

Partial Identification and Robust Treatment Choice: An Application to Young Offenders

Jörg Stoye
Department of Economics
New York University
j.stoye@nyu.edu

April 19, 2008

Abstract

This paper applies recently developed methods to robust assessment of treatment outcomes and robust treatment choice based on nonexperimental data. The substantive question is whether young offenders should be assigned to residential or nonresidential treatment in order to prevent subsequent recidivism. A large data set on past offenders exists, but treatment assignment was by judges and not by experimenters, hence counterfactual outcomes are not identified unless one imposes strong assumptions.

The analysis is carried out in two steps. First, I show how to compute identified bounds on expected outcomes under various assumptions that are too weak to restore conventional identification but may be accordingly credible. The bounds are estimated, and confidence regions that take current theoretical developments into account are computed. I then ask which treatment to assign to future offenders if the identity of the best treatment will not be learned from the data. This is a decision problem under ambiguity. I characterize and compute decision rules that are asymptotically efficient under the minimax regret criterion. The substantive conclusion is that both bounds and recommended decisions vary significantly across the assumptions. The data alone do not permit conclusions or decisions that are globally robust in the sense of holding uniformly over reasonable assumptions.

Keywords. Partial identification, bounds, statistical decision rules, treatment choice, treatment evaluation, minimax regret.

1 Introduction

This paper shows how recent results on identification, estimation, and treatment choice in situations of partial identification can be brought to bear on a practical example, using real-world data to analyze a question of real-world interest. The question is: What kind of treatment should be used on young offenders to prevent recidivism? The choice will be between confinement in residential treatment facilities and diversion to nonresidential treatment. The data are a large sample of young offenders from Utah who received one of the treatments and whose subsequent behavior was observed. If the offenders had been assigned to treatments at random, then one could trivially determine which treatment was more successful. The fundamental problem motivating the analysis is that in fact, treatment was assigned by judges whose intention was hardly to create a statistical experiment. Hence, there is a missing data problem: Counterfactual outcomes for the treatment group are unobserved and cannot be proxied for by observed outcomes for the control group, and vice versa.

Such problems are well known in the statistical and econometric literature on treatment evaluation.¹ The standard solution is to propose a model of the missing data. If the model is specified sufficiently tightly, then identification is restored and counterfactual outcomes can (asymptotically) be backed out. This approach has much merit but also an important limitation: The credibility of conclusions drawn from such analysis is bounded above by the credibility of the assumptions about missing data. The purpose of this paper is to demonstrate how recently developed methods can be used to generate interesting insights with no, or less, such assumptions. Of course, getting answers that do not rely on identifying assumptions comes at a price. When it comes to assessing the outcome distributions induced by different treatments, the price is that only bounds on these quantities can be estimated. When it comes to recommending treatment for future populations, the price is that many different assignment rules are admissible. To select one of them, one must commit to a potentially controversial principle of decision making under ambiguity. This principle could be Bayesianism or maximin utility, but I will emphasize minimax regret in the form recently proposed by [10].

While I provide some new theoretical results, these are not the paper's core contribution. The main point is rather to show how different very recent theoretical results, few of which have seen much application so far, can be combined in the analysis of a practical problem. Specifically, I rely on identification analysis by [14] as well as new results; I estimate partially identified parameters using meth-

¹See [6] and [18] for overviews and [4] for a recent, encyclopedic treatment and further references.

	nonresidential	residential	totals
no recidivism	4774	333	5107
recidivism	6977	1113	8090
totals	11751	1446	13197

Table 1: Contingency table of the data.

ods developed in [7] and [25]; I apply recent [11][21] analyses of treatment choice when identified population parameters are known; and I use [5] to generate asymptotically minimax regret efficient treatment rules. Regarding the substantive question, it turns out that treatment recommendations vary much with what one is willing to assume, as well as with decision theoretic commitments one is willing to make. The data alone are not able to inform decisions.

The most important precursor for this paper is [14], who use the same data; indeed, section 3 of this paper essentially updates their analysis in the light of subsequent developments. There is no decision theoretic component to [14] however. To my knowledge, the only other empirical applications of minimax regret are [13] and some passages in [12] for situations with partial identification and [2] for situations with conventional identification.

2 The Identification Problem

Substantively, this paper is about the impact of the juvenile justice system on delinquency and what to do about it, a question that has engendered lively academic and nonacademic debate; see [14] for an overview of the sociological literature. Specifically, should young offenders be assigned to residential facilities (i.e. prison-like, although typically not prisons in the usual sense) or to nonresidential treatment? The question will be analyzed by considering the outcomes actually experienced by young offenders who were assigned one of the two treatments in Utah. Thus, I exclude from analysis the fact that different treatments might have different incentive effects on potential future offenders.

I re-analyze the data used in [14]. They are lifted from the National Juvenile Court Data Archive and collect observations on male offenders born between 1970 and 1974 who came into contact with the criminal justice system before age 16 and who were eventually found guilty on a charge that would have been a criminal offense under adult law. The data reveal whether the offenders received residential or nonresidential treatment and whether they generated a new referral within the 24 months following the date of treatment. Generating such a referral is counted as recidivism. Table 1 summarizes key aspects of the data.

I will maintain the assumption that the sample is a random sample from the population of young offenders. If treatment had furthermore been assigned randomly as in an experiment, conclusions would be obvious: The probability of recidivism equals 59% in the nonresidential treatment group and 77% in the residential treatment group. Sampling errors are small compared to the difference between these numbers. Nonresidential treatment is preferable.

Of course, things are not so easy. The choice between residential and nonresidential treatment was made by judges, and it is unlikely that they intended to create an experiment. Therefore, outcomes experienced by the residential treatment group may not be a good proxy for outcomes that would have been experienced by the other group (and vice versa).

To analyze this formally, I use a *potential outcomes* framework [18]. Let the random variable $T \in \{r, n\}$ denote treatment, with a realization of $t = r$ corresponding to residential treatment.² The random variable Y_t denotes the potential outcome that would be experienced by a randomly selected young offender if assigned to treatment t . A realization of $Y_t = c$ (for “crime”) denotes recidivism; to avoid confusion with n for “nonresidential,” the other outcome will be denoted as $(Y_t \neq c)$. Importantly, there are two potential outcomes (Y_r, Y_n) for every young offender, but one can only observe the treatment that he actually received and the outcome he actually experienced, i.e. realizations (t, y_t) . Consequently, $\Pr(T = r)$, $\Pr(Y_r = c|T = r)$, and $\Pr(Y_n = c|T = n)$ are identified – they will be learned with arbitrary precision as samples grow large and are revealed with rather high precision by the existing sample. In contrast, the counterfactual probabilities $\Pr(Y_n = c|T = r)$, and $\Pr(Y_r = c|T = n)$ are unobservable, and no direct learning about them occurs in samples of arbitrary size. The obvious consequence is that the unconditional probabilities $\Pr(Y_n = c)$ and $\Pr(Y_r = c)$ are not identified either. Note that this is essentially a missing data problem, namely, observations on counterfactual outcomes are missing.

The problem is conventionally resolved by proposing a model of the missing data that restores identification. The most straightforward such model is to posit that the missing data are ignorable, i.e. that $\Pr(Y_t = c|T = r) = \Pr(Y_t = c|T = n), t = r, n$. This assumption is justified if the data were generated by an experiment, i.e. if T was chosen at random. In that case, T would be independent of (Y_r, Y_n) , hence ignorability would obtain. This is why randomized experiments make statistical analysis relatively easy.

With nonexperimental data, ignorability is typically not a plausible assumption. One may be willing to impose a different model of the missing data that is tight enough to ensure identification. In economics, a typical method would be to write

²I will generally use capital letters for random variables and small letters for realizations.

an explicit model that specifies the behavior of both judges and offenders up to some parameters which are then estimated. While this approach has a distinguished tradition in both statistics and economics, the credibility of conclusions obtained through it is bounded by the credibility of the identifying assumptions. The purpose of this paper is to showcase a recently developed set of tools that allows one to generate a lot of insight with no model of the missing data, or with a model that is very weak and therefore highly credible, but fails to restore identification. The core idea is that even such weak models generate *partial identification*: The data do not (asymptotically) reveal the exact parameter of interest, but they do reveal nontrivial information about it. The question is how to rigorously formalize this idea and how to base estimation on it. This will be done in section 3 of this paper, which will derive and estimate numerous bounds on the parameters of interest. I will also find out how much deviation from random assignment (in the sense of distortion of odds ratios) is needed to overturn the conclusions one would draw from ignoring the missing data problem. This part of the analysis can alternatively be interpreted as an exercise in global sensitivity analysis, or in bounding posterior expectations in a specific instance of a robust (multiple prior) Bayesian setup.

The insights generated by partially identified models are frequently insufficient to identify one treatment as unambiguously optimal; indeed, this case will occur here. Therefore, partial identification also raises new questions about treatment choice. In section 4 of the paper, I combine the minimax regret decision criterion with the partially identifying assumptions to arrive at specific treatment recommendations. As was previously found with minimax regret [11][20][21][22][24], these recommendation will often be to randomize treatment. Also, the assumptions one is willing to make strongly affect treatment recommendations. Section 5 concludes, and an appendix collects proofs of selected results.

3 Bounds on Parameters: Identification, Estimation, and Inference

In this section, I show what can be learned about $\Pr(Y_r = c)$ and $\Pr(Y_n = c)$ using no or weak assumptions about missing data. Thus, the analysis is one of *partial identification* as recently summarized by [9]. I first introduce partial identification by demonstrating assumption-free (with respect to the missing data) bounds. Then I motivate a number of partially identifying assumptions, derive what they would imply if identified quantities were known, and compute according estimates and confidence regions.

The basic idea behind partial identification is that even when quantities are not identified in the usual sense, the data generating process may reveal some infor-

mation about them. This is easy to see in the present example: Basic probability calculus implies that

$$\begin{aligned}\Pr(Y_r = c) &= \Pr(Y_r = c|T = n) \Pr(T = n) + \Pr(Y_r = c|T = r) \Pr(T = r) \\ \Pr(Y_n = c) &= \Pr(Y_n = c|T = n) \Pr(T = n) + \Pr(Y_n = c|T = r) \Pr(T = r),\end{aligned}$$

hence

$$\begin{aligned}\Pr(Y_r = c|T = r) \Pr(T = r) &\leq \Pr(Y_r = c) \\ &\leq \Pr(T = n) + \Pr(Y_r = c|T = r) \Pr(T = r) \\ \Pr(Y_n = c|T = n) \Pr(T = n) &\leq \Pr(Y_n = c) \\ &\leq \Pr(Y_n = c|T = n) \Pr(T = n) + \Pr(T = r)\end{aligned}$$

as originally derived in [8]. These bounds only depend on identified quantities, hence one will learn them from the data as samples grow large; in finite samples, they can be estimated.

Of course, such worst-case bounds may be rather wide. In this paper's application, replacing population expectations with sample means yields estimated bounds of

$$\begin{aligned}0.025 &\leq \Pr(Y_r = c) \leq 0.916 \\ 0.362 &\leq \Pr(Y_n = c) \leq 0.471,\end{aligned}$$

where the stark difference in length of bounds is due to the fact that most offenders received nonresidential treatment, thus its effect is better identified.

It is possible to create tighter bounds by committing to some assumption about how missing data were generated without going back all the way to identification. For example, [14] translate popular conjectures about decision making by judges into the following two assumptions.

Definition: Outcome Optimization

A judge is said to optimize outcomes if

$$\begin{aligned}\Pr(Y_r = c|T = r) &\leq \Pr(Y_n = c|T = r) \\ \Pr(Y_n = c|T = n) &\leq \Pr(Y_r = c|T = n).\end{aligned}$$

Definition: Skimming

A judge is said to practice skimming if

$$\begin{aligned}\Pr(Y_r = c|T = r) &\geq \Pr(Y_r = c|T = n) \\ \Pr(Y_n = c|T = r) &\geq \Pr(Y_n = c|T = n).\end{aligned}$$

Both assumptions presume that judges are able to assess individual offenders' types, and hence their potential outcomes, pretty well.³ They differ with respect to what judges are supposed to make of this information. The story behind outcome optimization is that they minimize recidivism. Thus, an offender is assigned to residential treatment iff this treatment minimizes his probability of recidivism. The story behind skimming is that judges are able to rank cases according to “toughness,” that tougher cases have the worse prognosis conditional on either treatment, and that judges assign the tougher cases to residential treatment, for example to deter future tough cases or to protect society from the present ones.

I will also analyze an assumption that is due to [17, chapter 4]:

Definition: Bounded Selection

Bounded selection with parameter $\kappa \in [1, \infty)$ obtains if

$$\frac{1}{\kappa} \frac{\Pr(Y_r = c|T = r)}{\Pr(Y_r \neq c|T = r)} \leq \frac{\Pr(Y_r = c|T = n)}{\Pr(Y_r \neq c|T = n)} \leq \kappa \frac{\Pr(Y_r = c|T = r)}{\Pr(Y_r \neq c|T = r)}$$

$$\frac{1}{\kappa} \frac{\Pr(Y_n = c|T = n)}{\Pr(Y_n \neq c|T = n)} \leq \frac{\Pr(Y_n = c|T = r)}{\Pr(Y_n \neq c|T = r)} \leq \kappa \frac{\Pr(Y_n = c|T = n)}{\Pr(Y_n \neq c|T = n)}.$$

Bounded selection constrains the divergence between observable and counterfactual odds ratios. Its stringency can be scaled by choosing κ : With $\kappa = 1$, it reduces to ignorable missing data, and it becomes entirely vacuous as $\kappa \rightarrow \infty$. [17] emphasizes its use as a local robustness (or sensitivity) tool that can reveal whether conclusions are knife-edge dependent on ignorability. This interpretation is less important here because ignorability is implausible to begin with. However, bounded selection with large or moderate κ is a potentially plausible assumption that contrasts with outcome optimization and skimming. The latter allow for very sharp differences between the two treatment groups, and at the same time presume that judges fully anticipated these differences. If either presumption appears dubious – i.e. differences between offenders are limited, or judges can differentiate between offenders only to a limited degree –, then bounded selection makes intuitive sense.⁴ While the assumption is not testable from this paper's data, one could calibrate plausible values for κ from auxiliary information, e.g. case studies of judges' behavior or data that record recidivism conditional on additional information that judges are likely to observe.

³Technically, one might want to assume that they know the individual realizations of potential outcomes, but it suffices that they observe characteristics of offenders that are not contained in the data but that allow for differentiated predictions.

⁴Bounded selection could be combined with either of the preceding assumptions to simultaneously constrain the direction and extent of selection effects, but this will not be pursued here.

The assumptions have in common that they fail to restore identification of previously unidentified quantities, but they imply restrictions on counterfactual outcomes and therefore induce a tightening of the worst-case bounds. At the same time, they can be argued to be more credible than fully identifying models because they merely impose some nonparametric restriction that often has clear substantive meaning (e.g., judges act in offenders' best interest), without the baggage of additional convenience assumptions. Their precise effect on bounds is as follows.

Proposition 1:

Bounds under Different Assumptions

Assume that outcome optimization holds, then

$$\begin{aligned} \Pr(Y_r = c|T = r) \Pr(T = r) + \Pr(Y_n = c|T = n) \Pr(T = n) \\ \leq \Pr(Y_r = c) \leq \Pr(Y_r = c|T = r) \Pr(T = r) + \Pr(T = n) \end{aligned}$$

$$\begin{aligned} \Pr(Y_r = c|T = r) \Pr(T = r) + \Pr(Y_n = c|T = n) \Pr(T = n) \\ \leq \Pr(Y_n = c) \leq \Pr(T = r) + \Pr(Y_n = c|T = n) \Pr(T = n) \end{aligned}$$

Assume that skimming holds, then

$$\Pr(Y_r = c|T = r) \Pr(T = r) \leq \Pr(Y_r = c) \leq \Pr(Y_r = c|T = r)$$

$$\Pr(Y_n = c|T = n) \leq \Pr(Y_n = c) \leq \Pr(T = r) + \Pr(Y_n = c|T = n) \Pr(T = n).$$

Assume that bounded selection holds, then

$$\begin{aligned} \Pr(Y_r = c|T = r) \Pr(T = r) + \frac{\Pr(Y_r = c|T = r)}{\kappa - (\kappa - 1) \Pr(Y_r = c|T = r)} \Pr(T = n) \\ \leq \Pr(Y_r = c) \leq \end{aligned}$$

$$\Pr(Y_r = c|T = r) \Pr(T = r) + \frac{\kappa \Pr(Y_r = c|T = r)}{1 + (\kappa - 1) \Pr(Y_r = c|T = r)} \Pr(T = n)$$

$$\begin{aligned} \frac{\Pr(Y_n = c|T = n)}{\kappa - (\kappa - 1) \Pr(Y_n = c|T = n)} \Pr(T = r) + \Pr(Y_n = c|T = n) \Pr(T = n) \\ \leq \Pr(Y_n = c) \leq \end{aligned}$$

$$\frac{\kappa \Pr(Y_n = c|T = n)}{1 + (\kappa - 1) \Pr(Y_n = c|T = n)} \Pr(T = r) + \Pr(Y_n = c|T = n) \Pr(T = n).$$

All of these bounds are tight, that is, they cannot be improved upon without further assumptions.

Assumption	LB on $\Pr(Y_n = c)$		UB on $\Pr(Y_n = c)$	
	95% CI	estimator	estimator	95% CI
worst-case	0.522	0.529	0.638	0.645
outcome optimization	0.606	0.613	0.638	0.645
skimming	0.586	0.594	0.638	0.645
bounded sel., $\kappa = 100$	0.523	0.533	0.638	0.645
bounded sel., $\kappa = 10$	0.535	0.543	0.631	0.638
bounded sel., $\kappa = 5$	0.546	0.553	0.625	0.632
bounded sel., $\kappa = 2$	0.567	0.574	0.610	0.618
bounded sel., $\kappa = 1$	0.585	0.594	0.594	0.603
	LB on $\Pr(Y_r = c)$		UB on $\Pr(Y_r = c)$	
	95% CI	estimator	estimator	95% CI
worst-case	0.080	0.084	0.975	0.977
outcome optimization	0.606	0.613	0.975	0.977
skimming	0.080	0.084	0.770	0.788
bounded sel., $\kappa = 100$	0.107	0.111	0.972	0.975
bounded sel., $\kappa = 10$	0.288	0.307	0.949	0.954
bounded sel., $\kappa = 5$	0.417	0.441	0.924	0.931
bounded sel., $\kappa = 2$	0.618	0.641	0.859	0.871
bounded sel., $\kappa = 1$	0.748	0.770	0.770	0.791

Table 2: Bounds on expected outcomes under different assumptions.

The first two sets of bounds in this proposition are due to [14]; the last one is new but similar to a result in [26]. Notice in particular that under outcome optimization, both $\Pr(Y_r = c)$ and $\Pr(Y_n = c)$ are bounded below by $\Pr(Y_r = c|T = r)\Pr(T = r) + \Pr(Y_n = c|T = n)\Pr(T = n)$, the aggregate probability of recidivism under the current treatment assignment scheme. This is as expected because judges are assumed to minimize expected recidivism. Also, bounds on the average treatment effect, the average treatment effect on the treated, and the average treatment effect on the untreated,

$$\begin{aligned}
ATE &\equiv \Pr(Y_r \neq c) - \Pr(Y_n \neq c) \\
ATT &\equiv \Pr(Y_r \neq c|T = r) - \Pr(Y_n \neq c|T = r) \\
ATU &\equiv \Pr(Y_r \neq c|T = n) - \Pr(Y_n \neq c|T = n)
\end{aligned}$$

as well as other parameters from the treatment evaluation literature follow immediately from proposition 1; they are not displayed for brevity.

To gauge the effect of partially identifying assumptions in the specific example, I estimate all bounds by replacing population expectations with sample means.

Table 2 presents the resulting estimates along with 95% confidence intervals for $\Pr(Y_n = c)$ respectively $\Pr(Y_r = c)$. These intervals were computed as suggested in [25]. Importantly, they are confidence regions for the true parameter and not for bounds on it, meaning that if θ_0 denotes the true probability of interest, and the interval $[\theta_l, \theta_u]$ denotes bounds on it, then the confidence intervals displayed fulfil $\Pr(\theta_0 \in CI) \rightarrow 95\%$ but *not* $\Pr([\theta_l, \theta_u] \subseteq CI) \rightarrow 95\%$. The latter would require larger confidence regions, but seems less relevant here because the quantities of ultimate (e.g., decision theoretic) interest are the probabilities and not the bounds.⁵

The conclusions one can draw about different treatments depend to a large degree on the assumptions one is willing to make about missing data. Notice in particular how bounds successively tighten as bounded selection is imposed with decreasing κ . Bounded selection with $\kappa = 2$ implies the conclusion one would also draw by assuming random assignment, namely, that nonresidential treatment is better.⁶ The other assumptions fail to identify the better of the two treatments. This raises the question which treatment to assign to future offenders. It is this question to which I now turn.

4 Treatment Choice: Maximin Utility and Minimax Regret

Analysis of treatment outcomes usually aims to inform treatment choice. In the present case, a social planner may have to assign treatments to future offenders. Formally, her problem is to pick a decision rule $(\delta_r, \delta_n) \in [0, 1]^2$, where δ_t specifies the probability with which to apply residential treatment conditional on $T = t$.⁷ In practice, this would mean that while judges continue to assign offenders to one of two treatment groups, the decision on what to do with those groups is taken away from them. Interior probabilities $\delta_t \in (0, 1)$ correspond to randomized decision

⁵The distinction between these types of confidence regions is due to [7]. Technically, I use the confidence interval CI_α^3 from [25]. By proposition 3 and lemma 3 in [25], this interval is valid if (i) there exist estimators $(\hat{\theta}_l, \hat{\theta}_u)$ of (θ_l, θ_u) that are \sqrt{N} -consistent, uniformly asymptotically jointly normal, and ordered (i.e. $\hat{\theta}_u \geq \hat{\theta}_l$ by construction) and (ii) their asymptotic variances and correlation coefficient $(\sigma_l^2, \sigma_u^2, \rho)$ can be uniformly consistently estimated by estimators $(\hat{\sigma}_l^2, \hat{\sigma}_u^2, \hat{\rho})$. Assuming that $[\theta_l, \theta_u]$ is in the interior of $[0, 1]$, it is easily verified that sample analogs fulfil (i) for all bounds considered here. Standard errors were computed as closed form functions of sample moments where feasible and bootstrapped with 100000 replications otherwise; these estimators are again easily verified to fulfil (ii). The tuning parameter b_N was chosen as in [3].

⁶It is easy to evaluate the “breakdown point,” that is, the smallest κ s.t. bounds on $\Pr(Y_r = c)$ and $\Pr(Y_n = c)$ fail to overlap. Ignoring estimation uncertainty, this point equals about 2.29.

⁷Being an example of a statistical decision rule, δ will generally depend on sample observations, but since I condition the analysis on a particular real-world sample, I suppress this dependence in notation. See [10][24] for more general treatments.

rules, i.e. to assigning treatment by coin tosses (where coins have bias δ). Non-randomized treatment choice corresponds to $\delta \in \{0, 1\}$; in particular, the decision rule that just implements judges' original assignments would be $(\delta_r, \delta_n) = (1, 0)$.

Randomized treatment assignment implies that different individuals that received the same assignment by judges may end up with different treatments a posteriori. This can be advantageous from a maximin utility or minimax regret perspective because it hedges risks; indeed, here as in other papers [11][20][21][22][24], minimax regret turns out to prescribe it. It might, however, raise legal or ethical concerns that are not explicated here (see [13] for a model that formally incorporates them). Furthermore, inviting judges to propose an assignment but then deviating from it might again encounter legal difficulties and might also invite strategic behavior by judges (more on this below). I therefore also consider a "pooled assignment" problem in which all offenders receive residential treatment with probability δ_{pooled} , independently of their realization of T . This corresponds to a mechanism where judges' assessments are not any more elicited and can, therefore, not be contradicted. Of course, it removes the social planner's ability to condition on whatever information was revealed by T .

Assume that the planner's objective is to minimize recidivism, hence if $\Pr(Y_r = c|T = n)$ were known, then she would set $\delta_r = 1$ if $\Pr(Y_r = c|T = r) > \Pr(Y_r = c|T = n)$ and $\delta_r = 0$ otherwise (and similar for δ_n ; δ_{pooled} would compare unconditional probabilities). If $\Pr(Y_r = c|T = n)$ were not known but identified, there would be an added layer of difficulty because it would have to be estimated; in typical cases, this would however be a routine problem. Without identification, things look different because the optimal treatment assignment might be unknown even in the limit. This would happen whenever identified bounds on probabilities fail to overlap, or equivalently, when the sign of the relevant treatment effect is not identified. From a decision theoretic point of view, it means that even the limit problem is one of ambiguity, i.e. it is characterized by uncertainty which cannot be described in probabilistic terms, and the planner will have to commit to some decision criterion that applies to such situations. I will consider two such criteria.

Definition: Maximin Utility and Minimax Regret

Assume that all identified quantities are known and let $\mathcal{S} \in [0, 1]^2$ collect the possible states of the world $s \equiv (\Pr(Y_r = c|T = n), \Pr(Y_n = c|T = r))$. Then a maximin utility treatment rule fulfils

$$\delta_t^{MU} \in \arg \min_{\delta \in [0,1]} \max_{s \in \mathcal{S}} \{ \delta \Pr(Y_r = c|T = t) + (1 - \delta) \Pr(Y_n = c|T = t) \}.$$

A minimax regret treatment rule fulfils

$$\delta_t^{MR} \in \arg \min_{\delta \in [0,1]} \max_{s \in \mathcal{S}} \{ \delta \Pr(Y_r = c|T = t) + (1 - \delta) \Pr(Y_n = c|T = t) \\ - \min \{ \Pr(Y_r = c|T = t), \Pr(Y_n = c|T = t) \} \}.$$

The optimal pooled rules δ_{pooled}^{MU} respectively δ_{pooled}^{MR} are defined analogously but using unconditional probabilities.

The definitions are not very general because they reflect adaptation to the present setting. For general formulations and historical, philosophical, as well as axiomatic discussion of the criteria, I refer to the large decision theoretic literature that has evolved around them.⁸ A brief intuition goes as follows. Maximin utility (respectively minimax loss) is designed to minimize the worst-case probability of recidivism. Minimax regret optimizes a worst-case scenario as well, but “worst-case” is defined differently: The loss incurred in state s is not the probability of recidivism, but the increase in this probability *relative to what could have been achieved given s* . Intuitively, minimax regret is not about how objectively bad a situation is, but about the damage that can be caused by making the wrong decision. Minimax regret was originally suggested by Savage [19] and recently attracted attention in econometrics mostly due to [10]. I will focus on minimax regret, partly because it is the lesser known criterion but recently received some revival, and partly because it is generally harder to compute, making feasibility of minimax regret analysis a current concern. Maximin utility rules for the present decision problem are easy to compute and are reported along the way to provide a comparison.

The obvious alternative to either decision rule is the Bayesian approach, namely to put a prior on \mathcal{S} . Technically, this is equivalent to imposing a missing data model, and it consequently leads to well-defined posterior distributions of (Y_r, Y_n) . Hence, it should be thought of less as a decision theoretic solution concept than as a way to restore identification. As such, it would be subject to the preceding section’s criticism of fully identifying assumptions. To avoid this criticism, one could opt for the robust Bayesian approach [1, chapter 4] and impose a set of priors. But whenever this approach fails to identify the better treatment, one is back to the decision problem considered in this section. Indeed, the difference between (nominally frequentist) minimax regret as considered here and the robust Bayesian Γ -minimax regret is largely semantic – partially identifying assumptions can be interpreted as characterizations of sets of priors. If the Bayesian approach is combined with a concern for robustness, then it does not avoid this section’s decision problem.

⁸The abstract treatment in [23] is tailored to applications in statistical decision theory. The definition of minimax regret in [22] is more general than but still close to the present one.

To compute the optimal rules, I first abstract from sampling uncertainty, i.e. I pretend that identified quantities are known. Minimax regret treatment rule for this idealized scenario can be derived from a result due to [13, proposition 4].

Lemma 1

Let bounds on $(\Pr(Y_r = c|T = t), \Pr(Y_n = c|T = t))$ be given by $\underline{\pi}_r \leq \Pr(Y_r = c|T = t) \leq \bar{\pi}_r$ and $\underline{\pi}_n \leq \Pr(Y_n = c|T = t) \leq \bar{\pi}_n$, then minimax regret treatment choice is

$$\delta_t^{MR} = \begin{cases} 0, & d_t^{MR} < 0 \\ d_t^{MR}, & 0 \leq d_t^{MR} \leq 1 \\ 1, & 1 < d_t^{MR} \end{cases}$$

$$d_t^{MR} = \frac{\bar{\pi}_n - \underline{\pi}_r}{\bar{\pi}_n + \bar{\pi}_r - \underline{\pi}_n - \underline{\pi}_r}.$$

The pooled decision rule δ_{pooled}^{MR} is obtained analogously but using bounds on unconditional probabilities.⁹

Minimax regret decision rules can be computed by substituting for $(\underline{\pi}_r, \bar{\pi}_r, \underline{\pi}_n, \bar{\pi}_n)$ in this lemma. The results are as follows.

Proposition 2:

Minimax Regret Treatment Choice under Different Assumptions

Assume that worst-case bounds apply, then

$$\begin{aligned} \delta_{pooled}^{MR} &= \Pr(Y_n = c|T = n) \Pr(T = n) + \Pr(Y_r \neq c|T = r) \Pr(T = r) \\ \delta_r^{MR} &= 1 - \Pr(Y_r = c|T = r) \\ \delta_n^{MR} &= \Pr(Y_n = c|T = n). \end{aligned}$$

Assume outcome optimization, then

$$\begin{aligned} \delta_{pooled}^{MR} &= \frac{\Pr(Y_r \neq c|T = r) \Pr(T = r)}{\Pr(Y_r \neq c|T = r) \Pr(T = r) + \Pr(Y_n \neq c|T = n) \Pr(T = n)} \\ \delta_r^{MR} &= 1 \\ \delta_n^{MR} &= 0. \end{aligned}$$

⁹Readers familiar with the treatment effect literature might find it interesting to note that if \underline{ATT} and \overline{ATT} denote bounds on ATT , then $d_r^{MR} = \overline{ATT}/(\overline{ATT} - \underline{ATT})$ and similarly for d_n^{MR} and d_{pooled}^{MR} (using ATU respectively ATE).

Assume skimming, then

$$\begin{aligned}
\delta_{pooled}^{MR} &= \\
\min &\left\{ \frac{\Pr(Y_n = c|T = n) \Pr(T = n) + \Pr(Y_r \neq c|T = r) \Pr(T = r)}{\Pr(Y_r = c|T = r) \Pr(T = n) + \Pr(Y_n \neq c|T = n) \Pr(T = r)}, 1 \right\} \\
\delta_r^{MR} &= \min \left\{ \frac{1 - \Pr(Y_r = c|T = r)}{1 - \Pr(Y_n = c|T = n)}, 1 \right\} \\
\delta_n^{MR} &= \min \left\{ \frac{\Pr(Y_n = c|T = n)}{\Pr(Y_r = c|T = r)}, 1 \right\}.
\end{aligned}$$

Define $p \equiv \Pr(T = r)$ and $\pi_t \equiv \Pr(Y_t = c|T = t)$ and assume bounded selection, then minimax regret is achieved by the projection onto $[0, 1]$ (as in lemma 1) of

$$\begin{aligned}
d_{pooled}^{MR} &= \frac{\pi_n \left(\frac{p\kappa}{1+(\kappa-1)\pi_n} + 1 - p \right) - \pi_r \left(p + \frac{1-p}{\kappa-(\kappa-1)\pi_r} \right)}{p \frac{(\kappa^2-1)(\pi_n-\pi_n^2)}{\kappa+(\kappa-1)^2(\pi_n-\pi_n^2)} + (1-p) \frac{(\kappa^2-1)(\pi_r-\pi_r^2)}{\kappa+(\kappa-1)^2(\pi_r-\pi_r^2)}} \\
d_r^{MR} &= \frac{\kappa^2 (\pi_n - \pi_n^2) + \kappa\pi_n^2 - \pi_r (\kappa + (\kappa - 1)^2(\pi_n - \pi_n^2))}{(\kappa^2 - 1)(\pi_n - \pi_n^2)} \\
d_n^{MR} &= \frac{\pi_n (\kappa + (\kappa - 1)^2(\pi_r - \pi_r^2)) - \pi_r - (\kappa - 1)\pi_r^2}{(\kappa^2 - 1)(\pi_r - \pi_r^2)}.
\end{aligned}$$

Some parts of proposition 2 have clear intuitions. Under outcome optimization, the conditional minimax regret treatment rule prescribes to implement the judges' decisions. This is to be expected because the assumption says that judges optimized outcomes to begin with. If skimming is assumed, the decision rule may assign all offenders to treatment r , but never assigns all offenders to treatment n . The reason is that if $T = r$ indicates the “tougher” population, then $\Pr(Y_r = c|T = r) \leq \Pr(Y_n = c|T = n)$ unambiguously implies that r is the better treatment.¹⁰ In contrast, no value of $(\Pr(Y_r = c|T = r), \Pr(Y_n = c|T = n))$ will unambiguously establish that n is better.

Proposition 2 identifies minimax regret treatment rules in the limit problem where identified quantities, and hence the bounds from proposition 1, are known. In practice, the rules must be estimated. I compute “plug-in estimators” $\widehat{\delta}_{pooled}^{MR}$, $\widehat{\delta}_r^{MR}$, and $\widehat{\delta}_n^{MR}$ by substituting sample means for population expectations. A justification for this approach is found in [5]: Under conditions fulfilled here, solving

¹⁰Proof: $\Pr(Y_r = c|T = n) \leq \Pr(Y_r = c|T = r) \leq \Pr(Y_n = c|T = n) \leq \Pr(Y_n = c|T = r)$, where the outer inequalities follow from skimming.

Assumption	$\widehat{\delta}_{\text{pooled}}^{\text{MR}}$	$\widehat{\delta}_r^{\text{MR}}$	$\widehat{\delta}_n^{\text{MR}}$
worst-case	0.55	0.23	0.59
outcome optimization	0.07	1	0
skimming	0.76	0.57	0.77
bounded sel., $\kappa = 100$	0.54	0.23	0.58
bounded sel., $\kappa = 10$	0.44	0.21	0.48
bounded sel., $\kappa = 5$	0.33	0.17	0.36
bounded sel., $\kappa = 2$	0	0	0
bounded sel., $\kappa = 1$	0	0	0

Table 3: Minimax regret treatment rules under different assumptions.

limit problems but with asymptotically efficient estimators substituted for population quantities leads to treatment rules that are asymptotically efficient as well. Sample means are maximum likelihood estimators of population expectations and, therefore, asymptotically efficient. Hence, plugging them into proposition 2 is an asymptotically minimax regret efficient procedure.

Results are displayed in table 3. They show that treatment choice much depends on assumptions made. An interesting transition can be seen with bounded selection, where $\kappa = 100$ essentially replicates the worst-case analysis, but non-residential treatment becomes the more attractive, the less severe the distortion through selective treatment assignment is assumed to be.¹¹

Two interesting features of the results are that (i) residential treatment is recommended with perhaps higher than expected frequency, and (ii) except for outcome optimization, $\widehat{\delta}_n^{\text{MR}} > \widehat{\delta}_r^{\text{MR}}$ and frequently even $\widehat{\delta}_n^{\text{MR}} > 0.5 > \widehat{\delta}_r^{\text{MR}}$, thus minimax regret tends to (stochastically) overturn the judges' decisions. The reason for both phenomena is that minimax regret values upside risks roughly symmetrically to downside risks. Since residential treatment is much less prevalent in the sample, its effect is less identified. Because empirical outcomes of either treatment are disappointing, this underidentification induces more upside risk than downside risk, leading to (i). The effect becomes more pronounced when one separates $\widehat{\delta}_r^{\text{MR}}$ and $\widehat{\delta}_n^{\text{MR}}$. In both cases, the expected outcome under the respective status quo treatment is identified, the other one is not. Since status quo outcomes are not very favorable, minimax regret tends to advise against the status quo treatment in both cases, leading to (ii).

Two implications of these features deserve comment. First, they are not nec-

¹¹Minimax regret “locks in” on treatment 0 as soon as κ is below the aforementioned threshold of 2.29.

essarily intuitive and may raise questions about minimax regret as a decision criterion. One possible interpretation of the result is that it is due to underspecification of prior information. Given the observed recidivism, one would not realistically expect either treatment to induce almost no recidivism if administered to the other group. Prior restrictions that reflect this will reduce the phenomena; this is illustrated by bounded selection with lower values of κ . Second, whenever the inversion occurs, minimax regret that conditions on judges' choice is not incentive compatible: If judges anticipate that their assessment of correct treatment will be (stochastically) overturned, they have an incentive to misrepresent it. As a result, it is hard to see how the decision rule could be implemented.

I conclude by briefly contrasting the maximin utility recommendation. Straightforward computations show that (plug-in estimators of) maximin utility treatment rules will completely focus on nonresidential treatment except if outcome optimization is assumed and treatment can condition on T , in which case maximin utility treatment choice is to implement the judges' assessment. Indeed, it is easily shown that under outcome optimization, $(\delta_r, \delta_n) = (1, 0)$ is the only admissible treatment rule – an unsurprising conclusion given that by assumption, judges know recidivism probabilities and optimize them already.

5 Conclusion

This paper illustrated the use of recent tools for robust estimation, inference, and decision making in a real-world application. I analyzed recidivism of young offenders in Utah as a problem of partial identification, derived bounds on the effect of different treatments under numerous assumptions, estimated these bounds, provided confidence regions, and analyzed the problem of assigning treatment to future offenders as a problem in robust decision making, specifically from a minimax regret perspective. The results are of some substantive interest, but part of the purpose was a proof of concept, namely to demonstrate that these new methods can be used on real-world problems.

I kept the analysis very simple on a number of dimensions. As a result, this paper can be replicated from the information given in table 1. To give the reader a better understanding of the range of the proof of concept, I will briefly comment on some simplifications that were not crucial and could easily be overcome.

First, the dataset contains a number of covariates, especially judicial district and number of prior offenses, that were used by [14] but not here. If they are believed to matter, a minimax regret decision maker can – and should, even in small samples [24, proposition 3] – condition on them by stratifying the sample. If one is willing to make specific assumptions about the effect of covariates, then one could

use them as “instrumental variables” or “monotone instrumental variables” as in [15]. To briefly elaborate the former idea, assume that conditional probabilities of recidivism do not vary across judicial districts, but judges’ behavior does. Then one could tighten the bounds by evaluating them separately across districts and taking the intersection of the resulting bounds. To understand monotone instrumental variables, assume that recidivism probability increases in the number of prior offenses. Then one could first evaluate bounds separately conditional on the number of prior offenses and then refine them by “ironing out” any nonmonotonicities. Both approaches would lead to sharper inference, although at the cost of introducing additional assumptions as well as nontrivial complications in estimation and inference.

Second, I compared two treatments, but many other applications will feature more treatment options. This is not a problem because the generalization of lemma 1 to finitely many treatments is known [22]. A caveat is that closed-form expressions along the lines of proposition 2 will not be possible any more, but numerical evaluation remains easy.

Third, I evaluated asymptotically efficient approximations to minimax regret treatment rules rather than exact finite sample rules. This may not strike the reader as a limitation at all, because it mirrors standard statistical practice. It is worth mentioning, however, that finite sample minimax regret treatment rules for closely related problems have been discovered in ongoing research [20][22][24]. Having said that, this paper’s decision problem is somewhat more complex and may not be amenable to finite sample analysis.

Fourth, the empirical example featured binary outcomes. More generally, one might be interested in the expected value of some outcome measure that is distributed on the unit interval, with high outcomes being desirable. For all results except the closed-form bounds induced by bounded selection, bounds on such expectations follow by substituting $\mathbb{E}(Y_t|T = t)$ for $\Pr(Y_t \neq c|T = t)$ throughout. Bounded selection induces bounds that can be expressed as solutions to optimization problems but not in closed form. All in all, the analysis easily generalizes to bounded, although not unbounded, outcomes. This is also true for the existing finite sample results.

As a final remark, this paper used frequentist language, but the substantive issues are orthogonal to the Bayesian vs. frequentist divide. Bayesian estimation under partial identification is analyzed in current research [16]. As for decision making, I reiterate that imposing a prior on missing data models amounts to making identifying assumptions and is, therefore, prone to the robustness critique. As soon as this critique is accommodated, be it by sets of priors or by interval probabilities, the treatment choice problem from section 4 will be encountered.

Acknowledgements

Special thanks to Chuck Manski for sharing his data. I also thank two anonymous referees and a guest editor for helpful comments. Of course, any and all errors are mine. Address: Jörg Stoye, Department of Economics, New York University, 19 W. 4th Street, New York, NY 10012, U.S.A.

Proofs

Proposition 1 I show the bounds on $\Pr(Y_r = c)$ given bounded selection, the other arguments are similar. By assumption,

$$\frac{1}{\kappa} \frac{\Pr(Y_r = c|T = r)}{\Pr(Y_r \neq c|T = r)} \leq \frac{\Pr(Y_r = c|T = n)}{\Pr(Y_r \neq c|T = n)} \leq \kappa \frac{\Pr(Y_r = c|T = r)}{\Pr(Y_r \neq c|T = r)}.$$

Separate algebraic transformation of the l.h. and r.h. inequalities yields

$$\begin{aligned} \frac{\Pr(Y_r = c|T = r)}{\kappa - (\kappa - 1) \Pr(Y_r = c|T = r)} &\leq \Pr(Y_r = c|T = n) \\ &\leq \frac{\kappa \Pr(Y_r = c|T = r)}{1 + (\kappa - 1) \Pr(Y_r = c|T = r)}. \end{aligned}$$

The bounds can be generated by substituting from this finding into $\Pr(Y_r = c) = \Pr(Y_r = c|T = n) \Pr(T = n) + \Pr(Y_r = c|T = r) \Pr(T = r)$. They are tight because either of the above bounds on $\Pr(Y_r = c|T = n)$ is contained in $[0, 1]$; hence, absent further restrictions, $\Pr(Y_r = c|T = n)$ can attain either value. ■

Proposition 2 The expressions for δ_{pooled}^{MR} follow by substituting from proposition 1 into lemma 1. The expressions for $(\delta_r^{MR}, \delta_n^{MR})$ follow similarly from bounds on $\Pr(Y_r = c|T = n)$ and/or $\Pr(Y_n = c|T = r)$ that follow immediately from the different assumptions. As an illustration, I show δ_n^{MR} under bounded selection. Note that $\Pr(Y_n = c|T = n)$ is identified and that bounds on $\Pr(Y_r = c|T = n)$ were stated in the preceding proof. Substituting into lemma 1 yields

$$d_n^{MR} = \frac{\Pr(Y_n = c|T = n) - \frac{\Pr(Y_r = c|T = r)}{\kappa - (\kappa - 1) \Pr(Y_r = c|T = r)}}{\frac{\kappa \Pr(Y_r = c|T = r)}{1 + (\kappa - 1) \Pr(Y_r = c|T = r)} - \frac{\Pr(Y_r = c|T = r)}{\kappa - (\kappa - 1) \Pr(Y_r = c|T = r)}}.$$

After some algebra, this yields the expression in the proposition. ■

References

- [1] J.O. Berger. *Statistical Decision Theory and Bayesian Analysis (2nd Edition)*. Springer Verlag, New York, 1985.
- [2] P. Eozenou, J. Rivas, and K.H. Schlag. “Minimax Regret in Practice – Four Examples on Treatment Choice.” Preprint, European University Institute, 2006.
- [3] Y. Fan and S. Park. “Confidence Sets for Some Partially Identified Parameters.” Preprint, Vanderbilt University, 2007.
- [4] J.J. Heckman and E.J. Vytlacil. “Econometric Evaluation of Social Programs, Part I: Causal Models, Structural Models, and Econometric Policy Evaluation.” J.J. Heckman and E.E. Leamer (Eds.) *Handbook of Econometrics Volume 6B*. Elsevier, Amsterdam, 2007a.
- [5] K. Hirano and J.R. Porter. “Asymptotics for Statistical Treatment Rules.” Preprint, University of Arizona and University of Wisconsin-Madison, 2006.
- [6] P.W. Holland. “Statistics and Causal Inference.” *Journal of the American Statistical Association* 81: 945-960, 1986.
- [7] G. Imbens and C.F. Manski. “Confidence Intervals for Partially Identified Parameters.” *Econometrica* 72: 1845-1857, 2004.
- [8] C.F. Manski. “Nonparametric Bounds on Treatment Effects.” *American Economic Review (Papers and Proceedings)* 80: 319-323, 1990.
- [9] — *Partial Identification of Probability Distributions*. Springer Verlag, New York, 2003.
- [10] —. “Statistical Treatment Rules for Heterogeneous Populations.” *Econometrica* 72: 1221-1246, 2004.
- [11] —. “Minimax-Regret Treatment Choice with Missing Outcome Data.” *Journal of Econometrics* 139: 105-115, 2007.
- [12] —. *Identification for Prediction and Decision*. Harvard University Press, Cambridge, MA, 2008a.
- [13] —. “Adaptive Partial Policy Innovation: Coping with Ambiguity through Diversification” Preprint, Northwestern University, 2008b.

- [14] C.F. Manski and D.S. Nagin. “Bounding Disagreements about Treatment Effects: A Case Study of Sentencing and Recidivism.” *Sociological Methodology* 29: 99-137, 1998.
- [15] C.F. Manski and J.V. Pepper. “Monotone Instrumental Variables: With an Application to the Returns to Schooling.” *Econometrica* 68: 997-1010, 2000.
- [16] H.R. Moon and F. Schorfheide. “A Bayesian Look at Partially-Identified Models.” Preprint, University of Southern California and University of Pennsylvania, 2007.
- [17] P.R. Rosenbaum. *Observational Studies (2nd Edition)*. Springer Verlag, New York, 2002.
- [18] D.B. Rubin. “Causal Inference Using Potential Outcomes: Design, Modeling, Decisions.” *Journal of the American Statistical Association* 100: 322-331, 2005.
- [19] L.J. Savage. “The Theory of Statistical Decisions.” *Journal of the American Statistical Association* 46: 55-67, 1951.
- [20] K.H. Schlag. “Eleven.” Preprint, European University Institute, 2006.
- [21] J. Stoye. “Minimax Regret Treatment Choice with Incomplete Data and Many Treatments.” *Econometric Theory* 23: 190-199, 2007a.
- [22] —. “Minimax Regret Treatment Choice with Finite Samples and Missing Outcome Data.” G. de Cooman, J. Veinárová, and M. Zaffalon (Eds.) *Proceedings of the Fifth International Symposium on Imprecise Probability: Theories and Applications*, Prague 2007b.
- [23] —. “Statistical Decisions under Ambiguity.” Preprint, New York University, 2007c.
- [24] —. “Minimax Regret Treatment Choice with Finite Samples.” Preprint, New York University, 2007d.
- [25] —. “More on Confidence Intervals for Partially Identified Parameters.” Preprint, Centre for Microdata Methods and Practice (CEMMAP), London, 2008.
- [26] L. Wasserman and J.B. Kadane. “Computing Bounds on Expectations.” *Journal of the American Statistical Association* 87: 516-522, 1992.