases (increased k), its effectiveness decreases. However, when the policymaker has utilitarian preferences, targeting based on biased treatment effect estimates still performs at least as well as risk-based targeting (and typically better) across all datasets, even for relatively severe levels of confounding. This indicates that using even relatively biased observational data to learn treatment rules is likely superior when the policymaker’s goal is just to maximize aggregate gain.

The second column of Figure 2 shows an alternative set of preferences, where the policymaker has a weighted utilitarian welfare function (for these plots, α is $2 \log 2$, the value at which the ratio of the 75th percentile weight to the 25th percentile weight is 2) which places greater weight on individuals with higher baseline risk, attaching a higher importance to the welfare of more vulnerable individuals. This welfare function decreases the gap between treatment effect and risk based targeting as individuals with higher risk now have more weight associated with them. However, targeting on (biased) treatment effects is still preferable to risk-based targeting across all datasets.

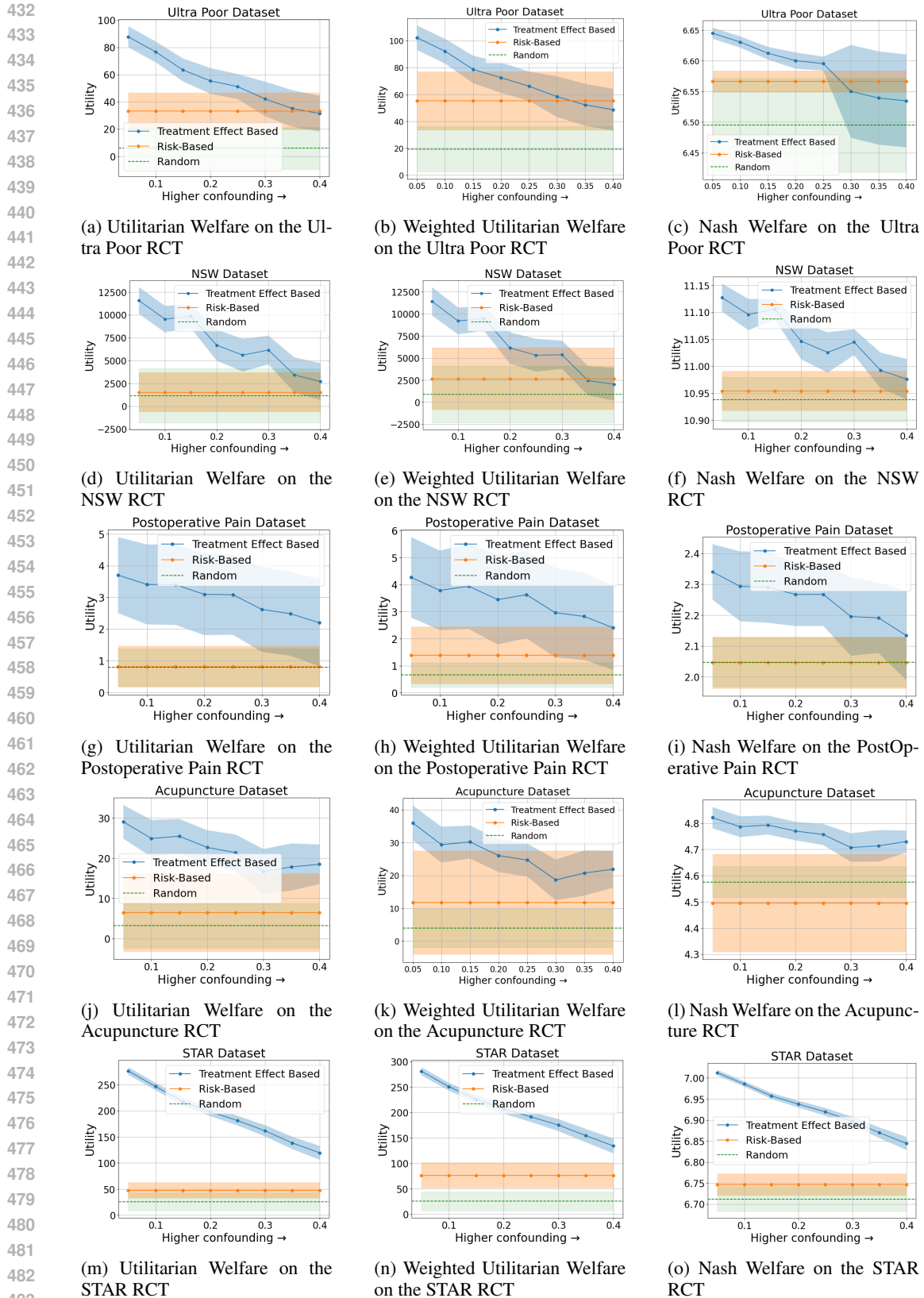


Figure 2: Comparison of risk-based targeting to biased treatment effect-based targeting by plotting the benefit offered by each policy against the amount of data systematically removed from the RCT to introduce confounding

This motivates us to ask: how much greater must the policymaker weight high-risk individuals in order to prefer risk-based targeting? In Table 1 we give the minimum value of α at which risk-based targeting finally outperformed treatment effect based-targeting at each level of systematic data removal. We limit α such that the ratio of the 75th percentile weight to the 25th percentile weight when sorted in descending order is less than 100; otherwise, only a few individuals have non-zero weights and estimating welfare gains becomes impossible. In general, we note that the required α values tend to be lower when we increase k , which follows directly from the fact that the treatment effect estimates become more biased at higher k . We note that most values in the table are at least $\alpha = 4$ at which the policymaker places over 7 times more weight on the welfare of an individual at the 75th percentile of baseline risk than an individual at the 25th percentile. High values of α act as an indicator that a policymaker prefers to target "higher risk" individuals more than others. Interestingly, for the STAR dataset, there is no α value (up to our upper bound) at which risk based targeting outperforms targeting based on treatment effects. We conclude that relatively extreme welfare weights are needed to rationalize risk-based targeting.

Table 1: Values of α for different k at which risk-based targeting outperforms treatment effect based targeting. 'na' indicates no such α was found

Dataset	5%	10%	15%	20%	25%	30%	35%	40%
Ultra Poor	na	8.5	5.5	4.5	3.5	2	1	0
NSW	na	6.5	8.25	4.5	4.25	3.5	1.5	1
Postoperative Pain	na	8	na	6.5	na	5.5	6.5	4
Acupuncture	na	7.5	na	6.75	6	3.75	4.5	5.75
STAR	na	na	na	na	na	na	na	na

Under the Nash social welfare function (third column of Figure 2), we observe some differences across our chosen settings but the general trend remains the same: as we increase confounding (increase the percentage of data we systematically remove), the benefit we accrue by following a treatment assignment policy based on biased treatment effect values decreases. At high levels of confounding, risk-based targeting accrues higher utility for a policymakers with a Nash social welfare function in the Ultra Poor setting in 2. However, across all other datasets, the policymaker prefers to target based on treatment effects even at high levels of confounding.

Collectively, these results indicate that policymakers are almost always better off targeting interventions based on treatment effect estimates rather than baseline risk values from a machine learning model, unless they are extremely inequality-averse.

6 CONCLUSION

This paper presents a systematic comparison between two of the most popular treatment assignment policies in use by policymakers today: risk-based targeting and treatment effect based targeting. We observe a tendency for risk-based targeting to produce higher welfare than a uniformly random allocation, confirming some of the intuition behind its widespread use by practitioners. However, targeting based on treatment effects results in substantially larger welfare gains. We then explore two potential considerations that may motivate risk-based targeting in practice: the threat of confounding in treatment effect estimates, and egalitarian preferences for assisting high-risk individuals. We find that either can narrow the gap between the two strategies, but relatively high levels of either factor (or typically both) are needed to fully rationalize risk-based targeting. Accordingly, the barriers to targeting based on treatment effects in practice may be overestimated at present. On some questions, our results are subject to greater uncertainty. For one, the RCTs we use are too small to quantify the entire curve giving treatment effects as a function of baseline risk with a high level of statistical precision. Future work could investigate settings with larger sample sizes, or strategies for pooling data across studies. Importantly though, we are able to give more precise conclusions (typically statistically significant) for our main results comparing the utility of risk based vs treatment effect targeting. Second, our investigation of egalitarian preferences assumes an essentially consequentialist perspective, where the policymaker's goal is to improve individuals' welfare as defined by their outcome. If policymakers have non-consequentialist preferences, for example viewing the assistance of those in need as an inherent good regardless of its effects, targeting directly on a measure of vulnerability may be more appropriate.

540 **Reproducibility Statement:** A rough version of the code is provided in the supplementary mate-
541 rial, which includes data preprocessing and experimentation for each of the datasets which we intend
542 to finalize and clean in the camera ready version. We also detail our procedures in the Appendix A
543 (dataset details) and in Section 4 (step by step experimental procedure).
544

545 REFERENCES

- 546
547 C.M. Achilles, Helen Pate Bain, Fred Bellott, Jayne Boyd-Zaharias, Jeremy Finn, John Folger, John
548 Johnston, and Elizabeth Word. Tennessee’s Student Teacher Achievement Ratio (STAR) project,
549 2008. URL <https://doi.org/10.7910/DVN/SIWH9F>.
- 550 Emily Aiken, Suzanne Bellue, Dean Karlan, Christopher R Udry, and Joshua Blumenstock. Machine
551 learning and mobile phone data can improve the targeting of humanitarian assistance. Working
552 Paper 29070, National Bureau of Economic Research, July 2021. URL <http://www.nber.org/papers/w29070>.
- 553
554 Eva Ascarza. Retention futility: Targeting high-risk customers might be ineffective. *Journal of*
555 *Marketing Research*, 55(1):80–98, 2018. doi: 10.1509/jmr.16.0163. URL <https://doi.org/10.1509/jmr.16.0163>.
- 556
557 Susan Athey and Stefan Wager. Policy learning with observational data, 2020. URL <https://arxiv.org/abs/1702.02896>.
- 558
559 Susan Athey, Niall Keleher, and Jann Spiess. Machine learning who to nudge: Causal vs predictive
560 targeting in a field experiment on student financial aid renewal, 2024. URL <https://arxiv.org/abs/2310.08672>.
- 561
562 Abhijit Banerjee, Esther Duflo, and Garima Sharma. Long-term effects of the targeting the ul-
563 tra poor program. *American Economic Review: Insights*, 3(4):471–86, December 2021. doi:
564 10.1257/aeri.20200667. URL <https://www.aeaweb.org/articles?id=10.1257/aeri.20200667>.
- 565
566 Ioannis Caragiannis, David Kurokawa, Hervé Moulin, Ariel D. Procaccia, Nisarg Shah, and Junxing
567 Wang. The unreasonable fairness of maximum nash welfare. *ACM Trans. Econ. Comput.*, 7(3),
568 September 2019. ISSN 2167-8375. doi: 10.1145/3355902. URL <https://doi.org/10.1145/3355902>.
- 569
570 HADI CHARKHGARD, KIMIA KESHANIAN, RASUL ESMAEILBEIGI, and PARISA
571 CHARKHGARD. The magic of nash social welfare in optimization: Do not sum, just multi-
572 ply! *The ANZIAM Journal*, 64(2):119–134, 2022. doi: 10.1017/S1446181122000074.
- 573
574 Wendy Coster. Making the best match: Selecting outcome measures for clinical trials and outcome
575 studies. *The American journal of occupational therapy : official publication of the American*
576 *Occupational Therapy Association*, 67:162–70, 03 2013. doi: 10.5014/ajot.2013.006015.
- 577
578 Rajeev Dehejia and Sadek Wahba. Propensity score-matching methods for nonexperimental causal
579 studies. *The Review of Economics and Statistics*, 84(1):151–161, 2002. URL <https://EconPapers.repec.org/RePEc:tpr:restat:v:84:y:2002:i:1:p:151-161>.
- 580
581 Rajeev H. Dehejia and Sadek Wahba. Causal effects in nonexperimental studies: Reevaluating the
582 evaluation of training programs. *Journal of the American Statistical Association*, 94(448):1053–
583 1062, 1999. ISSN 01621459. URL <http://www.jstor.org/stable/2669919>.
- 584
585 Jared C. Foster, J. M. Taylor, and Stephen J. Ruberg. Subgroup identification from randomized clin-
586 ical trial data. *Statistics in Medicine*, 30, 2011. URL <https://api.semanticscholar.org/CorpusID:24046082>.
- 587
588 Frederic Gerdon, Ruben L Bach, Christoph Kern, and Frauke Kreuter. Social impacts of algo-
589 rithmic decision-making: A research agenda for the social sciences. *Big Data & Society*, 9(1):
590 20539517221089305, 2022. doi: 10.1177/20539517221089305. URL <https://doi.org/10.1177/20539517221089305>.

- 594 Donald P. Green and Holger L. Kern. Modeling Heterogeneous Treatment Effects in Survey Ex-
595 periments with Bayesian Additive Regression Trees. *Public Opinion Quarterly*, 76(3):491–511,
596 09 2012. ISSN 0033-362X. doi: 10.1093/poq/nfs036. URL [https://doi.org/10.1093/](https://doi.org/10.1093/poq/nfs036)
597 [poq/nfs036](https://doi.org/10.1093/poq/nfs036).
- 598 Jennifer Hill and Yu-Sung Su. Assessing lack of common support in causal inference using Bayesian
599 nonparametrics: Implications for evaluating the effect of breastfeeding on children’s cognitive
600 outcomes. *The Annals of Applied Statistics*, 7(3):1386 – 1420, 2013. doi: 10.1214/13-AOAS630.
601 URL <https://doi.org/10.1214/13-AOAS630>.
- 602 Jennifer L. Hill. Bayesian nonparametric modeling for causal inference. *Journal of Computational*
603 *and Graphical Statistics*, 20(1):217–240, 2011. doi: 10.1198/jcgs.2010.08162. URL <https://doi.org/10.1198/jcgs.2010.08162>.
- 604 Keisuke Hirano, Guido W Imbens, and Geert Ridder. Efficient estimation of average treatment
605 effects using the estimated propensity score. Working Paper 251, National Bureau of Economic
606 Research, March 2000. URL <http://www.nber.org/papers/t0251>.
- 607 Nathan Kallus and Angela Zhou. Minimax-Optimal Policy Learning Under Unobserved Confound-
608 ing. *Management Science*, 67(5):2870–2890, May 2021. doi: 10.1287/mnsc.2020.3699. URL
609 <https://ideas.repec.org/a/inm/ormnsc/v67y2021i5p2870-2890.html>.
- 610 Edward H. Kennedy. Towards optimal doubly robust estimation of heterogeneous causal effects,
611 2023a. URL <https://arxiv.org/abs/2004.14497>.
- 612 Edward H. Kennedy. Semiparametric doubly robust targeted double machine learning: a review,
613 2023b. URL <https://arxiv.org/abs/2203.06469>.
- 614 Edward H Kennedy, Shreya Kangovi, and Nandita Mitra. Estimating scaled treatment effects with
615 multiple outcomes. *Statistical Methods in Medical Research*, 28(4):1094–1104, 2019. doi:
616 10.1177/0962280217747130. URL <https://doi.org/10.1177/0962280217747130>.
617 PMID: 29254442.
- 618 Sören R Künzel, Jasjeet S Sekhon, Peter J Bickel, and Bin Yu. Metalearners for estimating heteroge-
619 neous treatment effects using machine learning. *Proceedings of the national academy of sciences*,
620 116(10):4156–4165, 2019.
- 621 Sören R. Künzel, Jasjeet S. Sekhon, Peter J. Bickel, and Bin Yu. Metalearners for estimating het-
622 erogeneous treatment effects using machine learning. *Proceedings of the National Academy of*
623 *Sciences*, 116(10):4156–4165, 2019. doi: 10.1073/pnas.1804597116. URL [https://www.](https://www.pnas.org/doi/abs/10.1073/pnas.1804597116)
624 [pnas.org/doi/abs/10.1073/pnas.1804597116](https://www.pnas.org/doi/abs/10.1073/pnas.1804597116).
- 625 Robert LaLonde. Evaluating the econometric evaluations of training programs with experimen-
626 tal data. *American Economic Review*, 76(4):604–20, 1986. URL [https://EconPapers.](https://EconPapers.repec.org/RePEc:aea:aecrev:v:76:y:1986:i:4:p:604-20)
627 [repec.org/RePEc:aea:aecrev:v:76:y:1986:i:4:p:604–20](https://EconPapers.repec.org/RePEc:aea:aecrev:v:76:y:1986:i:4:p:604-20).
- 628 Fiona Mchardy and Frances F Chung. Postoperative sore throat: cause, prevention and treat-
629 ment. *Anaesthesia*, 54, 1999. URL [https://api.semanticscholar.org/CorpusID:](https://api.semanticscholar.org/CorpusID:31521199)
630 [31521199](https://api.semanticscholar.org/CorpusID:31521199).
- 631 Smitha Milli, Luca Belli, and Moritz Hardt. Causal inference struggles with agency on on-
632 line platforms. In *Proceedings of the 2022 ACM Conference on Fairness, Accountability, and*
633 *Transparency*, FAccT ’22, pp. 357–365, New York, NY, USA, 2022. Association for Com-
634 puting Machinery. ISBN 9781450393522. doi: 10.1145/3531146.3533103. URL <https://doi.org/10.1145/3531146.3533103>.
- 641 Ziad Obermeyer, Brian Powers, Christine Vogeli, and Sendhil Mullainathan. Dissecting racial bias
642 in an algorithm used to manage the health of populations. *Science*, 366(6464):447–453, 2019.
643 doi: 10.1126/science.aax2342. URL [https://www.science.org/doi/abs/10.1126/](https://www.science.org/doi/abs/10.1126/science.aax2342)
644 [science.aax2342](https://www.science.org/doi/abs/10.1126/science.aax2342).
- 645 Kanad Shrikar Pardeshi, Itai Shapira, Ariel D. Procaccia, and Aarti Singh. Learning social welfare
646 functions, 2024. URL <https://arxiv.org/abs/2405.17700>.

- 648 Judea Pearl. *Causality*. Cambridge University Press, 2 edition, 2009.
- 649
- 650 Juan C. Perdomo, Tolani Britton, Moritz Hardt, and Rediet Abebe. Difficult lessons on social predic-
651 tion from wisconsin public schools, 2023. URL <https://arxiv.org/abs/2304.06205>.
- 652 Paul R. Rosenbaum and Donald B. Rubin. The central role of the propensity score in observational
653 studies for causal effects. *Biometrika*, 70(1):41–55, 1983. ISSN 00063444, 14643510. URL
654 <http://www.jstor.org/stable/2335942>.
- 655 Donald B. Rubin. Estimating causal effects of treatments in randomized and nonrandomized
656 studies. *Journal of Educational Psychology*, 66:688–701, 1974. URL [https://api.
657 semanticscholar.org/CorpusID:52832751](https://api.semanticscholar.org/CorpusID:52832751).
- 658
- 659 Marybeth Shinn and Molly K. Richard. Allocating homeless services after the withdrawal of the
660 vulnerability index–service prioritization decision assistance tool. *American Journal of Public
661 Health*, 112(3):378–382, 2022. doi: 10.2105/AJPH.2021.306628. URL [https://doi.org/
662 10.2105/AJPH.2021.306628](https://doi.org/10.2105/AJPH.2021.306628). PMID: 35196047.
- 663 Andrea Skelly, Joseph Dettori, and Erika Brodt. Assessing bias: the importance of considering con-
664 founding. *Evidence-based spine-care journal*, 3:9–12, 02 2012. doi: 10.1055/s-0031-1298595.
- 665
- 666 Jerzy Splawa-Neyman, D. M. Dabrowska, and T. P. Speed. On the application of probability theory
667 to agricultural experiments. essay on principles. section 9. *Statistical Science*, 5(4):465–472,
668 1990. ISSN 08834237, 21688745. URL <http://www.jstor.org/stable/2245382>.
- 669 Lu Tian, Ash A. Alizadeh, Andrew J. Gentles, and Robert Tibshirani. A simple method for es-
670 timating interactions between a treatment and a large number of covariates. *Journal of the
671 American Statistical Association*, 109(508):1517–1532, 2014. ISSN 01621459. URL [http:
672 //www.jstor.org/stable/24247388](http://www.jstor.org/stable/24247388).
- 673 Ravi Varadhan and John D. Seeger. Estimation and reporting of heterogeneity of treatment effects.
674 2013. URL <https://api.semanticscholar.org/CorpusID:9040373>.
- 675
- 676 Stefan Wager and Susan Athey. Estimation and inference of heterogeneous treatment effects using
677 random forests. *Journal of the American Statistical Association*, 113(523):1228–1242, 2018.
- 678 Bryan Wilder and Pim Welle. Learning treatment effects while treating those in need, 2024. URL
679 <https://arxiv.org/abs/2407.07596>.
- 680
- 681 Steve Yadlowsky, Scott D. Fleming, Nigam Haresh Shah, Emma Brunskill, and Stefan Wager. Eval-
682 uating treatment prioritization rules via rank-weighted average treatment effects. 2021. URL
683 <https://api.semanticscholar.org/CorpusID:244117089>.

685 A APPENDIX

687 A.1 DATASETS

- 689 • Targeting the Ultra Poor (TUP) in India ((Banerjee et al., 2021)): This RCT was conducted
690 to study the long-term effects of providing large one-time capital grants to low income-
691 families and observing how family income and overall consumption changed over a period
692 of 7 years. We consider a family’s total expenditure as the outcome, which is positively
693 affected by treatment. We filter the dataset before use by removing null values and per-
694 forming feature selection to limit the number of covariates. The dataset consists of 796 ex-
695 amples post filtering. We quantify baseline risk $b(X)$ as an estimate of baseline expenditure
696 $E[Y(0)|X]$ from a machine learning model, with low values of $E[Y(0)|X]$ corresponding
697 to high baseline risk and vice versa. This follows the hypothesis that households with low
698 expenditure at baseline will benefit most from the treatment.

698 While constructing a doubly robust estimator to estimate pseudo outcomes for this dataset,
699 we found the estimated propensity scores to be very high/low for certain examples, which
700 would consequently scale pseudo outcome estimates to unusually large values. Therefore,
701 we manually set propensity scores uniformly according to the treated:untreated ratio in the
RCT.

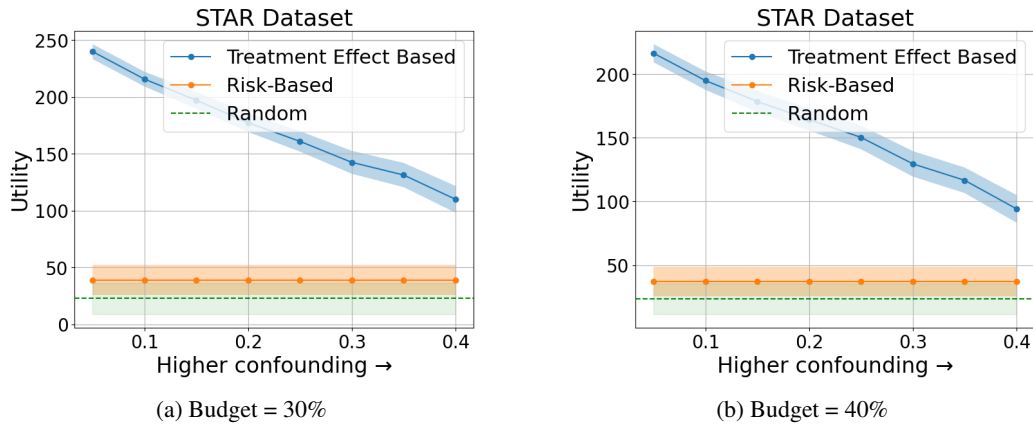
- 702 • NSW (National Supported Work demonstration) Dataset ((Dehejia & Wahba, 1999; 2002;
703 LaLonde, 1986): This is a popular causal inference dataset that was used to estimate the
704 impact of the National Supported Work Demonstration, a job training program, on ben-
705 eficiaries' income in 1978. The covariates include demographic variables like age, race,
706 marital status and academic background, along with the beneficiary's income in 1975 prior
707 to the experiment as a baseline. The dataset consists of 722 examples (297 treated and
708 425 control). Here too, we use an estimate of an individuals baseline income as a measure
709 of risk, following the hypothesis that low income individuals will benefit more from the
710 treatment.
- 711 • Tennessee's Student Teacher Achievement Ratio (STAR) project ((Achilles et al., 2008)):
712 The Tennessee State Department of Education conducted a comprehensive four-year study
713 called the Student/Teacher Achievement Ratio (STAR) to examine the effects of class size
714 on student performance. This research, backed by the Tennessee General Assembly, in-
715 volved 11601 students across 79 schools. The study design included three different class-
716 room configurations:
 - 717 – Small classes with 13-17 students per teacher
 - 718 – Regular classes with 22-25 students per teacher
 - 719 – Regular classes with 22-25 students plus a full-time teacher's aide

720 To ensure unbiased results, both students and teachers were randomly assigned to these dif-
721 ferent classroom types. The experiment began when the participants entered kindergarten
722 and continued through their third-grade year, allowing for a longitudinal analysis of the
723 impact of class size on educational outcomes. In this paper, we only focus on the first two
724 types of classes mentioned above, so as to stay consistent with treatment value being binary
725 in other RCTs. This large-scale research project aimed to provide empirical evidence on
726 the relationship between class size and student achievement. Again, we filter the dataset
727 before use by removing null values and performing feature selection to limit the number of
728 covariates. We focus on students in kindergarten and a cumulative measure of their scores
729 on various tests as the outcome under consideration. The filtered dataset consists of 3712
730 students examples. Since we do not have the students' test scores at baseline, we train a
731 random forest model on rows corresponding to students who did not receive the treatment
732 with their test scores at endline being the outcome variable. The prediction offered by this
733 model for every student is then used as a proxy for their baseline test scores and the nega-
734 tive of this value is used as baseline risk. This follows the general hypothesis that students
with low test scores need the treatment more.

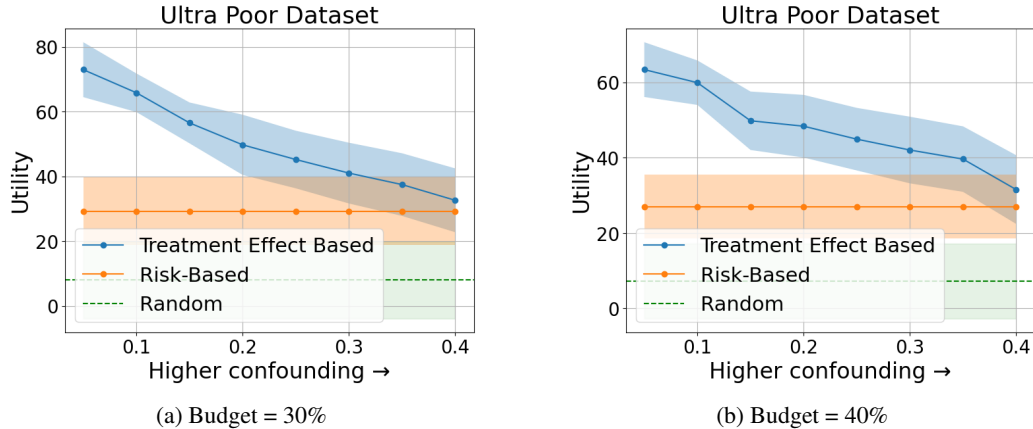
- 735 • Postoperative Pain Dataset: Patients undergoing operations like tracheal intubations often
736 experience throat pain following treatment(Mchardy & Chung, 1999). This RCT was con-
737 ducted to test the efficacy of gargling with licorice solution prior to endotracheal intubation
738 on reducing postoperative sore throat, which is a common side-effect of the procedure.
739 The investigation involved 236 adult participants scheduled for elective thoracic surgeries
740 necessitating the use of double-lumen endotracheal tubes. The outcome we focus on is a
741 patient's throat pain 4 hours after surgery, measured on a discrete Likert scale from 0 to 7.
742 Additional covariates include a patient's gender, BMI, age, Mallampati score, preoperative
743 pain, surgery size and smoking status. Here, the effect of the treatment is to reduce the
744 amount of throat pain, hence the treatment effect is negative. In order to maintain consis-
745 tency with other plots, we present results with the sign for treatment effect reversed. Since
746 we do not have a measured value for throat pain at baseline, we again train a random for-
747 est model on rows corresponding to patients who did not receive the treatment with their
748 throat pain at endline being the outcome variable. The prediction offered by this model for
749 every patient is then used as a proxy for their baseline throat pain and consequently as the
750 baseline risk. This follows the intuition that patients with more severe throat pain require
the treatment more than their co-patients.
- 751 • Acupuncture Dataset: This RCT aimed to determine the effect of acupuncture therapy on
752 headache severity in patients with chronic headaches. These measures were assessed at
753 randomization, 3 months post-randomization, and 1 year post-randomization. We focus on
754 headache severity 1 year post-randomization. Headache severity is measured on a discrete
755 Likert scale from 0 to 5. The dataset consists of data from 401 patients with covariates
including patient age, sex, chronicity(number of years of headache severity) and whether

756 the headaches were diagnosed as migraines or not. Here again, the effect of the treatment
 757 is to reduce the severity of headaches, hence the treatment effect is negative. In order to
 758 maintain consistency with other plots, we present results with the sign for treatment effect
 759 reversed. We estimate headache severity at baseline $E(Y(0)|X]$ using a machine learning
 760 model and use it as a proxy for baseline risk, following the intuition that patients with more
 761 severe headaches need the treatment more.
 762

763 A.2 ADDITIONAL PLOTS



779 Figure 3: Comparison of risk-based assignment to biased treatment effect based assignment for the
 780 STAR dataset, with fixed budget of 30% and 40% of the population respectively.
 781



800 Figure 4: Comparison of risk-based assignment to biased treatment effect based assignment for the
 801 Ultra Poor dataset, with fixed budget of 30% and 40% of the population respectively.
 802
 803
 804
 805
 806
 807
 808
 809

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

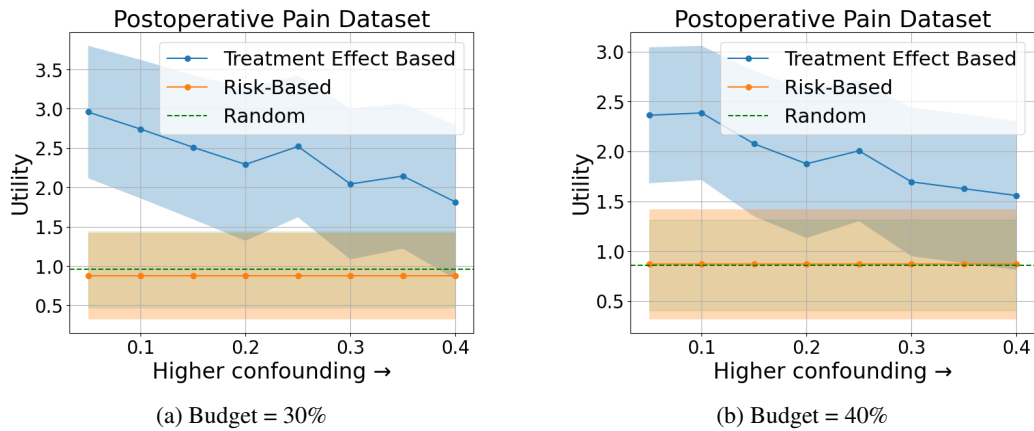


Figure 5: Comparison of risk-based assignment to biased treatment effect based assignment for the Postoperative Pain dataset, with fixed budget of 30% and 40% of the population respectively.

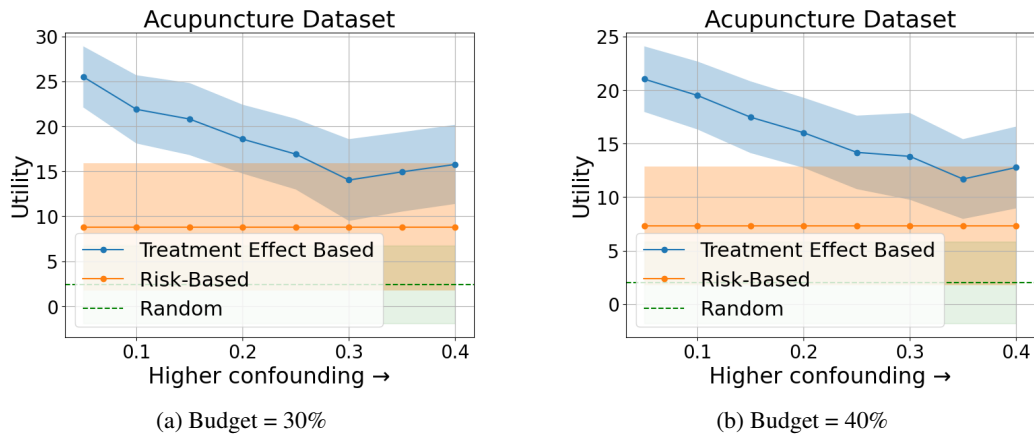


Figure 6: Comparison of risk-based assignment to biased treatment effect based assignment for the Acupuncture dataset, with fixed budget of 30% and 40% of the population respectively.

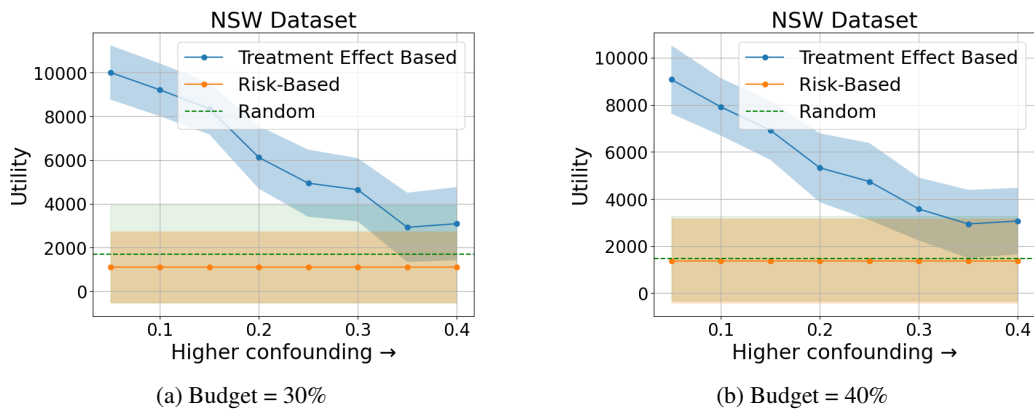


Figure 7: Comparison of risk-based assignment to biased treatment effect based assignment for the NSW dataset, with fixed budget of 30% and 40% of the population respectively.

B ROBUSTNESS CHECK: DISJOINT NUISANCE FUNCTION ESTIMATION

In the main paper, we employ a sample splitting approach, in line with the literature on doubly robust CATE estimation and previous work on policy optimization/comparison. In particular, our strategy is equivalent to the cross-validation strategy used in Athey & Wager (2020) to evaluate learned policies on RCT data. One potential concern with this procedure is that overlapping sets of data are used both to train the treatment effect-based targeting policy and to fit the nuisance functions used for evaluation. In theory this should not be an issue because, with RCT data, the DR estimator used for evaluation will be unbiased even if the outcome regression is mis-estimated (because the propensity score is guaranteed to be well-specified).

However, as an additional robustness check, in this section we completely separate the data used for training the CATE estimator for the treatment policy and the estimator used for policy evaluation. In particular, we set aside half the data for use only for evaluation and one half only for training allocation policies. The training portion of the data is used to fit a treatment assignment policy which treats individuals with the highest estimated treatment effects via a DR learner (up to the budget constraint). Then, we evaluate the average treatment effect on the treated of each policy using the evaluation split. Within the evaluation split, we also fit a DR estimator, but the nuisance functions for the this DR estimator are now fit on an entirely disjoint set of data from that used to optimize the policy.

This procedure has two drawbacks, both related to the reduction in sample size. First, it substantially reduces the amount of data available for training the allocation policy, and so potentially underestimates the effectiveness of causal targeting. Second, it makes less efficient use of data for evaluation as well, which tends to inflate the variance of the evaluation and the size of confidence intervals. We only implement this process for the STAR and NSW datasets because the other datasets are too small to credibly train the allocation policy after setting aside half the data. Because of these drawbacks, we do not expect the results from this experiment to be identical to our main analysis, but we include it to check if the qualitative conclusions are similar.

We indeed see observe similar high-level conclusions with the main analysis (Figure 8, 9, 10). The effectiveness of causal targeting is reduced, but it still outperforms risk-based targeting by a substantial margin at low levels of confounding even when the decision maker has egalitarian preferences.

Table 2 shows the value of α for each level of confounding at which the point estimates for utility for risk-based targeting exceed those for treatment effect based targeting for the two datasets we consider, similar to Table 1

Table 2: Values of α for different k at which risk-based targeting outperforms treatment effect based targeting. 'na' indicates no such α was found

Dataset	5%	10%	15%	20%	25%	30%	35%	40%
NSW	na	4	5	2.5	0	0	0	0
STAR	4	3.5	3.5	2.5	2	1.5	0.5	1

B.1 ADDITIONAL PLOTS

Additional plots with different budgets using disjoint nuisance functions are shown in Figures 9, 10.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

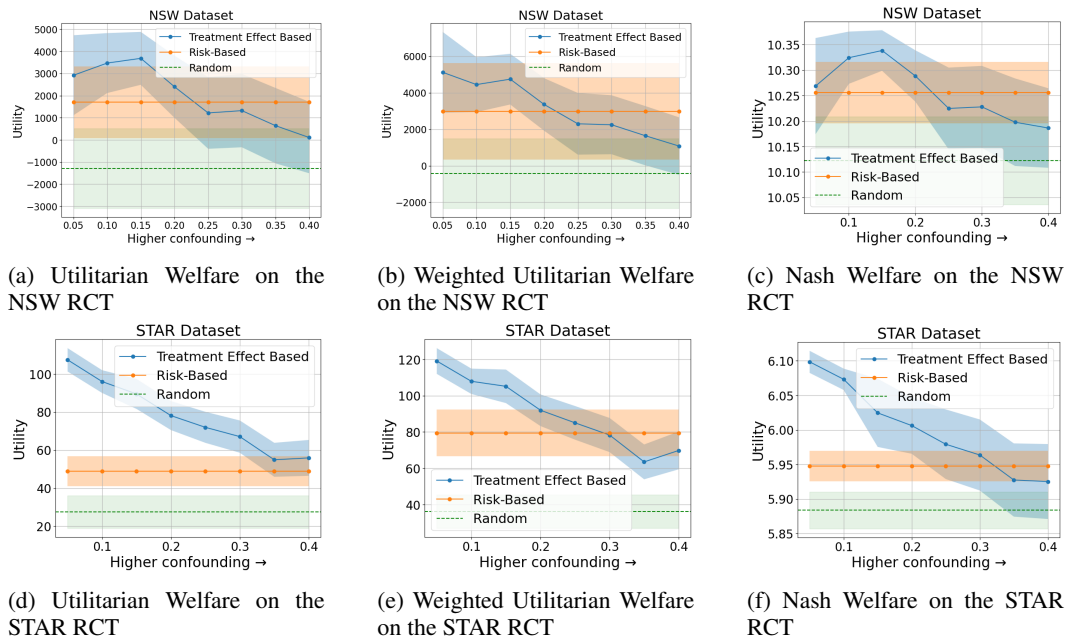


Figure 8: Comparison of risk-based targeting to biased treatment effect-based targeting by plotting the benefit offered by each policy against the amount of data systematically removed from the RCT to introduce confounding. The DR estimator is trained on a dedicated training sample and targeting decisions are made on a separate sample.

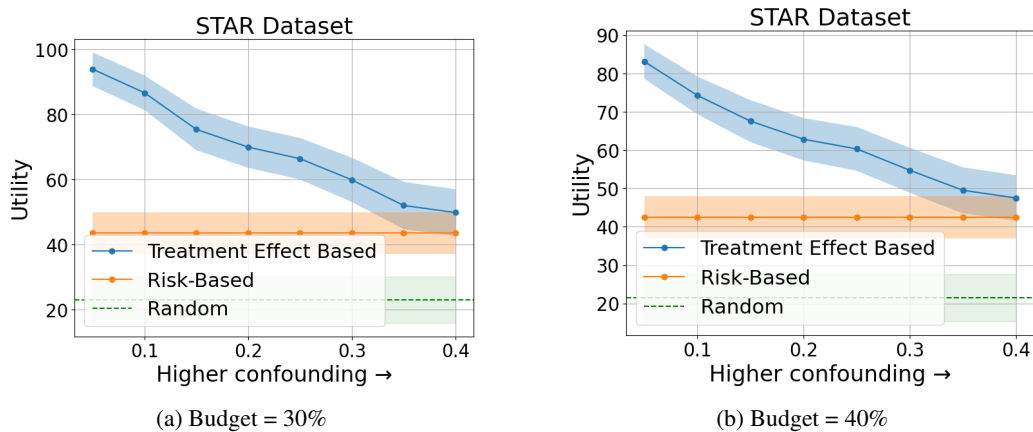


Figure 9: Comparison of risk-based assignment to biased treatment effect based assignment for the STAR dataset, with fixed budget of 30% and 40% of the population respectively. The DR estimator is trained on a dedicated training sample and targeting decisions are made on a separate sample.

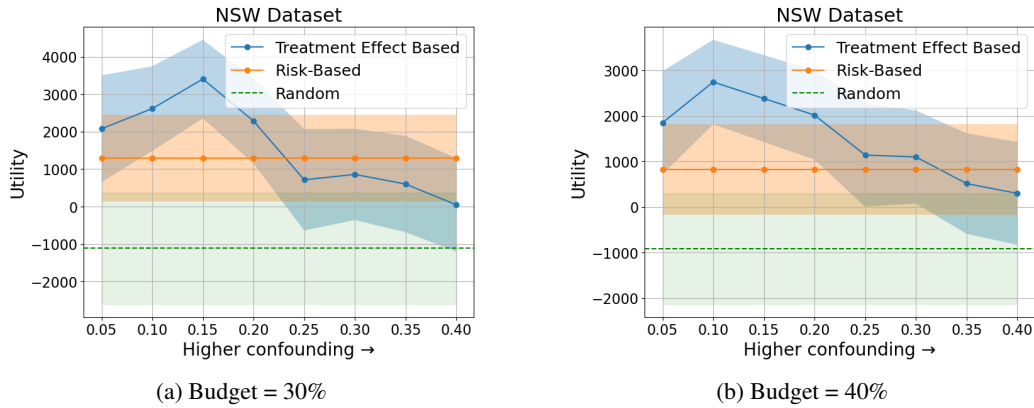


Figure 10: Comparison of risk-based assignment to biased treatment effect based assignment for the NSW dataset, with fixed budget of 30% and 40% of the population respectively. The DR estimator is trained on a dedicated training sample and targeting decisions are made on a separate sample.

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025