

Statistical Decision Rules in Econometrics*

Keisuke Hirano

Department of Economics
Pennsylvania State University
hirano@psu.edu

Jack R. Porter

Department of Economics
University of Wisconsin
jrporter@ssc.wisc.edu

February 28, 2020

Abstract

Statistical decision rules map data into actions. Point estimators, inference procedures, and forecasting methods can be viewed as statistical decision rules. However, other types of rules are possible, such as rules for assigning individuals to treatments based on covariates, and methods for designing auctions. We discuss heuristics for constructing statistical decision rules, and survey results that characterize the properties of various classes of decision rules. Particular attention is paid to developing large-sample approximations to the distributions and associated risk properties of statistical decision rules.

Keywords: statistical decision theory; treatment assignment rules; limit experiments; risk.

*We are grateful to the Editor, the Referee, and Gary Chamberlain for helpful comments.

Outline

1. Introduction
2. General Setup and Evaluation of Decision Rules
 - 2.1 Setup
 - 2.2 Counterfactuals in the Wald Framework
 - 2.3 Classes of Statistical Decision Rules
 - 2.4 Evaluating Statistical Decision Rules
 - 2.5 Bounds and Large Sample Approximations for Decision Rules
3. Point Decisions
 - 3.1 Shift Equivariance in Point Estimation
 - 3.2 Asymptotics for Point Estimators
 - 3.3 Limits for Risk Functions
 - 3.4 Point Decisions for General Loss Functions
 - 3.5 Variable Selection and Shrinkage
4. Treatment Assignment Rules
 - 4.1 Treatment Assignment as a Decision Problem
 - 4.2 Welfare and Risk Analysis of Treatment Assignment Rules
 - 4.3 Local Asymptotics for Treatment Rules
 - 4.4 Other Treatment Assignment Problems
5. Other Topics
 - 5.1 Nonstandard Functionals
 - 5.2 Partial Identification
 - 5.3 Confidence Intervals and Sets
 - 5.4 Experimental and Data Collection Design
6. Conclusion
- * References

1 Introduction

A statistic can be defined simply as a quantity computed from a sample of data. Thus any statistical procedure is a mapping from the sample space of the statistical “experiment” into some space of “actions.” Standard procedures such as parameter estimators, hypothesis tests, and confidence intervals, can be viewed in this way, but so can other types of procedures, such as forecasting rules, allocation decisions based on past data, and methods for designing experiments. If we can also specify a criterion with which to evaluate such decision rules, then we can draw upon a rich set of results from statistical decision theory to analyze, compare, and choose among feasible rules.

Decision theoretic concepts and methods permeate statistics and econometrics with varying degrees of formality, and we can only touch on some aspects of their use in this chapter. Our goal is to survey some recent work in econometrics that adopts this framework to study different types of empirical problems. We review some classical problems, such as point estimation and confidence interval construction, but our goal is to highlight a richer set of decision problems that arise naturally in economic applications and can be handled within the statistical decision theory framework.

We primarily focus on analyzing rules from a frequentist perspective, calculating their *ex ante* expected performance, where the expectation is with respect to the sampling distribution of the data. In some cases, it is possible to characterize the finite sample properties of decision rules and obtain exact optimality results. However, in many economic applications, the richness of the hypothesized set of probability distributions, the complexity of the action space, or the nature of the evaluation criterion, make exact analysis infeasible. We consider how large sample approximations can be used to obtain approximate measures of performance and, in some cases, approximate optimality results for statistical decision rules. Exact and large sample analyses are often complementary. Finite sample results for carefully chosen special cases of the general decision problem can provide intuition and form the technical basis for large sample theory.

Statistical decision theory perspective adds, to the usual statistical model, a specific decision problem that the researcher faces. This requires specifying a set of actions (an action space) and a utility or loss function that pins down benefit or loss from taking a certain action given the (possibly unknown) state of nature. The need to incorporate some notion of utility or loss into probability calculations was recognized very early, for example in Daniel Bernoulli’s solution to the St. Petersburg paradox. Ramsey (1926) proposed to calculate expected utility with respect to subjective probability distributions. Subjective expected utility theory was developed more fully by later work, especially Savage (1972). The Savage formulation can be interpreted as providing an axiomatic justification for Bayesian statistical decision theory, which chooses actions to maximize posterior expected utility (or equivalently, minimize posterior expected loss). There is a long tradition of Bayesian analysis in statistics and econometrics, in which unknown model parameters are treated as random variables, whose distributions are updated based on observed data. The Bayesian approach is well motivated on decision-theoretic grounds, and

accommodates different action spaces and different loss or utility functions in a conceptually straightforward way. The development of computational methods including Markov chain Monte Carlo algorithms, and conceptual advances in Bayesian modeling in high-dimensional parameter spaces, have enabled researchers to apply Bayesian methods to a wide array of empirical problems in economics. For a recent survey of Bayesian econometrics, see Geweke, Koop, and van Dijk (2011). Other useful references for Bayesian statistics and econometrics include Bernardo and Smith (1994), Gelman, Carlin, Stern, Dunson, Vehtari, and Rubin (2013), and Lancaster (2004).

Another line of work adopted a frequentist view of probability, in which decision rules are evaluated on an *ex ante* basis according to their repeated sampling properties. The frequentist and subjectivist perspectives on decision rules are often fairly compatible; for example the frequentist properties of rules derived from Bayesian principles can be analyzed, and some methods derived from other considerations (such as certain types of shrinkage estimators) can be viewed as approximate or exact Bayes estimators under certain prior distributions for the parameters. Wald (1950) developed a framework for statistical decision theory in which rules are evaluated according to their sampling properties under different possible values for the unknown parameters of the statistical model. In this respect, it can be viewed as an extension of the Neyman-Pearson paradigm. In the Wald framework, a decision rule is characterized by its risk function – its expected loss as a function of the model parameter.

As with frequentist statistical theory in the Fisher or Neyman-Pearson paradigms, it may be difficult to obtain exact (finite-sample) results in complicated settings. A key insight of Le Cam (1972) was that asymptotic approximations could be applied to Wald's statistical decision theory framework. Rather than derive large-sample distributions for a specific estimator or inference procedure under a single value of the parameter, we approximate entire risk functions, and more generally we approximate the entire decision problem. To the extent that these approximations simplify the original problem, it may be possible to more easily compare different procedures and find optimal ones. This type of analysis has played an important role in developing asymptotic optimality theory for point estimators and hypothesis testing procedures, but its potential for analyzing a much wider range of statistical decision problems is perhaps underappreciated, and is a central motivation for this chapter.

The next section reviews the general statistical decision theory framework, which can be used to study many applications of statistical decision rules in economics. We set up the basic components of the framework, including the statistical model for the data, the space of possible actions of the decision-maker, and the criterion (in the form of a loss or welfare function) by which outcomes will be evaluated. We discuss how the Wald framework can accommodate policy analysis based on potential outcomes or counterfactuals. We examine some useful heuristics for constructing decision rules, and discuss methods for evaluating rules based on their risk or expected loss functions. We also introduce some of the large sample theory that will be used in the remainder of the chapter. We focus on obtaining approximate characterizations of risk or expected loss; for this purpose it is often useful to consider local parameter sequences that are chosen to ensure that the limiting version of the decision problem is nondegenerate and reflects the key features of the finite-sample decision problem.

The remaining sections then consider various applications of the general approach. Section 3 considers “point” decision problems, where the action space is a subset of a Euclidean space. The leading case is point estimation in parametric models, but we also want to extend classical results on point estimation to handle other decision problems, for example the problem of choosing a reserve price for a first-price auction. We focus on finite-sample theory for decision problems that have a translation-equivariant form, and large-sample results that rely on an asymptotic translation structure. This encompasses the classic locally asymptotically normal case but also some “nonregular” problems.

Section 4 considers treatment assignment, where the action is a treatment protocol that specifies how to allocate a treatment among individuals in some population. This problem has gained renewed interest in economics and other fields in recent years, with a growing number of empirical applications. Section 5 considers other applications, including problems involving nonregular functionals, partially identified models, and experimental design, where the decision-theoretic perspective can provide useful insights.

2 General Setup and Evaluation of Decision Rules

2.1 Setup

Our setup broadly follows the frequentist statistical decision theory approach of Wald (1950). Many texts, such as Ferguson (1967) and Berger (1993), cover statistical decision theory in detail.¹ We will consider a *statistical model* to be a collection of probability measures \mathcal{M} on some measurable space \mathcal{Z} . We will usually parametrize the probability measures $P \in \mathcal{M}$ as P_θ , where the parameter θ lies in a parameter space Θ :

$$\mathcal{M} = \{P_\theta : \theta \in \Theta\}$$

In nonparametric and semiparametric settings, the set of possible probability measures \mathcal{M} is infinite-dimensional (and therefore so is Θ). In this case it is sometimes more convenient to drop the parametric notation and simply identify θ with P .

The decision maker observes a random variable $Z \sim P_\theta$, with support in \mathcal{Z} , and chooses an action based on the observation. The action space \mathcal{A} is a measurable space. A (nonrandomized) statistical decision rule is a measurable mapping $\delta : \mathcal{Z} \rightarrow \mathcal{A}$. A randomized decision rule is a Markov kernel with source \mathcal{Z} and target \mathcal{A} .² In practice, we can view a randomized decision rule as a function $\delta(Z, U)$, where U is a random variable independent of Z .³ The additional randomness induced by U is usually unnecessary from the standpoint of minimizing risk, but we sometimes need the additional generality to characterize procedures used in practice. For example, Section 3.5 discusses some shrinkage estimators that are

¹See also Blackwell and Girshick (1979), Strasser (1985), Le Cam (1986), Le Cam and Yang (2000), Liese and Miescke (2008), and Manski (2019).

²Let \mathcal{S}_Z and \mathcal{S}_A denote σ -algebras for \mathcal{Z} and \mathcal{A} . A Markov kernel is a map $\Delta : \mathcal{S}_A \times \mathcal{Z} \rightarrow [0, 1]$ such that $z \mapsto \Delta(S, z)$ is measurable for every $S \in \mathcal{S}_A$ and $S \mapsto \Delta(S, z)$ is a probability measure for every $z \in \mathcal{Z}$.

³See Pfanzagl (1994), Theorem 1.10.33.

asymptotically equivalent to randomized estimators.

We evaluate decision rules based on a welfare function $W : \Theta \times \mathcal{A} \rightarrow \mathbb{R}$, where we interpret higher values of welfare as more desirable. The expected welfare of a (nonrandomized) rule δ under distribution P_θ is

$$E_\theta [W(\theta, \delta(Z))] = \int W(\theta, \delta(z)) dP_\theta(z).$$

(For a randomized rule $\delta(Z, U)$, we further integrate W over the appropriate distribution of U .) Note that this quantity depends on the value of θ , which is not known.⁴ Therefore, to evaluate decision rules we will need to aggregate the expected welfare over possible values of θ in some way, as we will discuss further below. In some cases it is convenient to work instead with a loss function $L(\theta, a)$, where smaller values of loss are preferred. The *risk* of a decision rule is its expected loss:

$$R(\theta, \delta(Z)) = \int L(\theta, \delta(z)) dP_\theta(z).$$

Example 1 (a) Suppose that after observing $Z \sim P_\theta$, it is desired to produce a point estimate of θ . Then the action space can be taken as $\mathcal{A} = \Theta$, and we can denote a (nonrandomized) point estimator as a function $\hat{\theta} : \mathcal{Z} \rightarrow \Theta$. Suppose θ is scalar and we evaluate the performance of the point estimator by squared error loss $L(\theta, a) = (\theta - a)^2$. Then its risk is

$$R(\theta, \hat{\theta}) = \int L(\theta, \hat{\theta}(z)) dP_\theta(z) = \int (\theta - \hat{\theta}(z))^2 dP_\theta(z).$$

This fits into our framework by setting $W(\theta, a) = -L(\theta, a)$. We will discuss point estimation and other point decision problems, where the action space is a subset of a Euclidean space, in more detail in Section 3.

(b) Alternatively, consider a hypothesis testing problem within this setup. Suppose the researcher needs to choose between the following null and alternative hypotheses:

$$H_0 : \theta \leq 0 \text{ vs. } H_1 : \theta > 0.$$

Then the action space could be specified as $\mathcal{A} = \{H_0, H_1\}$. Rather than specifying a classical hypothesis testing setup where power is maximized subject to controlling size, we consider a loss function that incorporates a trade-off between Type I and Type II errors by assigning a fixed penalty to each type of error:

$$L(\theta, a) = \begin{cases} 1 & \text{if } \theta \leq 0 \text{ and } a = H_1 \\ K & \text{if } \theta > 0 \text{ and } a = H_0 \\ 0 & \text{otherwise} \end{cases}$$

where $K > 0$ captures the trade-off between the two types of errors. The risk associated with a testing

⁴In many cases, only some elements of θ , or some function of θ , is relevant for welfare. For convenience we will sometimes use $\kappa : \Theta \rightarrow \mathcal{K} \subset \mathbb{R}^d$ to denote the decision-relevant quantity. Some specific examples of such functionals κ are considered in Section 5.1.

decision rule $\delta(z)$ is

$$R(\theta, \delta) = \mathbf{1}\{\theta \leq 0\} \Pr_\theta(\delta(Z) = H_1) + K \mathbf{1}\{\theta > 0\} \Pr_\theta(\delta(Z) = H_0).$$

In Section 4, we introduce a closely related loss function in the discussion of treatment assignment decision problems.

Later, to develop large sample results we will work with a sequence of statistical models:

$$\mathcal{M}_n = \{P_\theta^{(n)} : \theta \in \Theta\},$$

with $Z^{(n)} \sim P_\theta^{(n)}$ supported on $\mathcal{Z}^{(n)}$. Here n usually represents the sample size of the data. A leading case is when we observe an i.i.d. sample $Z^{(n)} = Z^n := (Z_1, \dots, Z_n)$ where $Z_i \sim P_\theta$. Then $P_\theta^{(n)} = P_\theta^n$, the n -fold product measure of P_θ . We can then modify the notation for decision rules and welfare accordingly.

2.2 Counterfactuals in the Wald Framework

The Wald framework is quite general and can handle decision problems involving counterfactual outcomes, and more generally structural models with latent variables. To do so, we need to specify the parameter space Θ and associated observation model $Z \sim P_\theta$, the action space \mathcal{A} , and the welfare function $W(\theta, a)$ or loss function $L(\theta, a)$.

Suppose there are possible treatments or policies $t \in \mathcal{T}$ that are interpreted as interventions that could affect an outcome of interest $Y \in \mathcal{Y}$. There is a primitive utility function $u : \mathcal{Y} \rightarrow \mathbb{R}$ that evaluates final outcomes. If Y has a distribution F , then expected utility is $\int u(y) dF(y)$. More generally, there could be some welfare functional $w(F)$ that evaluates different possible distributions of Y .

Since Y is thought to be affected by the action, it is convenient to posit the existence of latent variables $Y(t)$ for $t \in \mathcal{T}$. These *potential outcomes* reflect counterfactual scenarios under different interventions. Each potential outcome has a marginal distribution $F_{Y(t)}$ in the population. If this distribution were known, we could evaluate the intervention t according to its welfare $w(F_{Y(t)})$.

This policy problem becomes a statistical decision problem if, in addition, we have some data Z that are informative about the potential outcome distributions. Structural models in econometrics typically involve latent variables that are invariant to interventions. In causal (potential outcome) models the latent variables are the potential outcomes themselves, but other types of structural models may specify more primitive objects from which potential outcomes may be derived. Suppose that the structural model specifies latent variables

$$S \sim G_\theta, \quad \theta \in \Theta.$$

Here θ is the structural parameter, and we suppose that there is some mapping from θ to the potential

outcome distributions. Let $F_{Y(t),\theta}$ denote the potential outcome distribution under intervention t , given the structural parameter θ .

The latent variables of the structural model are linked to the observed data Z by some observation function $Z = \text{Obs}_\theta(S)$, so

$$Z \sim P_\theta := G_\theta \circ \text{Obs}_\theta^{-1}.$$

We have allowed the observation function $\text{Obs}(\cdot)$ to depend on θ , to accomodate cases like the auction model in Example 3 below. In principle we could also allow the observation process to be noisy, by extending the setup to allow for an additional source of randomness.

In the classical econometric terminology, the distribution P_θ is the *reduced form* distribution of the data under the structural parameter θ . In general, different values of θ may lead to the same observable distribution P_θ , in which case θ is not point identified. (We discuss partial identification in Section 5.2.) The set of possible reduced form distributions $\{P_\theta : \theta \in \Theta\}$ can always be reparametrized as $\{P_\gamma : \gamma \in \Gamma\}$ such that the reduced form parameter γ is point identified. (In the sequel, we will use θ as the general symbol for the parameter of the model, which could be a structural or reduced form parameter depending on the context.)

Finally, we need to specify the action space \mathcal{A} and the welfare function $W(\theta, a)$. In the simplest case, the action space corresponds to the set of treatments, $\mathcal{A} = \mathcal{T}$, and the welfare function (in the expected utility case) is

$$W(\theta, a) = \int u(y) dF_{Y(a),\theta}(y). \quad (1)$$

The action space could be more complicated, for example allowing for randomization over treatments (in which case \mathcal{A} is a set of distributions over \mathcal{T}), or conditional or dynamic treatments (in which case \mathcal{A} is an appropriate set of functions from some set into \mathcal{T}). The welfare function would then need to be defined appropriately relative to \mathcal{A} .

Example 2 (*Simple Randomized Experiment*) Suppose that there are two possible treatments, one of which we will assign to the entire target population. Then we can take $\mathcal{A} = \mathcal{T} = \{0, 1\}$. There is an outcome of interest Y , and we define potential outcomes $Y(0)$ and $Y(1)$. We have a random sample of size n from the target population and randomly assign the treatment to individuals. The latent variables are $S = (S_1, \dots, S_n)$, with

$$S_i = (T_i, Y_i(0), Y_i(1)).$$

Random assignment implies that T_i is independent of $(Y_i(0), Y_i(1))$. Then $S \sim G_\theta$, where θ indexes the joint distribution of all the latent variables. We observe the treatment T_i and the potential outcome corresponding to the treatment assigned to i , $Y_i = Y_i(T_i)$. This defines the observation function $\text{Obs}(\cdot)$, which takes

$$S_i = (T_i, Y_i(0), Y_i(1)) \mapsto (T_i, Y_i(T_i)) = Z_i.$$

Then P_θ is the joint distribution of $Z = (Z_1, \dots, Z_n)$.

For the treatment assignment problem, the relevant components of θ are the marginal distributions $F_{Y(0)}$ and $F_{Y(1)}$. By randomization of treatment, we have

$$Y_i | T_i = 1 \sim F_{Y(1)}, \quad Y_i | T_i = 0 \sim F_{Y(0)}.$$

In this case, the decision relevant components of θ are point-identified from Z . Welfare could be specified as in (1) to correspond to expected utility associated with outcomes. A common example would take welfare to be the expected outcome itself, $u(y) = y$.

The simple binary decision problem in Example 2 can be extended to the case where individuals have observable characteristics $X \in \mathcal{X}$, and we allow the assignment to depend on X . A treatment assignment rule is a mapping $a : \mathcal{X} \rightarrow \{0, 1\}$, interpreted as the decision whether or not to assign some treatment or policy to individuals with characteristics $x \in \mathcal{X}$. Let \mathcal{A} be a set of possible mappings, which could be the set of all possible mappings $2^{\mathcal{X}}$ or some subset of $2^{\mathcal{X}}$. Hence the action space is a subset of a functional space. Then the statistical treatment rule $\delta : \mathcal{Z} \rightarrow \mathcal{A}$ selects a conditional treatment assignment rule based on data Z . We will discuss this problem in detail in Section 4.

The following example illustrates how the framework can handle decision problems involving a structural econometric model of individual behavior.

Example 3 (*Empirical Auction Design*) Suppose there are m bidders for a single object. The latent variables of the structural model are the bidders' private valuations for the object $S = (v_1, \dots, v_m)$, drawn jointly from a distribution G_θ :

$$S = (v_1, \dots, v_m) \sim G_\theta.$$

For example, these valuations could be drawn i.i.d. from a common distribution, corresponding to the independent private values case, or they could arise from an affiliated private values model as in Li, Perrigne, and Vuong (2002).

Given the rules of the auction and an equilibrium notion, bidders will choose bids b_1, \dots, b_m based on their individual valuations. Depending on the application, the observed data could consist of all the bids, or just the winning bid, or some other function of the bid vector. In any case, there will be an observation equation

$$Z = \text{Obs}_\theta(S),$$

where $\text{Obs}_\theta(\cdot)$ captures the mapping from the valuation draws to the observed data. Since the agents' equilibrium bidding functions could depend on the value of θ , which is unknown to the econometrician but typically assumed to be common knowledge among the bidders, the observation function may depend on θ .

Suppose we wish to use the past data to redesign the auction for future instances. For example, we could impose a reserve price, or choose among different auction formats. Let t represent an auction design, and

let \mathcal{T} be the set of possible auction designs under consideration. Let $Y(t)$ be the revenue generated by an auction of type t , when bidders' valuations are drawn from G_θ . If we measure welfare by the expected revenue of the auction, then the welfare function would be

$$W(\theta, t) = \int y dF_{Y(t), \theta}(y).$$

where the dependence of the revenue distribution on G_θ is captured by the θ subscript in $F_{Y(t), \theta}$. In Example 5 below, we further specialize this example to explore large-sample approximations for empirical auction analysis and design.

2.3 Classes of Statistical Decision Rules

Having set up the basic framework for a statistical decision problem, the decision-maker's problem is to choose a statistical decision rule from the set of possible rules. While this choice will be dictated by the specific problem at hand, there are some general classes of rules that are useful in many problems. We discuss a few of them next.

ML Plug-in Rules: Suppose that $\Theta \subset \mathbb{R}^k$ and the probability measures P_θ admit densities p_θ with respect to some fixed measure on \mathcal{Z} . Let $\hat{\theta}_{ML} = \hat{\theta}_{ML}(Z)$ be the maximum likelihood estimator of θ :

$$\hat{\theta}_{ML}(Z) = \arg \max_{\theta \in \Theta} p_\theta(Z).$$

A maximum likelihood (ML) plug-in rule chooses an action that maximizes welfare taking θ at the estimated value:

$$\delta_{ML}(z) = \arg \max_{a \in \mathcal{A}} W(\hat{\theta}_{ML}(z), a).$$

More generally, for any estimator $\hat{\theta}$ we could define a plug-in rule analogously.

Bayes Rules: Let π be a (prior) probability measure over Θ . (We can also allow π to be a general, not necessarily unitary, measure, provided that the following expressions remain well defined.) Let the *Bayes welfare* for rule δ be

$$\overline{W}(\pi, \delta) = \int E_\theta [W(\theta, \delta(Z))] d\pi(\theta) = \int \left[\int W(\theta, \delta(z)) p_\theta(z) dz \right] d\pi(\theta),$$

where $p_\theta(z)$ is the likelihood function. The Bayes welfare averages the expected welfare, which depends on θ , with respect to the prior measure π . A Bayes rule $\delta_B(z)$ maximizes the Bayes welfare

$$\delta_B(z) = \arg \max_{\delta} \overline{W}(\pi, \delta).$$

When the order of integration can be switched, Bayes welfare can be expressed as

$$\overline{W}(\pi, \delta) = \int \left[\int W(\theta, \delta(z)) p_\theta(z) d\pi(\theta) \right] dz.$$

In this case, the Bayes rule is simply the action that maximizes the inner integral in the Bayes welfare expression for each z :

$$\delta_B(z) = \arg \max_a \int W(\theta, a) p_\theta(z) d\pi(\theta).$$

The Bayesian posterior distribution of θ given z is proportional to the product of the likelihood and prior,

$$d\pi(\theta|z) \propto p_\theta(z) d\pi(\theta),$$

so the Bayes rule maximizes posterior expected welfare:

$$\delta_B(z) = \arg \max_a E[W(\theta, a) | z] := \int W(\theta, a) d\pi(\theta|z).$$

Empirical Welfare Maximization: Suppose that $Z^n = (Z_1, Z_2, \dots, Z_n)$ where Z_i are i.i.d. $P \in \mathcal{M}$. Here we identify θ with P so that the welfare function can be written as $W(P, a)$. A natural nonparametric estimator of P is the empirical distribution \mathbb{P}_n which puts probability $\frac{1}{n}$ on each observed value of Z_i . The empirical welfare maximizer solves:

$$\delta_{EW}(z) = \arg \max_{a \in \mathcal{A}_n} W(\mathbb{P}_n, a),$$

where \mathcal{A}_n is a class of actions which could depend on sample size. This can be viewed as a type of plug-in rule.

An alternative formulation is useful in the case where W has the form of an expectation. If

$$W(P, a) = \int u(z, a) dP(z),$$

for some function $u(\cdot, \cdot)$, then

$$W(\mathbb{P}_n, a) = \frac{1}{n} \sum_{i=1}^n u(Z_i, a),$$

and the empirical welfare maximizer simply maximizes $W(\mathbb{P}_n, a)$ with respect to a .

2.4 Evaluating Statistical Decision Rules

For a given decision rule δ , its expected welfare $E_\theta[W(\theta, \delta(Z))]$ is a function of $\theta \in \Theta$. Typically, we cannot learn θ perfectly from the available data, so it is not possible to construct a rule that achieves the best possible expected welfare uniformly over θ . A rule δ is *admissible* if no other rule does as well for all

possible values of θ and strictly improves upon it for some θ : there does not exist $\tilde{\delta}$ with

$$E_{\theta}[W(\theta, \delta(Z))] \leq E_{\theta}[W(\theta, \tilde{\delta}(Z))] \quad \forall \theta \in \Theta, \text{ with}$$

$$E_{\theta}[W(\theta, \delta(Z))] < E_{\theta}[W(\theta, \tilde{\delta}(Z))] \quad \text{for some } \theta.$$

In practice, many rules may be admissible. To make finer comparisons among rules, one approach is to aggregate their expected welfare over $\theta \in \Theta$ in some way, for example by averaging with respect to some measure on Θ or considering some notion of worst-case performance.

2.4.1 Bayes Welfare

Recall that for a prior distribution π over Θ , the Bayes welfare of a rule δ is

$$\overline{W}(\pi, \delta) = \int E_{\theta}[W(\theta, \delta(Z))] d\pi(\theta),$$

and, given π , the corresponding Bayes decision rule maximizes $\overline{W}(\pi, \delta)$. Thus Bayes decision rules automatically maximize average expected welfare. The Bayes welfare criterion can be motivated axiomatically; see for example Anscombe and Aumann (1963) and Savage (1972).

2.4.2 Maxmin and Minmax Regret

The minimum expected welfare of a decision rule δ is

$$\inf_{\theta \in \Theta} E_{\theta}[W(\theta, \delta)].$$

A maxmin rule maximizes this quantity.⁵ Gilboa and Schmeidler (1989) develop an axiomatic justification of the maxmin expected welfare criterion.

An alternative criterion is minmax regret. Consider the action $a^*(\theta)$ that maximizes welfare under θ :

$$a^*(\theta) \in \arg \max_a W(\theta, a)$$

with corresponding risk $W^*(\theta) = \max_a W(\theta, a)$.

Welfare loss regret is defined as the difference between the welfare of this ideal action and the welfare of the given action: $W^*(\theta) - W(\theta, a)$, and a minmax regret rule minimizes:

$$\sup_{\theta \in \Theta} E_{\theta}[W^*(\theta) - W(\theta, \delta)].$$

⁵More generally, we could consider the worst case Bayes welfare with respect to a set of prior distributions over Θ , leading to what is called the Γ -minmax criterion.

The minmax regret criterion seems to have been first proposed by Savage (1951), in a discussion of the minmax criterion used by Wald (1950). Axioms for minmax-regret were proposed by Milnor (1954); see also Hayashi (2008) and Stoye (2011).

2.4.3 Optimality under Restrictions on Decision Rules

In some applications, we may wish to restrict the class of decision rules to satisfy some constraints. For example, there may be practical or institutional restrictions on the kinds of allocation rules allowed in the treatment assignment problems introduced following Example 2 above and further considered in Section 4. A more familiar example is the classical Neyman-Pearson approach to inference, where one imposes a constraint that the decision rule controls size (and possibly satisfies some other restrictions), and then seeks to find the “best” rule (according to some criterion) within this class. Similarly, in point estimation one could restrict attention to unbiased estimators and then seek a minmax or average risk minimizing rule. One could also consider rules that approximately satisfy some constraint, as in Müller and Wang (2019). Sharp results are typically only available in very simple models, though in certain key cases they can be extended via the local asymptotic approximation approach discussed in the next subsection.

2.5 Bounds and Large Sample Approximations for Decision Rules

A statistical decision rule δ can be viewed as a Markov transition from the data space \mathcal{Z} to the action space \mathcal{A} . Thus it takes each probability distribution $\{P_\theta\}$ over \mathcal{Z} into a probability distribution over \mathcal{A} . In principle, given a rule δ we could work out its distribution under every θ , and then compute its expected welfare $E_\theta[W(\theta, \delta(Z))]$ as a function of θ . However, in practice this may be computationally infeasible, especially if we want to evaluate and compare a large set of possible decision rules.

Nevertheless, it may be possible to obtain useful approximations to the distributions and expected welfare properties of decision rules, in a form that facilitates comparisons of rules. This is especially the case for expected welfare regret, because the regret criterion recenters expected welfare conveniently for bounding and limiting arguments. In some cases it is possible to obtain a simple characterization of rules that are approximately optimal with respect to a specific criterion.

2.5.1 Welfare Bounds via Concentration Inequalities

It is often difficult to explicitly calculate worst-case expected welfare, but in some cases it is possible to obtain finite-sample bounds on expected welfare regret. Recent applications of this approach in economics include Manski (2004) and Kitagawa and Tetenov (2018).

To give a flavor of these arguments, suppose that Z_i are i.i.d. P and welfare has the form $W(P, a) =$

$E_P[u(Z, a)]$ for some function u . The empirical welfare based on a sample of size n is

$$W(\mathbb{P}_n, a) = \frac{1}{n} \sum_{i=1}^n u(Z_i, a).$$

The empirical welfare maximizer is

$$\delta_{ew} = \arg \max_{a \in \mathcal{A}_n} W(\mathbb{P}_n, a),$$

but the true welfare of this rule is $W(P, \delta_{ew})$. Since P is not known, we do not know the true welfare of the rule.

One possibility is to obtain bounds on the regret

$$W(P, a^*) - W(P, \delta_{ew}),$$

where a^* solves $\max_{a \in \mathcal{A}} W(P, a)$. Here \mathcal{A} could be equal to \mathcal{A}_n , the same class of possible rules considered in the construction of the EWM rule, in which case the welfare of δ_{ew} is compared to the welfare of the best possible rule in the same class. Or \mathcal{A} could be the class of all possible rules, so that the welfare of the EWM rule is compared to the best possible welfare among all rules.

Useful finite-sample bounds on the regret can sometimes be obtained from concentration inequalities for empirical processes. The key step is to bound the empirical process

$$W(\mathbb{P}_n, a) - W(P, a) = \frac{1}{n} \sum_{i=1}^n u(Z_i, a) - E_P[u(Z_i, a)],$$

uniformly over a . Under some conditions, one can guarantee that this quantity is small with high probability. See Bousquet, Boucheron, and Lugosi (2004) for a survey of techniques for obtaining these types of bounds.

This approach yields worst-case bounds on regret which typically depend on sample size n . Embedding the fixed-sample decision problem in a sequence where sample size n increases, one then obtains a rate of convergence of minmax regret towards zero. A decision rule is (minmax regret) rate-optimal if its minmax regret expected welfare achieves the best possible rate of convergence.

2.5.2 Large Sample Approximations via Local Asymptotics

Large sample distributional approximations play an important role in statistics and econometrics, especially in constructing inference procedures. Here, we show some ways to use large sample theory to obtain approximate characterizations of the welfare properties of decision rules and comparisons between decision rules. We use the framework developed by Le Cam (1972, 1986). Our notation and key results are adapted from van der Vaart (1991a) and van der Vaart (1998). Other useful treatments include Ibragimov and Hasminskii (1981) and Le Cam and Yang (2000). In this section we focus on the case with

i.i.d. data from a smooth parametric model to illustrate the main concepts, but local asymptotic approximations can also be used in some nonstandard parametric models, settings with dependent data, and settings where the parameter space is infinite-dimensional.

Consider a random sample $Z^n = (Z_1, \dots, Z_n)$, where the Z_i are i.i.d. with distribution P_θ for $\theta \in \Theta \subset \mathbb{R}^k$. So Z^n is distributed P_θ^n . We will take approximations as the sample size n goes to infinity. If the model is point identified, typically there will exist point estimators $\hat{\theta}$ of θ that are consistent in the sense that $\hat{\theta} \xrightarrow{P} \theta$ as $n \rightarrow \infty$. Then a plug-in decision rule based on $\hat{\theta}$ will typically have expected welfare converging to the infeasible optimal welfare that would obtain if θ were known. Of course, this type of crude approximation does not capture the welfare loss due to estimation error in $\hat{\theta}$, nor does it lead to useful comparisons between, say, plug-in rules based on different consistent estimators of θ .

Thus, we wish to obtain finer approximations that capture the role of parameter uncertainty in the performance of a decision rule. Local asymptotic approximations aim to do this by considering sequences of parameter values that are difficult to distinguish from each other even as the sample size grows. Suppose that the parametric model $\{P_\theta\}$ satisfies a standard smoothness condition around a centering value $\theta_0 \in \Theta$:

Assumption 1 (a) *Differentiability in quadratic mean: there exists a function $s : \mathcal{Z} \rightarrow \mathbb{R}^m$, the score function, such that*

$$\int \left[dP_{\theta_0+h}^{1/2}(z) - dP_{\theta_0}^{1/2}(z) - \frac{1}{2} h' \cdot s(z) dP_{\theta_0}^{1/2}(z) \right]^2 = o(\|h\|^2) \quad \text{as } h \rightarrow 0;$$

(b) *The Fisher information matrix $J_0 = E_{\theta_0}[ss']$ is nonsingular.*

Assumption 1 is a sufficient condition for the model to be locally asymptotically normal (LAN) at θ_0 (van der Vaart, 1998). Given this assumption, it will be useful to adopt the usual local parametrization around a point θ_0 ,

$$\theta_{n,h} = \theta_0 + \frac{h}{\sqrt{n}}.$$

Differentiability in quadratic mean ensures that the normalized log-likelihood function of the model is approximately quadratic in a local neighborhood of θ_0 . This turns out to be a key property that leads to the local asymptotic normality property. More generally, the limiting distributions of natural estimators and the limiting properties of the decision problem will depend crucially on the local behavior of the likelihood function of the model. In “non-regular” models the likelihood may have a different limiting form, but this form may be sufficiently tractable to lead to useful risk approximations and comparisons of decision rules.

In general, we fix $\theta_0 \in \Theta$ and consider sequences of local alternatives $\theta_0 + \phi_n h$, where $h \in \mathbb{R}^k$ and $\phi_n \rightarrow 0$ is a normalizing sequence of matrices. In regular parametric models satisfying Assumption 1, the appropriate norming is $\phi_n = \frac{1}{\sqrt{n}} I_k$ where I_k is the k -dimensional identity matrix. Then, we would approximate

the distribution of a root- n consistent estimator $\hat{\theta}$ by taking the distributional limit of

$$\phi_n^{-1}(\hat{\theta} - (\theta_0 + \phi_n h)) = \sqrt{n}(\hat{\theta} - \theta_0) - h$$

under the sequence of measures $P_{\theta_0 + h/\sqrt{n}}^n$. We will denote this convergence by $\theta_0 + h_0/\sqrt{n} \rightsquigarrow$. However, in other applications a different norming sequence may be appropriate, and different subcomponents of θ could have different norming rates.

Consider the sequence of statistical models $\mathcal{E}_n = \{P_{\theta_0 + \phi_n h}^n : h \in \mathbb{R}^k\}$. Their likelihood ratio processes are defined as

$$\Lambda_{n,h_0} = \left(\frac{dP_{\theta_0 + \phi_n h}^n}{dP_{\theta_0 + \phi_n h_0}^n} \right)_{h \in \mathbb{R}^k}. \quad (2)$$

When working with Λ_{n,h_0} , we take its distribution under the local sequence of probability measures $P_{\theta_0 + \phi_n h_0}^n$. The sequence of experiments \mathcal{E}_n is said to converge weakly to the experiment \mathcal{E} , supported on some measurable space $\{\mathcal{X}, \mathcal{B}\}$, with probability measures $\{F_h : h \in \mathbb{R}^k\}$, if the likelihood ratio processes in (2) converge (in the sense of finite-dimensional weak convergence) to

$$\left(\frac{dF_h}{dF_{h_0}} \right).$$

In other words, the likelihood ratio processes of the model are approximated by the likelihood ratio processes of the model $\{F_h : h \in \mathbb{R}^k\}$.

For example, under the regularity conditions in Assumption 1, it can be shown that

$$\log \left[\frac{dP_{\theta_0 + \phi_n h}^n}{dP_{\theta_0}^n} \right] = h' \Delta_n - \frac{1}{2} h' J_0 h + o_p(1),$$

where J_0 is the Fisher information matrix defined in Assumption 1(b), and $\Delta_n \rightsquigarrow N(0, J_0^{-1})$, where the weak convergence is under the measures $P_{\theta_0}^n$. From this we can show that, for each $h_0 \in \mathbb{R}^k$ and every finite subset $I \subset \mathbb{R}^k$, the vectors

$$\left(\frac{dP_{\theta_0 + h/\sqrt{n}}^n}{dP_{\theta_0 + h_0/\sqrt{n}}^n} \right)_{h \in I} \xrightarrow{\theta_0 + h_0/\sqrt{n} \rightsquigarrow} \left(\exp \left[(h - h_0)' \Delta - \frac{1}{2} (h - h_0)' J_0 (h - h_0) \right] \right),$$

where $\Delta \sim N(h_0, J_0^{-1})$. The limit on the right is the log likelihood ratio associated with the experiment of observing a single draw from the shifted normal distribution $N(h, J_0^{-1})$, where J_0 is known and h is the parameter of interest.

The fact that the likelihood ratio process of the model of interest, after suitable normalization, converges to the likelihood ratio of the simple shifted normal model, suggests that the normal model serves as a kind of canonical model, characterizing the limit distributions of feasible decision rules. This intuition is made precise in asymptotic representation theorems, which state that, for any sequence of decision rules

δ_n that have limiting distributions under the local parameter sequence $\theta_0 + h/\sqrt{n}$, there exists a (possibly randomized) decision rule based on a shifted normal $N(h, J_0^{-1})$ whose exact distributions under every h correspond to the limit distributions of δ_n . This enables us to characterize the limiting distributions, and associated welfare and risk properties, of decision rules, and in some cases leads to tractable welfare- and risk-optimality results. We will examine specific applications of asymptotic representation theorems in later sections.

Other limits besides the shifted normal experiment are also possible, and we examine some examples below. If the set of probability distributions under consideration is infinite-dimensional, corresponding to a nonparametric or semiparametric model, local asymptotic approximations are also possible by an appropriate construction of the local parameter space; see Bickel, Klaassen, Ritov, and Wellner (1993), van der Vaart (1991a), and Section 4 below. (See also Müller (2011) for an alternative approach to local asymptotic efficiency without parametric restrictions.)

2.5.3 Approximating Risk and Expected Welfare

Recall that we evaluate decision rules by their expected welfare functions, or their risk functions. The expected welfare of a rule δ_n is

$$E_{\theta_0 + \phi_n h} [W(\theta_0 + \phi_n h, \delta_n)],$$

where the expectation is with respect to the distribution of δ_n under $\theta_0 + \phi_n h$. For a given δ_n , we can view this as a function of the local parameter h .

Taking $n \rightarrow \infty$, we could take the limit of the above expression and compare different decision rules by their limiting expected welfare or risk functions. However, not all choices for the utility or loss functions will lead to useful limits. In later sections we will discuss some conditions that ensure that the limiting expected welfare and risk functions are nondegenerate and lead to meaningful comparisons of decision rules; in some cases these conditions may be quite stringent.

A distinct but closely related approach is to work directly with the approximate distribution of the decision rule. Under local parametrizations of the type described above, decision rules δ_n will often have limit distributions after some normalization:

$$r_n(\delta_n - c_n) \rightsquigarrow Q_h,$$

where Q_h are the limiting laws under different values of the local parameter h , and the norming r_n and centering c_n may depend on the underlying model, the nature of the decision problem, and the specific decision rule. Fixing n , we regard this as an approximation to the finite-sample distribution of δ_n :

$$\delta_n \stackrel{a}{\sim} c_n + Q_h / r_n.$$

With this approximation, we can then obtain heuristic approximations to the expected welfare or risk of δ_n , given a choice for the utility or loss function. We will discuss the relationship between these two approximation strategies in some of the applications below.

3 Point Decisions

In this section, we consider decision problems where the action space \mathcal{A} is a subset of a finite-dimensional Euclidean space. This applies to point estimation in parametric and semiparametric models, but also includes other decision problems. For example, data from auctions can be used to learn about bidder preferences, and in turn to choose the reserve price in future auctions to increase expected revenue to the seller. This can be viewed as a point decision problem where the action is the future value of the reserve price. Another class of point decision problems are statistical portfolio choice problems, where the goal is to choose an allocation vector a . Each element of the vector a is the fraction of wealth that is to be allocated to a particular asset. Data are used to learn about the joint probability distribution of asset returns, but the ultimate goal is to choose the future allocation.

For point estimation and related problems, there is a well developed theory for standard loss functions. We review some key aspects of this theory, emphasizing the use of invariance arguments. We then consider asymptotic theory for point decisions. Under standard loss functions, the finite-sample theory carries over nicely, leading to local asymptotic optimality results.⁶ However, we also wish to consider other loss functions which may be less tractable but are natural in economic applications. We may also consider procedures that involve some preliminary model selection or model averaging. For these more complicated problems, it is more difficult to obtain optimality results, even in large samples. But tools are available to simplify the characterization of the welfare properties of specific rules and to compare rules.

3.1 Shift Equivariance in Point Estimation

To introduce some key ideas, we first consider a very simple point estimation problem. Suppose we observe $Z \sim N(\theta, 1)$ (where the variance is known). A point estimator for θ is a decision rule $\delta : \mathbb{R} \rightarrow \mathbb{R}$ with the interpretation that $\delta(z)$ is the estimate of θ when $Z = z$. Consider squared error loss $L(\theta, a) = (a - \theta)^2$. The risk of the estimator δ is

$$R(\theta, \delta) = E_{\theta} [(\delta(Z) - \theta)^2] = \int (\delta(z) - \theta)^2 \phi(z - \theta) dz$$

where $\phi(\cdot)$ is the standard normal density. As we noted above, this fits into our framework with welfare W equal to the negative of loss L . Expected welfare then equals the negative of risk.

⁶Some of the results in this section are drawn from the working paper Hirano and Porter (2003b).

Consider the point estimator $\delta(Z) = Z$. It has a number of attractive properties (which will follow from the results below). It is the maximum likelihood estimator (and hence the ML plug-in rule); it is also the Bayes decision rule under a diffuse prior; and it is the empirical welfare maximizer. This rule minimizes Bayes risk with respect to the diffuse prior. It can also be shown to be minmax (and minmax regret), and its risk is constant over the parameter space.

These basic finite-sample results can be generalized in a number of ways. One useful extension is to settings that preserve the “shift” structure of the simple normal model. Suppose we observe

$$Z = \theta + V,$$

where $\theta \in \mathbb{R}^k$ and V has a known, continuous probability distribution on \mathbb{R}^k , independent of θ . We are interested in estimating θ based on a single observation of Z . Formally:

Parametric Shift Model *Let the sample space be $(\mathbb{R}^k, \mathcal{B})$, where \mathcal{B} is the Borel sigma-algebra. Let $p(\cdot)$ be a fixed probability density with respect to Lebesgue measure on \mathbb{R}^k . The probability measures $\{P_\theta : \theta \in \mathbb{R}^k\}$ specify that Z has density $p(z - \theta)$.*

For the parametric shift model, let the action space \mathcal{A} be the parameter space \mathbb{R}^k , corresponding to point estimation or point forecasting. Suppose that the loss function has the form $L(a - \theta)$. In other words, the loss only depends on the parameter and the action (i.e. the estimate or forecast) through their difference $a - \theta$. We assume L has the following properties:

Assumption 2 *The loss function $L(a - \theta)$ satisfies:*

1. $L(\cdot)$ is continuous and $L(a - \theta) \geq 0$.
2. $L(0) = 0$;
3. The sets $\{w \in \mathbb{R}^k : L(w) \leq \tau\}$ are compact for all $\tau \geq 0$.

This covers a number of popular loss functions besides squared error loss. We could work with absolute error loss $L(a - \theta) = |a - \theta|$, or asymmetric loss functions such as “check function” loss:

$$L(a - \theta) = \begin{cases} c(a - \theta) & \text{if } a - \theta \geq 0 \\ -(1 - c)(a - \theta) & \text{if } a - \theta < 0 \end{cases}$$

for some $c \in (0, 1)$. Another asymmetric loss function, used in forecasting and other applications, is the linex (linear-exponential) loss (Varian (1974); see also Zellner (1986) and Christoffersen and Diebold (1997)):

$$L(a - \theta) = \exp(c(a - \theta)) - c(a - \theta) - 1, \quad c \neq 0. \quad (3)$$

This loss function is approximately linear on one side of zero and approximately exponential on the other side, with the degree of asymmetry controlled by c .

In this extension of the normal mean problem, the statistical model, the action space, and the loss function have a similar “shift” or translation group structure. Formally, we can consider the translation group in \mathbb{R}^k as defining groups of transformations over the parameter and sample spaces, and the loss function is symmetric with respect to these groups. Eaton (1989) provides an extensive discussion of statistical methods for models with a group structure.

It will be useful to consider decision rules that are also symmetric with respect to translations in \mathbb{R}^k . These are the *equivariant* estimators in this setting. A nonrandomized estimator is equivariant to translation if, for all $z, g \in \mathbb{R}^k$,

$$\delta(z + g) = \delta(z) + g.$$

In other words, if the observation z is shifted by g , then the corresponding estimate will also be shifted by g . Letting $g = -z$, it follows that $\delta(z) = \delta(0) + z$, so any such rule must have the form $\delta(Z) = Z + s$, where s is some constant.

We can extend the concept of equivariance to randomized estimators. A randomized estimator is *equivariant in law* if the distribution of $\delta - \theta$ under P_θ is invariant to θ : letting $\mathcal{L}_\theta(\cdot)$ denote the law of a given random variable under P_θ , there is a law \mathcal{L}_0 such that

$$\mathcal{L}_0 = \mathcal{L}_\theta(\delta - \theta), \quad \forall \theta \in \mathbb{R}^k.$$

The following result, known as the Convolution Theorem, provides a simple characterization of equivariant-in-law estimators. (This result follows from Proposition 1 in Le Cam (1986), Chapter 8.3.)

Proposition 1 *Let δ be an equivariant-in-law randomized estimator in a parametric shift model, and let \mathcal{L}_0 be its null law as defined above. Then \mathcal{L}_0 can be written as*

$$\mathcal{L}_0 = \mathcal{L}(V + S),$$

where $V = Z - \theta$ and S is a random variable independent of V .

This result implies that any equivariant-in-law estimator can be written as the simple convolution of the observation and another random variable:

$$\delta = Z + S.$$

Let μ_S denote the distribution of S . Then the risk of an equivariant-in-law estimator can be written as

$$R(\theta, \delta) = E_\theta[L(\delta - \theta)] = \int \int L(z + s - \theta) p(z - \theta) dz d\mu_S(s).$$

Setting $v = z - \theta$, we have

$$R(\theta, \delta) = \int \int L(v + s) p(v) dv d\mu_S(s),$$

which does not depend on θ . Since equivariant-in-law rules have constant risk, they can be ordered

and we can seek to minimize risk across this class. In most cases, this minimum can be attained by a nonrandomized equivariant rule of the form $\delta = Z + s$, for which the risk is

$$R(\theta, \delta) = \int L(v + s)p(v)dv.$$

To streamline the analysis, we make the following assumption:

Assumption 3 *The minimization problem*

$$\min_{s \in \mathbb{R}^k} \int L(v + s)p(v)dv$$

has a unique solution s^ .*

In practice it is usually straightforward to check whether the assumption holds and to calculate s^* . Then we have the following result:

Proposition 2 *Suppose that Assumptions 2 and 3 hold. Then the estimator $\delta = Z + s^*$ is best equivariant.*

Proof: Since s^* minimizes $\int L(v + s)p(v)dv$, it follows that for any law $\mu_S(s)$,

$$\int \int L(v + s)p(v)dv d\mu_S(s) \geq \int \int L(v + s^*)p(v)dv d\mu_S = \int L(v + s^*)p(v)dv.$$

□

Working with equivariant rules simplifies the analysis considerably, but at first glance may appear to be unnecessarily restrictive. It turns out, however, that the best equivariant rule possesses some appealing optimality properties.

First, consider the average risk criterion. Recalling the definition of Bayes welfare in Section 2, define the Bayes risk of a rule δ for a prior measure π on the parameter space as

$$\bar{R}(\pi, \delta) = \int R(\theta, \delta)d\pi(\theta).$$

Consider the case where π is Lebesgue measure on \mathbb{R}^k . Under this “flat” prior, the best equivariant rule minimizes Bayes risk among all possible randomized estimators:

Proposition 3 *Under Assumptions 2 and 3, the best equivariant rule $\delta = Z + s^*$ minimizes $\bar{R}(\pi, \delta)$ for π equal to Lebesgue measure on \mathbb{R}^k .*

Proof: The Bayes risk equals

$$\int R(\theta, \delta) d\theta = \int \int L(\delta(z) - \theta) p(z - \theta) dz d\theta.$$

Rearranging the order of integration, it is enough to choose $\delta_B(z)$ for each z to minimize

$$\int L(\delta(z) - \theta) p(z - \theta) d\theta.$$

Setting $v = z - \theta$, we can write the minimand as

$$\int L(\delta(z) + v - z) p(v) dv,$$

and setting $\alpha(z) = \delta(z) - z$,

$$\int L(v + \alpha(z)) p(v) dv.$$

By Assumption 3, this is minimized by setting $\alpha(z) = s^*$ for all z , implying that $\delta(z) = z + s^*$ minimizes Bayes risk.

□

The best equivariant estimator is also minmax. The following result is a consequence of the generalized Hunt-Stein theorem in Wesler (1959).

Proposition 4 *Suppose that Assumptions 2 and 3 hold. Then $\delta = Z + s^*$ is minmax among all randomized estimators.*

Finally, let us also consider the maximum likelihood estimator. In the shift model, the MLE solves

$$\hat{\theta}(z) = \arg \max_{\theta \in \mathbb{R}^k} p(z - \theta).$$

Suppose that the density $p(v)$ has a unique maximum v^* . Then it is clear that

$$\hat{\theta}(z) = z - v^*.$$

So the MLE is equivariant, and by shifting the MLE we obtain the best equivariant estimator:

$$\delta^*(z) = \hat{\theta}(z) + (s^* + v^*).$$

While we have seen that some natural estimators, such as the MLE and the Bayes estimator under a flat prior, are equivariant, it may also be desirable to consider the behavior of a broader class of estimators. We will return to this issue in discussing model choice and shrinkage procedures below.

3.2 Asymptotics for Point Estimators

The exact results surveyed in the previous subsection are powerful but require that the parametric model, the action space, and the loss function share a special structure. This rarely obtains in realistic applications, but the shift structure, or at least some aspects of it, often hold approximately in more complicated problems, a point first shown by Le Cam in a number of influential papers including Le Cam (1970, 1972). Here we consider some ways to use Le Cam's local asymptotic theory to simplify the analysis of point decision rules.

As in Subsection 2.5.2, consider a random sample $Z^n = (Z_1, \dots, Z_n)$, where the Z_i are i.i.d. with distribution P_θ for $\theta \in \Theta \subset \mathbb{R}^k$. Fix $\theta_0 \in \Theta$ and suppose there is a local parametrization $\theta_0 + \phi_n h$ for $h \in \mathbb{R}^k$ and some norming sequence $\phi_n \rightarrow 0$. The norming sequence is intended to make the parameter $\theta = \theta_0 + \phi_n h$ difficult to distinguish from θ_0 asymptotically, in the sense of convergence of likelihood ratios below. Typically, there will exist estimators $\hat{\theta}$ such that $\phi_n^{-1}(\hat{\theta} - \theta)$ has nondegenerate limiting distributions under local sequences of parameter values.

Suppose the likelihood ratio processes satisfy convergence to a shift experiment: for every finite subset $I \subset \mathbb{R}^k$ and every $h_0 \in \mathbb{R}^k$,

$$\Lambda_{n,h_0} = \left(\frac{dP_{\theta_0 + \phi_n h}^n}{dP_{\theta_0 + \phi_n h_0}^n} \right)_{h \in I} \xrightarrow{\theta_0 + \phi_n h_0} \left(\frac{f_{\theta_0}(X - h)}{f_{\theta_0}(X - h_0)} \right)_{h \in I}, \quad (4)$$

where $f_{\theta_0}(\cdot)$ is a probability density function, and X has density $f_{\theta_0}(x - h_0)$. The limit is the likelihood ratio process of the experiment of observing $X = h + V$, where V has density $f_{\theta_0}(v)$.

Example 4 Recall that if the differentiability in quadratic mean condition, Assumption 1, holds, then the model is locally asymptotically normal (LAN). In this case, $\phi_n = n^{-1/2}$ and f_{θ_0} can be taken to be a multivariate normal density:

$$f_{\theta_0}(\cdot) = dN(\cdot | 0, J_0^{-1}),$$

where J_0 is the Fisher information matrix at θ_0 .

A closely related case occurs when J_0 is random. Such cases are called local asymptotic mixed normal (LAMN), and arise in some time series applications and other settings. (See Davies (1985), Jeganathan (1982), and Jeganathan (1995), among others.) LAMN models are shift experiments conditional on J_0 . We will not pursue this case explicitly below, but many of the results for shift experiments generalize to the case of conditional shift experiments.

Example 5 Structural econometric models with parameter-dependent support are non-regular, but can have interesting and tractable limit distributions. To illustrate this we consider a simplified version of the parametric auction model from Paarsch (1992).

Consider a first-price procurement auction with m bidders under the independent private values paradigm. Bidders are symmetric with cost $c \sim \text{Exp}(\theta)$, where $\text{Exp}(\theta)$ denotes an exponential distribution with mean θ . In the symmetric equilibrium, bidders bid

$$b = c + \frac{\theta}{m-1}$$

Note that the support of the distribution of bids depends on the parameter θ .

We have a random sample of n auctions (each with m bidders). Suppose we observe the winning (lowest) bid \underline{b}_i for each auction. Then for $i = 1, \dots, n$,

$$\underline{b}_i \sim \text{Exp}\left(\frac{\theta}{m}\right) + \frac{\theta}{m-1}.$$

Although the distribution of winning bids is shifted by $\theta/(m-1)$, its shape also depends on the parameter, so this is not a shift experiment. However, the structure of the model simplifies under local asymptotics.

Fix a centering value θ_0 . The appropriate norming turns out to be $\phi_n = n^{-1}$. Then it can be shown that the experiments converge to a shifted exponential experiment with

$$f_{\theta_0}(\cdot) = d\text{Exp}\left(\cdot \mid \frac{m-1}{m}\theta_0\right). \quad (5)$$

See Hirano and Porter (2003a) for further details.

A consequence of the limiting shift experiment is an asymptotic representation theorem for estimators (and other point decisions). The following result is due to van der Vaart (1991a).

Theorem 5 Suppose that the sequence of likelihood ratio processes converges to the likelihood ratio process of a shift experiment as in Equation (4). Let $\hat{\theta}_n$ be a sequence of estimators that have limits under the local parameter sequences $\theta_0 + \phi_n h$:

$$\phi_n^{-1}(\hat{\theta}_n - \theta_0) \overset{\theta_0 + \phi_n h}{\rightsquigarrow} L_h.$$

Then there exists a (possibly randomized) estimator $T(X, U)$ such that

$$T(X, U) \sim L_h$$

when $X = h + V$, where V has density $f_{\theta_0}(v)$, and U has a standard uniform distribution independent of X .

This result states that any estimator possessing limits under the local parameter sequences can be represented asymptotically as some (possibly randomized) estimator in the simple shift model $X = h + V$.

Moreover, under additional regularity conditions, the form of the matching limit experiment estimator

T can be deduced from the nature of the original estimator in a number of important cases. In particular, suppose that $\hat{\theta}_n$ is the maximum likelihood estimator. Then, under conditions for an argmax theorem (see, e.g., Theorem 3.2.2 in van der Vaart and Wellner (1996)), its limiting representation T is the maximum likelihood estimator in the limit experiment:

$$T(X) = \arg \max_h f_{\theta_0}(X - h).$$

Here T does not depend the auxiliary random term U that appears in the asymptotic representation theorem. Since T is shift-equivariant in the limit experiment, $\hat{\theta}_n$ is locally asymptotically shift-equivariant, i.e. regular. We will discuss regular estimators in more detail below.

3.3 Limits for Risk Functions

We want to use large sample approximations to study the approximate risk properties of point estimators. To do this we need to examine the limiting properties of risk functions.

Suppose that the loss function for point estimation depends only on the difference $\hat{\theta} - \theta$. Then the risk of the estimator is the expected loss

$$R(\theta, \hat{\theta}_n) = E_{\theta} [L(\hat{\theta}_n - \theta)].$$

If the point estimator is consistent, i.e. if $\hat{\theta}_n \xrightarrow{p} \theta$ for all θ , then the risk function will typically converge to 0 as $n \rightarrow \infty$. This type of limit is too crude to make useful asymptotic comparisons between consistent point estimators. Thus we will consider rescaling the risk so that its limit is nondegenerate. This is most straightforward when the loss function is homogeneous.

3.3.1 Homogeneous Loss Function

Suppose $L(\cdot)$ is homogenous of degree j . For example, squared error loss $L(w) = w^2$ is homogenous of degree 2 while absolute error loss $L(w) = |w|$ is homogenous of degree 1. Then it is natural to scale up the risk by the factor ϕ_n^{-j} , so that

$$\begin{aligned} \phi_n^{-j} R(\theta, \hat{\theta}_n) &= E \left[\phi_n^{-j} L(\hat{\theta}_n - \theta) \right] \\ &= E \left[L(\phi_n^{-1}(\hat{\theta}_n - \theta)) \right]. \end{aligned}$$

For example, in LAN models, $\phi_n = n^{-1/2}$, and $n^{j/2} R(\theta, \hat{\theta}_n) = E [L(\sqrt{n}(\hat{\theta}_n - \theta))]$. For estimators such that $\sqrt{n}(\hat{\theta}_n - \theta)$ has a nondegenerate limiting distribution, the limit of the rescaled risk expression will typically be nonzero, facilitating asymptotic risk comparisons.

In Theorem 5, the limiting behavior of point estimators is considered under local parametrizations $\theta_0 +$

$\phi_n h$. This suggests that we should consider the limiting risk as a function of the local parameter h . Suppose the estimator satisfies the conditions of Theorem 5. Then, under mild additional conditions,

$$\begin{aligned}\phi_n^{-j} R(\theta_0 + \phi_n h, \hat{\theta}_n) &= E_{\theta_0 + \phi_n h} [L(\phi_n^{-1} [\hat{\theta}_n - (\theta_0 + \phi_n h)])] \\ &= E_{\theta_0 + \phi_n h} [L(\phi_n^{-1} (\hat{\theta}_n - \theta_0) - h)] \\ &\rightarrow E_h [L(T(X, U) - h)] \\ &=: R_\infty(h, T).\end{aligned}$$

In this form, limiting risk captures finer variation in estimator behavior than limits taken pointwise in θ . Local asymptotic risk often approximates the finite sample risk of the estimator well. This finite-sample behavior is also important for shrinkage-type estimators; see Section 3.5.

Having transformed the estimation problem to the local parameter space, and having normalized the loss by a factor based on the localization and the degree of homogeneity of the loss function, we can now compare estimators by the normalized limiting risks. This amounts to considering estimators in the limiting version of the problem, using the same loss function as in the original problem.

Thus we can connect this limiting problem to the finite sample theory in Section 3.1. Theorem 5 gives an asymptotic characterization of all estimators that possess limit distributions under the local parameter sequences $\theta_0 + \phi_n h$. An important subclass of these estimators are the regular estimators. An estimator $\hat{\theta}_n$ is *regular* (at θ_0) if the limit distributions of $\phi_n^{-1} [\hat{\theta}_n - (\theta_0 + \phi_n h)]$ are the same for all h . Equivalently, regular estimators satisfy

$$\phi_n^{-1} (\hat{\theta}_n - \theta_0) \overset{\theta_0 + \phi_n h}{\rightsquigarrow} L_0 + h,$$

where L_0 is a fixed law that does not depend on h . By Theorem 5, for each such regular estimator $\hat{\theta}_n$, there exists an estimator $T(X, U)$ with the distribution $L_0 + h$ in the limit experiment. Such estimators are equivariant in law, and in Section 3.1 we saw that equivariant-in-law estimators are easy to analyze. In particular, if the loss function satisfies Assumption 2, the risk of $T(X, U)$ is constant in h , making risk calculations and comparisons relatively straightforward. Under some conditions, there is a unique best equivariant estimator, and this estimator is also minmax among all estimators. The risk of the best equivariant estimator in the limiting version of the problem can then serve as a local asymptotic risk bound for the original problem.

The limiting shift experiment given in (4) involves a random variable X with density $f_{\theta_0}(x - h)$. Following the argument in Section 3.1, if $f_{\theta_0}(\cdot)$ has a unique maximum, then the maximum likelihood estimator $T_{ML} = T_{ML}(X)$ can be written as

$$T_{ML} = X - v_{\theta_0}^*,$$

where

$$v_{\theta_0}^* = \arg \max_v f_{\theta_0}(v).$$

Let $s_{\theta_0}^*$ be the minimizer of $\int L(v + s) f_{\theta_0}(v) dv$. Then the shifted MLE $T_{ML} + v_{\theta_0}^* + s_{\theta_0}^*$ is the best equivariant

estimator in the limit experiment.

We can then seek to construct an estimator sequence in the original problem that matches the optimally shifted MLE in the limit. As noted above, under regularity conditions, if $\hat{\theta}_{ML}$ is the MLE in the original problem, then $\phi_n^{-1}(\hat{\theta}_{ML} - \theta_0) \overset{\theta_0 + \phi_n h}{\rightsquigarrow} T_{ML}(X)$. Hence, to achieve asymptotic optimality, we can locally shift $\hat{\theta}_{ML}$ to match the shifted MLE, $T_{ML} + v_{\theta_0}^* + s_{\theta_0}^*$, in the limit experiment. Note that $v_{\theta_0}^*$ and $s_{\theta_0}^*$ depend on f_{θ_0} which is generally unknown. Under further conditions, we can replace these terms by the estimates $v_{\hat{\theta}_{ML}}^*$ and $s_{\hat{\theta}_{ML}}^*$, yielding the following estimator:

$$\hat{\theta}_{BR} = \hat{\theta}_{ML} + \phi_n(v_{\hat{\theta}_{ML}}^* + s_{\hat{\theta}_{ML}}^*).$$

This estimator is best regular, in the sense that its limiting risk is minimal among regular estimators, and it is locally asymptotically minmax in the sense that its worst case limiting risk (suitably defined, see for example van der Vaart (1998)) is minimal among all estimators.

It is important to note that optimality in the local asymptotic minmax sense is conceptually different from finite sample minmaxity, and neither condition implies the other. In the asymptotic analysis, we are considering the normalized risk of the estimator for values of θ in a shrinking neighborhood of θ_0 , rather than over the original parameter space Θ .

Example 4 (LAN example, continued) *In the LAN case these considerations lead to familiar asymptotic optimality results. Specifically, the limit experiment consists of observing $X \sim N(h, J_0^{-1})$, and we have $f_{\theta_0}(\cdot)$ equal to the multivariate normal density with mean zero and variance equal to J_0^{-1} . The density $f_{\theta_0}(\cdot)$ is maximized at $v_{\theta_0}^* = \underline{0}$, so the limiting distribution of the MLE is the same as the distribution of X . For squared error loss, and other symmetric loss functions satisfying Assumption 2, the optimal shift is $s_{\theta_0}^* = 0$. Thus X is best equivariant and minmax under squared error loss (and other symmetric loss functions), and it follows that the MLE is asymptotically best equivariant and minmax.*

Example 5 (Auction example, continued) *Given data on winning bids \underline{b}_i from independent auctions with m bidders, the maximum likelihood estimator is*

$$\hat{\theta}_{ML} = (m-1) \left[\min_i \underline{b}_i \right]. \quad (6)$$

In this case, it is straightforward to obtain the exact distribution of the MLE, which is

$$\hat{\theta}_{ML} \overset{\theta}{\sim} \text{Exp}\left(\frac{(m-1)}{mn} \theta\right) + \theta.$$

As noted above, the limit experiment is a shifted exponential distribution where f_{θ_0} is the density of an exponential distribution, as given in (5). The exponential density is maximized at zero for all θ_0 , so $v_{\theta_0}^ = 0$.*

It follows that the MLE has

$$n(\hat{\theta}_{ML} - \theta_0) \overset{\theta_0 + h/n}{\rightsquigarrow} X,$$

where X has density $f_{\theta_0}(x - h)$. The limiting distribution could alternatively be seen directly from the exact distribution in (6). As discussed in this section, the limiting distribution is the same as the distribution of the MLE in the shifted exponential experiment. However, in that experiment, the optimal estimator under a location-equivariant loss L is given by $X + s_{\theta_0}^*$, where $s_{\theta_0}^*$ is defined above and depends on both L and f_{θ_0} . For example, under squared error loss the optimal shift can be calculated to be

$$s_{\theta_0}^* = -\frac{\theta_0(m-1)}{m},$$

while under absolute error loss the optimal shift is

$$s_{\theta_0}^* = -\frac{\theta_0(m-1)}{m} \log 2.$$

Putting together these pieces, $n^{-1}(v_{\theta_0}^* + s_{\theta_0}^*)$ can be viewed as a kind of bias-correction for the MLE, where “bias” is defined relative to the loss function L . In practice, we can replace $s_{\theta_0}^*$ by a suitable estimator by replacing θ_0 with $\hat{\theta}_{ML}$ in the expression above. Noting that $v_{\theta_0}^* = 0$, our locally asymptotically optimal estimator is

$$\tilde{\theta}_n = \hat{\theta}_{ML} + \frac{s_{\hat{\theta}_{ML}}^*}{n}.$$

This estimator will have limiting risk function $R_{\infty}(h, X + s_{\theta_0}^*)$ equal to the limiting risk of the best equivariant estimator in the shifted exponential limit experiment.

3.3.2 Nonhomogeneous Loss

The theory outlined in the previous subsection is the standard approach to showing local asymptotic optimality of point estimators in settings with limiting shift form, including the LAN case. The homogeneity (and location-equivariant form) of the loss function allows us to work with the same loss function in the renormalized limiting version of the estimation problem.

Next, we consider the problem of choosing an estimator when the loss function is not homogeneous. An example of a nonhomogeneous loss function is linex loss $L(a - \theta) = \exp(c(a - \theta)) - c(a - \theta) - 1$, as defined in (3). Depending on the value of c , linex loss can be strongly asymmetric, but it is smooth and approximately quadratic at the origin. As a result, the $j = 2$ norming is appropriate, because, for $(\hat{\theta}_n - \theta)$ small,

$$\phi_n^{-2} R(\theta, \hat{\theta}_n) \approx \frac{c^2}{2} E \left[(\phi_n^{-1} (\hat{\theta}_n - \theta))^2 \right]. \quad (7)$$

Linex loss is approximately quadratic near zero, so if one applies the appropriate correction under squared error loss, the estimator will have asymptotically optimal limiting risk. For example in LAN models, the optimal correction is zero, so that no correction is needed for the MLE to be optimal with respect to local

asymptotic risk. This is somewhat unsatisfying, because the asymmetry (at the original global scale) of the linex loss function no longer plays a role in the analysis.

We may wish to apply a small sample correction that reflects the global properties of the nonhomogeneous, possibly asymmetric loss function, but also achieves the usual first-order, local asymptotic optimality for the resulting estimator. We outline one possible approach. For concreteness, we base the corrected estimator on the MLE, but the approach could be applied to other estimators that are asymptotically equivalent to the MLE up to a local shift, such as the Bayesian posterior mean.

Suppose we have convergence to a shift experiment as in Theorem 5. Further, assume that the maximum likelihood estimator satisfies $\phi_n^{-1}(\hat{\theta}_{ML} - \theta_0) \overset{\theta_0 + \phi_n h}{\rightsquigarrow} T_{ML}(X)$, where $T_{ML}(X) = X - v_{\theta_0}^*$ as defined in Section 3.3.1. Since $X \overset{h}{\rightsquigarrow} V + h$ as in Section 3.2, we use the approximation $\hat{\theta}_{ML} - \theta_0 \underset{h}{\approx} \phi_n V + \phi_n(h - v_{\theta_0}^*)$.

Now consider a shifted version of MLE: $\hat{\theta}_{ML} + \phi_n s$. We want to use the approximation to the distribution of the MLE to obtain an approximately optimal choice of s . The risk of $\hat{\theta}_{ML} + \phi_n s$ is:

$$\begin{aligned} E_{\theta_0 + \phi_n h} [L(\hat{\theta}_{ML} + \phi_n s - (\theta_0 + \phi_n h))] &= E_{\theta_0 + \phi_n h} [L(\hat{\theta}_{ML} - \theta_0 - \phi_n h + \phi_n s)] \\ &\approx E_V [L(\phi_n V + \phi_n(h - v_{\theta_0}^*) - \phi_n h + \phi_n s)] \\ &= E_V [L(\phi_n V - \phi_n v_{\theta_0}^* + \phi_n s)] \\ &= \int [L(\phi_n v - \phi_n v_{\theta_0}^* + \phi_n s)] f_{\theta_0}(v) dv. \end{aligned}$$

Choose $s_{\theta_0}^*$ to minimize the above expression, which approximates the risk of locally shifted MLE. Given the shift structure of the limit experiment and the regularity of the MLE, it follows that the approximate risk expression does not depend on the local parameter h . Hence, the approximately optimal local shift $s_{\theta_0}^*$ also does not depend on h . As above, one could implement this correction as: $\hat{\theta}_{ML} + \phi_n s_{\hat{\theta}_{ML}}^*$. This estimator would be locally asymptotically optimal under loss L but would include a correction to capture the global curvature of the loss function.

Example 4 (LAN example, continued) *In the LAN case,*

$$\sqrt{n} \left(\hat{\theta}_{ML} - \left(\theta_0 + \frac{h}{\sqrt{n}} \right) \right) \overset{\theta_0 + h/\sqrt{n}}{\rightsquigarrow} V,$$

where $V \sim N(0, J_0^{-1})$. Consider linex loss L as introduced in (3). We seek the finite sample correction factor described above. For simplicity, take θ_0 to be scalar. To obtain $s_{\theta_0}^*$, solve

$$\min_s E \left[L \left(\frac{1}{\sqrt{n}} (V + s) \right) \right]$$

which yields

$$s_{\theta_0}^* = -\frac{1}{2\sqrt{n}} cJ_0^{-1}.$$

We can then use the estimator

$$\tilde{\theta} = \hat{\theta}_{ML} - \frac{1}{2n} c\hat{J}^{-1},$$

which reflects the asymmetry in linex loss. The correction to MLE is of order n^{-1} , so $\tilde{\theta}$ is first order asymptotically equivalent to MLE, and hence asymptotically efficient.

Example 5 (Auction example, continued) The maximum likelihood estimator in the first price auction is regular, with

$$n\left(\hat{\theta}_{ML} - \left(\theta_0 + \frac{h}{n}\right)\right) \overset{\theta_0 + h/n}{\rightsquigarrow} V,$$

where $V \sim \text{Exp}\left(\frac{m-1}{m}\theta_0\right)$. Again we consider linex loss L as introduced in (3) and want to find the corresponding finite sample correction. To obtain $s_{\theta_0}^*$, solve

$$\min_s E\left[L\left(\frac{1}{n}(V + s)\right)\right]$$

which yields

$$s_{\theta_0}^* = \frac{n}{c} \ln\left(1 - \frac{c(m-1)\theta_0}{nm}\right).$$

Note that the approximate linex risk exists if $\frac{m}{(m-1)\theta_0} > \frac{c}{n}$, which holds for large enough n . The resulting shifted MLE is $\tilde{\theta} = \hat{\theta}_{ML} + \frac{1}{c} \ln\left(1 - \frac{c(m-1)\hat{\theta}_{ML}}{nm}\right)$. For finite n , this differs from the bias correction derived above under squared error loss, but since

$$\frac{n}{c} \ln\left(1 - \frac{c(m-1)\theta_0}{nm}\right) \approx -\frac{(m-1)\theta_0}{m}$$

for large n , the small sample correction for linex loss is asymptotically equivalent to the bias corrected estimator under squared error loss.

3.4 Point Decisions for General Loss Functions

So far in this section we have focused on estimation problems, where the loss function depends only on the difference $a - \theta$. This was convenient because it was compatible with the shift or translation-equivariant structure that often characterizes limit experiments under local asymptotics. We now discuss more general decision problems where the action space is still a finite-dimensional Euclidean space but the loss (or welfare) function may not take a translation-invariant form. This situation occurs, for example, in the choice of reserve auction price in Example 2, and the choice of portfolio weights discussed below in Example 6 below.

For notational simplicity we consider the case with a scalar action space, and emphasize heuristic derivations. Let $\mathcal{A} = \mathbb{R}$ and let the loss function be $L(\theta, a)$. Here we do not restrict the loss function to have the shift form $L(a - \theta)$. It will be useful to compare action and decisions to the ideal (but typically infeasible) “action rule” $a^*(\theta) = \operatorname{argmin}_a L(\theta, a)$. Assume that $a^*(\theta)$ is unique for all θ and is smooth in θ . Also, assume that loss is normalized so that $L(\theta, a^*(\theta)) = 0$ for all θ , as would be true for regret loss. And finally assume that $0 = \frac{\partial}{\partial a} L(\theta, a^*(\theta)) [a - a^*(\theta)]$ for $a \in \mathcal{A}$, which would typically follow by $0 = \frac{\partial}{\partial a} L(\theta, a^*(\theta))$, as would be the case in an interior solution to the minimization problem defining a^* . The additional flexibility of this last assumption can be helpful in cases of constrained action spaces, see Example 6 below.

We consider two related, but distinct ways to adjust plug-in rules, which generalize the discussion in the previous subsections. Our first method is based on limiting versions of all the components of the decision problem under local parameter sequences. In particular, it involves a quadratic local approximation to the loss function.

Method 1

Consider a local parameter sequence $\theta_n(h) := \theta_0 + \phi_n h$. Under suitable conditions, we can expand the loss function as follows:

$$\begin{aligned} L(\theta_n(h), a) &= L(\theta_n(h), a^*(\theta_n(h))) + \frac{\partial}{\partial a} L(\theta_n(h), a^*(\theta_n(h))) [a - a^*(\theta_n(h))] \\ &\quad + \frac{1}{2} \frac{\partial^2}{\partial a^2} L(\theta_n(h), a^*(\theta_n(h))) [a - a^*(\theta_n(h))]^2 + o(|a - a^*(\theta_n(h))|^2) \\ &= \frac{1}{2} \frac{\partial^2}{\partial a^2} L(\theta_n(h), a^*(\theta_n(h))) [a - a^*(\theta_n(h))]^2 + o(|a - a^*(\theta_n(h))|^2) \end{aligned}$$

because the first two terms are zero by the assumptions on loss above.

Suppose that we have an estimator $\hat{\theta}$ and consider the plug-in rule $a^*(\hat{\theta})$. Suppose that the plug-in rule is regular in the sense that, for all h , there is a rate $\tilde{\phi}_n$ and a random variable D_{θ_0} such that

$$\tilde{\phi}_n^{-1} (a^*(\hat{\theta}) - a^*(\theta_n(h))) \overset{\theta_0 + \phi_n h}{\rightsquigarrow} D_{\theta_0}.$$

If the estimator $\hat{\theta}$ is regular and a^* is sufficiently smooth, this condition will typically hold. We consider shifted versions of the plug-in rule of the form

$$a^*(\hat{\theta}) + \tilde{\phi}_n s,$$

and seek to find the optimal value of s . The scaled local asymptotic risk of the shifted rule is

$$\begin{aligned}\tilde{\phi}^{-2}R(\theta_n(h), a^*(\hat{\theta}) + \tilde{\phi}_n s) &\approx \frac{1}{2} \frac{\partial^2}{\partial a^2} L(\theta_n(h), a^*(\theta_n(h))) E_{\theta_n(h)} \left[[\tilde{\phi}_n^{-1} (a^*(\hat{\theta}) - a^*(\theta_n(h))) + s]^2 \right] \\ &\approx \frac{1}{2} \frac{\partial^2}{\partial a^2} L(\theta_0, a^*(\theta_0)) E_{D_{\theta_0}} [(D_{\theta_0} + s)^2],\end{aligned}\quad (8)$$

where the last line exploits the regularity of the plug-in rule. This expression is minimized by setting $s = s_{\theta_0}^*$, where

$$s_{\theta_0}^* = -E[D_{\theta_0}]. \quad (9)$$

This suggests that under suitable conditions that ensure the uniformity of the approximation in (8), the rule

$$a^*(\hat{\theta}) + \tilde{\phi}_n s_{\hat{\theta}}^* \quad (10)$$

will minimize the local asymptotic risk among shifted plug-in rules. In essence, we can take a quadratic approximation to the loss function $L(\theta, a)$ and use the regular asymptotic distribution of $a^*(\hat{\theta})$. The resulting rule is first-order locally asymptotically optimal (within the class of shifted plug-in rules) from both a minimax and Bayes perspective.

On the other hand, similar to Section 3.3.2, this type of local asymptotic approximation leads us to disregard the global structure of the loss function. Thus it may be desirable to consider other procedures for constructing adjustments to plug-in rules. Our second method does so by working with the original loss function, but working in the original rather than the local parameter space.

Method 2

Again, given an estimator $\hat{\theta}$ and plug-in rule $a^*(\hat{\theta})$, we seek a shifted plug-in rule $a^*(\hat{\theta}) + t$. Setting $h = 0$ in the local asymptotic parameter sequence, we use the limiting distribution approximation $a^*(\hat{\theta}) \approx a^*(\theta_0) + \tilde{\phi}_n D_{\theta_0}$ to approximate risk $R(\theta_0, a^*(\hat{\theta}) + t) \approx R(\theta_0, a^*(\theta_0) + \tilde{\phi}_n D_{\theta_0} + t)$, and solve the following minimization problem:

$$t_{\theta_0}^* = \arg \min_t E_{\theta_0} [L(\theta_0, a^*(\theta_0) + \tilde{\phi}_n D_{\theta_0} + t)].$$

The solution $t_{\theta_0}^*$ is derived from a risk approximation under the original loss function. Then we can plug in $\hat{\theta}$ to obtain the feasible shifted rule

$$a^*(\hat{\theta}) + t_{\hat{\theta}}^*. \quad (11)$$

Under regularity conditions, the rule (11) will be first-order equivalent to the locally asymptotically optimal shifted rule in (10). In finite samples, it will generally differ because it uses the original loss function instead of its local quadratic approximation. We cannot claim that the rule (11) is optimal in a formal, higher-order sense, because it relies on the first-order approximate distribution of $a^*(\hat{\theta})$. But if this distributional approximation is accurate, (11) may provide finite-sample gains by better accounting for the

shape of the original loss function.

Example 6 *Portfolio allocation based on historical returns data can also be viewed as a point decision problem. Suppose X is an m -vector of risky assets with mean $\mu = E(X)$ and variance $\Sigma = \text{Var}(X)$. For allocation $a \in \mathbb{R}^m$ s.t. $\sum_{j=1}^m a_j = 1$, a mean-variance objective function is given by*

$$W(\theta, a) = a' \mu - c a' \Sigma a,$$

for some $c \geq 0$. The parameters of the model are $\theta = (\mu, \Sigma)$.

The ideal action rule can be solved straightforwardly:

$$\begin{aligned} a^*(\theta) &= \arg \max_{a: \sum a_j = 1} W(\theta, a) \\ &= \frac{1}{2c} \Sigma^{-1} \mu - \frac{1}{2c} \Sigma^{-1} \mathbf{1} (\mathbf{1}' \Sigma^{-1} \mathbf{1})^{-1} \mathbf{1}' \Sigma^{-1} \mu + \Sigma^{-1} \mathbf{1} (\mathbf{1}' \Sigma^{-1} \mathbf{1})^{-1} \end{aligned}$$

where $\mathbf{1}$ denotes a vector of ones. Given an estimator $\hat{\theta}$, such as the MLE or a shrinkage type estimator, the corresponding plug-in rule is $a^*(\hat{\theta})$. Plug-in rules are widely used, but it is well known that they may work poorly due to estimation error (estimation risk) in $\hat{\theta}$.

There is a large literature on Bayesian portfolio choice, including Zellner and Chetty (1965), Bawa, Brown, and Klein (1979), Kandel and Stambaugh (1996), Barberis (2000), and Avramov and Zhou (2010). Mori (2001, 2004) and Kan and Zhou (2007) work with regret loss

$$L^R(\theta, \delta) = W(\theta, a^*(\theta)) - W(\theta, \delta) \quad (12)$$

and analyze the ex ante expected welfare properties of different rules. Mori (2004) shows that MLE and shrinkage plug-in rules are inadmissible, and considers certain adjustments to plug-in rules.

We consider shifted plug-in rules of the form $a^*(\hat{\theta}) + t$. To satisfy the constraint that the allocation vector adds up to one, we will require that the shift t satisfies $\mathbf{1}' t = 0$. (Note that the action rule $a^*(\cdot)$ satisfies $a^*(\theta)' \mathbf{1} = 1$ for any θ .) In this section we have considered two different ways to form the shift. Method 1 takes a quadratic approximation to the loss function, leading to the shift given in (10). Method 2, leading to the shifted rule given in (11), uses the original loss function. In this case, the two approaches are identical because mean-variance utility is already quadratic. Specifically, we can solve for the shift as follows:

$$\min_{t: \mathbf{1}' t = 0} E_{\theta}[L^R(\theta, a^*(\hat{\theta}) + t)].$$

Using either approach developed above, some calculations lead to the following form of the shift:

$$t_{\theta}^* = a^*(\theta) - E_{\theta}[a^*(\hat{\theta})],$$

resulting in the feasible rule

$$a^*(\hat{\theta}) + t_{\hat{\theta}}^* = 2a^*(\hat{\theta}) - E_{\theta}[a^*(\hat{\theta})] \parallel_{\theta=\hat{\theta}}.$$

Here, the optimal shift can be seen to be a bias correction to the plug-in rule. Given that $a^*(\cdot)$ is in the constrained action space, it follows automatically that the shift $t_{\hat{\theta}}^*$ satisfies $\mathbf{1}' t_{\hat{\theta}}^* = 0$ for all θ , and the new shifted rule satisfies the constraint.

The derivation above did not depend on the specific estimator $\hat{\theta}$ used to form the plug-in rule, nor did it require knowledge of the limiting distribution for the estimator. In practice, the term $E_{\theta}[a^*(\hat{\theta})]$ would usually need to be approximated based on the limiting distribution of $\hat{\theta}$. These shifted plug-in rules differ from those in Mori (2004), who develops proportional adjustments to the plug-in rule based on the MLE and Gaussian data to obtain unbiased rules and further considers admissibility of the derived rule.

In Sections 3.2 and 3.3, we considered parameter estimation and focused on corrections or shifts of point estimators. In this subsection which considers general point decisions, Methods 1 and 2 focus on shifting the plug-in action rule, leading to rules of the form $a^*(\hat{\theta}) + t$. We note here that we can also consider shifting the parameter estimator before plugging it into the optimal action rule. That is, we can seek a rule of the form $a^*(\hat{\theta} + t)$ that minimizes an approximation to its risk. We briefly examine this approach to indicate how the ideas in Methods 1 and 2 can be extended and can take advantage of our earlier results on point estimators.

Method 1'

Suppose $\hat{\theta}$ is regular with $\phi_n^{-1}(\hat{\theta} - \theta_n(h)) \overset{\theta_n(h)}{\rightsquigarrow} V_{\theta_0}$, and we search for a shift s such that $a^*(\hat{\theta} + \phi_n s)$ minimizes local asymptotic risk. If a^* is smooth, we can use the Delta Method to obtain a quadratic approximation to the risk of the shifted rule that is analogous to Equation (8), with $D_{\theta_0} + s$ replaced by $\frac{da^*(\theta_0)}{d\theta}[V_{\theta_0} + s]$. We can then choose the shift s to minimize this approximate risk.

A solution to this alternate minimization problem is obtained by setting

$$s = s_{\theta_0}^* = -E_{\theta_0}[V_{\theta_0}],$$

This amounts to removing the (local asymptotic) bias of $\hat{\theta}$ before passing into the function $a^*(\cdot)$, rather than bias-correcting the action rule directly as in Equation (9). A feasible version of this rule is

$$a^*(\hat{\theta} + \phi_n s_{\hat{\theta}}^*),$$

which will minimize local asymptotic risk within the class of rules considered under suitable conditions.

Method 2'

In Method 2, we used the limiting distribution of the plug-in action rule and the original loss function to obtain a shifted version of the plug-in rule. We can modify this approach by shifting the parameter estimator $\hat{\theta}$ instead of $a^*(\hat{\theta})$, using the limiting distribution V_{θ_0} of the parameter estimator. In particular, using $\theta_0 + \phi_n V_{\theta_0}$ as an approximate distribution of $\hat{\theta}$, the risk of a rule $a^*(\hat{\theta} + t)$ can be approximated by $E_{\theta_0} [L(\theta_0, a^*(\theta_0 + \phi_n V_{\theta_0} + t))]$. Then the infeasible optimal rule of the form $a^*(\hat{\theta} + t)$ has t equal to:

$$t_{\theta_0}^* = \arg \min_t E_{\theta_0} [L(\theta_0, a^*(\theta_0 + \phi_n V_{\theta_0} + t))].$$

A feasible version of this rule is

$$a^*(\hat{\theta} + t_{\hat{\theta}}^*).$$

This modified rule will account for the curvature of the loss function beyond its quadratic approximation, and could be particularly convenient when an approximation to the distribution of the estimator is readily available.

Example 5 (Auction example, continued) *In the auction application, we can consider other point decisions besides estimation. For example, we may wish to use data from past auctions to design future auctions. A simple auction design problem is to select the reserve price (in the case of procurement auctions, the maximum acceptable bid). Our discussion is based on Kim (2010).*

Let $a \in \mathcal{A} = [0, \infty)$ be the reserve price. Suppose the auctioneer faces a cost of c_0 and wants to maximize profits from the auction. With m bidders and exponentially distributed costs for each bidder with density g_θ and c.d.f. G_θ , the expected profit from the auction with a reserve price of a is

$$\begin{aligned} W(\theta, a) &= c_0[1 - G_\theta(a)]^m + m \int_0^a [c g_\theta(c) + G_\theta(c)] (1 - G_\theta(c))^{m-1} dc \\ &= e^{-a(\frac{m}{\theta})} \left[c_0 - a + \theta \left(\frac{m-1}{m} \right) \right] - e^{-a(\frac{m-1}{\theta})} \left[\theta \left(\frac{m-1}{m} \right) \right] + \theta \left(\frac{2m-1}{m(m-1)} \right) \end{aligned} \quad (13)$$

where the last expression employs the exponential distribution of bidder costs assumed above.

Kim (2010, 2013, 2015) develops Bayesian methods for auction design, which are Bayes-welfare optimal by construction. Aryal and Kim (2013) considers Γ -maxmin expected utility (set of priors) auction design in a partially identified setting. Morgenstern and Roughgarden (2015) construct nearly optimal auction design by bounding revenue loss. See also Chawla, Hartline, and Nikipelov (2017) who consider adaptive auction design for revenue maximization.

From (13), the ideal reserve price action rule can be explicitly solved:

$$\begin{aligned} a^*(\theta) &= \arg \max_a W(\theta, a) \\ &= \theta + c_0 - \theta w_0 \left(\frac{\theta + c_0}{\theta} \right) \end{aligned} \quad (14)$$

where $w_0(\cdot)$ denotes the principal branch of Lambert's W function.⁷ Given $\hat{\theta}_{ML}$, a feasible plug-in reserve price rule is then $\hat{a}^* = a^*(\hat{\theta}_{ML})$.

We showed that the estimator $\hat{\theta}_{MLE}$ has a shifted exponential form, and is therefore biased, both in finite sample and asymptotically. This suggests that there may be scope to improve upon the MLE plug-in reserve price. As discussed above, Method 1' describes how previous results on the optimally shifted MLE for quadratic loss can be used directly to improve the plug-in rule. Quadratic loss yields a bias-corrected MLE of the form:

$$\tilde{\theta} = \hat{\theta}_{ML} \left(1 - \frac{m-1}{nm} \right).$$

Then, the optimally shifted reserve price is equivalent to

$$\tilde{a}^* = \arg \max_a W(\tilde{\theta}, a) = \tilde{\theta} + c_0 - \tilde{\theta} w_0 \left(\frac{\tilde{\theta} + c_0}{\tilde{\theta}} \right).$$

Alternatively, we could use Method 2' to find an adjustment to the plug-in rule. We start with the ML plug-in rule $a^*(\hat{\theta}_{ML})$ and shift the estimator to be plugged in, as described by Method 2', to maximize profits. Recall $V \sim \text{Exp}(\frac{m-1}{m}\theta)$, the limiting distribution for the MLE under θ , and let $\tilde{V}_\theta = \frac{V}{n} + \theta$. Then, following Method 2', find t_θ^* that maximizes:

$$\begin{aligned} E_\theta [W(\theta, a^*(\hat{\theta}_{ML} + t))] \\ = E_{\tilde{V}_\theta} \left[\exp \left(\left(-\frac{m}{\theta} \right) \left\{ (\tilde{V}_\theta + t) + c_0 - (\tilde{V}_\theta + t) \cdot w_0 \left(\frac{(\tilde{V}_\theta + t) + c_0}{\tilde{V}_\theta + t} \right) \right\} \right) \left((\tilde{V}_\theta + t) + 2c_0 - (\tilde{V}_\theta + t) \cdot w_0 \left(\frac{(\tilde{V}_\theta + t) + c_0}{\tilde{V}_\theta + t} \right) + \frac{m-1}{m} \theta \right) \right. \\ \left. - \frac{m-1}{m} \theta \exp \left(\left(-\frac{m-1}{\theta} \right) \left\{ (\tilde{V}_\theta + t) + c_0 - (\tilde{V}_\theta + t) \cdot w_0 \left(\frac{(\tilde{V}_\theta + t) + c_0}{\tilde{V}_\theta + t} \right) \right\} \right) + \frac{2m-1}{m(m-1)} \theta \right] \end{aligned}$$

where equations (13) and (14) for profits and the plug-in rule are used to obtain the explicit expression above.

Here we have made use of the limiting distribution, \tilde{V}_θ , which is also the exact distribution of the MLE. So, in this case the expected profit expression is exact, and the (infeasible) shift obtained t_θ^* is the exact finite-sample optimal shift to ML for the plug-in rule. A feasible version of this rule would then be given by $a^*(\hat{\theta}_{ML} + t_{\hat{\theta}_{ML}}^*)$. This highlights an advantage of Method 2 and 2': if we have better approximations or the exact distribution of the underlying estimator or plug-in action rule, we can use those in place of their first-order asymptotic approximations to evaluate risk under the original loss function.

3.5 Variable Selection and Shrinkage

So far we have mainly focused on point decision rules that are translation equivariant or regular (i.e., locally asymptotically translation equivariant). Such procedures can have minmax or average risk opti-

⁷The principal branch of Lambert's W function is the real inverse of $h(w) = we^w$ for positive arguments.

mality properties. However, there are often good reasons to consider estimators and point decision rules that are not equivariant. In this section we review some results for shrinkage estimators, which are used in a variety of guises in empirical work in economics, including recent applications of machine learning regression techniques.

As in Section 3.1, we first consider a simple finite-sample setting. Suppose that Z is a k -dimensional normal vector with

$$Z \sim N(\mu, \sigma^2 I_k),$$

where σ^2 is known. Equivalently, we observe independent random variables Z_j for $j = 1, \dots, k$ with $Z_j \sim N(\mu_j, \sigma^2)$. A natural point estimator for μ is $\hat{\mu} = Z$, with components $\hat{\mu}_j = Z_j$. When the parameter space for μ is unrestricted, this estimator is the maximum likelihood estimator and the Bayesian posterior mean under a flat prior. Under squared error loss, it is best equivariant, minmax, and average risk optimal under a flat prior.

Consider estimators of the form

$$\tilde{\mu}_j = Z_j \cdot \Gamma_j(Z),$$

where $0 \leq \Gamma_j(z) \leq 1$. This estimator shrinks each component of $\hat{\mu} = Z$ towards the origin, by a factor that could depend on the data. (We could also consider estimators that shrink towards some other point in \mathbb{R}^k with minor notational modifications.) For example, if

$$\tilde{\mu}_j = Z_j \mathbf{1}(|Z_j| > c),$$

for some c , the estimator $\tilde{\mu}_j = 0$ when the initial estimate Z_j is sufficiently close to 0. Pre-test estimators (also called hard thresholding estimators) that zero out components for which a standard test does not reject the hypothesis that the component is zero, are of this form. Alternatively, the shrinkage towards zero could be more smooth. For example, the positive-part James-Stein estimator for $k \geq 3$ has

$$\tilde{\mu}_j = Z_j \cdot \left(1 - \frac{k-2}{\|Z\|^2}\right)^+,$$

and shrinks $\hat{\mu}_j$ towards zero by a factor that depends on the sum of squares of the elements of Z . (See James and Stein (1961) and Baranchik (1964).) Other point estimators, such as the Bayesian posterior mean under some choices of the prior, ridge regression (Hoerl and Kennard, 1970), and parametric empirical Bayes methods (Morris, 1983), are also of the smooth shrinkage form. Lasso-type estimators (Tibshirani, 1996) combine variable selection with smooth shrinkage. Bagging (Breiman, 1996) converts hard thresholding estimators into a smooth estimator, by averaging the thresholding estimators over bootstrap replications.

For point estimators with action space $\mathcal{A} = \mathbb{R}^k$, consider the sum-of-squares loss function

$$L(\mu, a) = \sum_{j=1}^k (a_j - \mu_j)^2.$$

The estimator $\hat{\mu} = Z$ has constant risk $R(\mu, \hat{\mu}) = k\sigma^2$, and is minmax. In contrast, the risk functions of shrinkage estimators typically vary with μ . The risk of the positive-part James-Stein estimator, and a number of other shrinkage estimators, can be shown to satisfy

$$R(\mu, \tilde{\mu}) \leq k\sigma^2 \quad \forall \mu,$$

with strict inequality for some values of μ , when $k \geq 3$.⁸ In other words, $\tilde{\mu}$ dominates $\hat{\mu}$, and $\hat{\mu}$ is inadmissible.

Thus it is possible, in principle, to improve upon the conventional maximum likelihood estimator in this setting by using some type of shrinkage estimator. The potential gains are large when μ is high-dimensional and there are natural restrictions on the parameter space. For this reason, shrinkage methods may be useful in nonparametric and high-dimensional regression methods, forecasting methods using high dimensional vector autoregressions as in Doan, Litterman, and Sims (1984), and problems involving many units with limited data on each unit as in Graham and Hirano (2011). The shrinkage may be towards specific points in the parameter space suggested by economic theory, as in Ingram and Whiteman (1994), Del Negro and Schorfheide (2004), and Fessler and Kasy (2019), among others.

On the other hand, it is more difficult to obtain sharp optimality results that hold in wide generality, particularly if the assumption of homoskedasticity is dropped.⁹ In general different shrinkage estimators may have risk functions that cross, because they trade off performance across the parameter space differently. Abadie and Kasy (2019) explore the risk properties of different popular shrinkage estimators under different assumptions on sparsity. It can also be difficult to calculate the exact risk function for some procedures based on variable selection or shrinkage, at least under a realistically wide range of data generating processes.

In the remainder of this section we will focus on obtaining large-sample approximations to the risk functions of shrinkage-type procedures using local asymptotics. We seek approximations that are sufficiently tractable to allow comparisons between procedures, and in some cases suggest ways to improve upon them. Local asymptotic approximations of the type we consider here have been used by a number of authors to explore model selection and shrinkage estimators; see for example Knight and Fu (2000), Bühlmann and Yu (2002), Inoue and Kilian (2006), Claeskens and Hjort (2008), and Hansen (2016a).¹⁰ Some of the following results are based on Hirano and Wright (2017).

To fix ideas, consider the classical regression model, where we have i.i.d. draws from a distribution for (y_i, x_i) , with

$$y_i = \beta' x_i + u_i, \quad E[u_i | x_i] = 0, V[u_i | x_i] = \sigma^2,$$

⁸For results on the risk of James-Stein and related estimators, see James and Stein (1961), Baranchik (1964), Egerton and Laycock (1982), Robert (1988), Hansen (2015), and Hansen (2016b), among others.

⁹Recent work on shrinkage estimation in the heteroskedastic case includes Xie, Kou, and Brown (2012).

¹⁰Other types of approximations, especially those that allow the dimension of the parameter space to increase with sample size, are also possible and can provide guidance on the choice of the shrinkage estimator and the selection of the tuning parameters. See Hansen (2016a) and Abadie and Kasy (2019) for some recent results that employ increasing-dimension asymptotics.

where x_i is $k \times 1$. Assume $E[x_i x_i'] = I_k$ for simplicity. (Many of these restrictions can be relaxed.) Let $\hat{\beta}$ be the usual LS estimator. In this setting, the positive-part James-Stein estimator is

$$\tilde{\beta} = \hat{\beta} \times \max \left\{ 1 - \frac{k-2}{n \hat{\beta}' \hat{V}^{-1} \hat{\beta}}, 0 \right\},$$

where \hat{V} is an estimate of $V(\hat{\beta})$. This shrinks all elements of $\hat{\beta}$ towards zero by an equal factor between 0 and 1.

Post-model selection estimators can also be viewed as shrinkage estimators, which set some elements of β to zero based on the data. One approach popular in forecasting applications is to use a pseudo out-of-sample criterion to select elements of β to set to zero. One version of the out-of-sample (OOS) model selection approach is the following. For each model m which picks a subset of the k regressors, estimate the model recursively starting a fraction $\pi \in (0, 1)$ of the way through the sample. At each recursion calculate a one step ahead point forecast, and compare the forecast to the realized value of y . Calculate the empirical mean squared error of the one step ahead forecasts for observations $y_{[\pi n]}, \dots, y_n$. The model m with the lowest empirical mean squared error is chosen, and re-estimated using the entire sample.

The pseudo out of sample approach is intuitive and popular in practice. But it is somewhat difficult to analyze because of its recursive structure. To analyze it and other procedures that use different types of sample splitting to select the model, such as ν -fold cross-validation, it is helpful to use local asymptotics and a representation result tailored to the problem.

We assume the standard LAN conditions and adopt the following local parametrization:

$$\beta_n = \underline{0} + \frac{h}{\sqrt{n}}.$$

Here we have centered the local sequence at $\beta_0 = \underline{0}$ to capture the situation where it is difficult to tell which elements of β are nonzero. Then it can be shown that the OLS estimator (using all k regressors) has the following limit under $\beta_n = h/\sqrt{n}$:

$$\sqrt{n} \hat{\beta} \rightsquigarrow Y, \quad \text{where } Y \sim N(h, \Omega).$$

Thus the OLS estimator acts like an observation in the limiting $N(h, \Omega)$ experiment, and it is regular.

The James-Stein estimator has the following limit:

$$\sqrt{n} \tilde{\beta} \rightsquigarrow Y \times \max \left\{ 1 - \frac{k-2}{Y'Y}, 0 \right\}.$$

So $\tilde{\beta}$ acts like a James-Stein estimator in the limit experiment. It is not regular because its distribution depends on h , but its distributions under h (and associated risks) can be numerically evaluated.

The estimator based on pseudo out-of-sample model selection is less simple to approximate, because it

cannot be written as a function solely of the unrestricted OLS estimator $\hat{\beta}$. However, it (and many other post-model selection and shrinkage estimators) can be viewed as functions of the partial sums

$$\sum_{i=1}^s x_i x_i', \quad \sum_{i=1}^s x_i y_i, \quad s = 1, \dots, T.$$

If we can obtain useful approximations to the partial sums under the local sequences $\beta_n = h/\sqrt{n}$, we can approximate the distributions of a wide range of estimators. To this end, consider the partial sums of $x_i y_i$. For $r \in [0, 1]$, let $[nr]$ be the largest integer less than nr . Then, under standard regularity conditions, a functional central limit theorem gives

$$n^{-1/2} \sum_{i=1}^{[nr]} x_i y_i \rightsquigarrow r b + \sigma B(r),$$

where $B(r)$ is a standard k -dimensional Brownian motion. By a standard construction of the Brownian bridge, we can re-express the limit on the right as

$$r b + \sigma B(r) \sim r Y + \sigma U(r),$$

where $Y \sim N(h, \Omega)$ and $U(r)$ is a Brownian bridge, independent of Y and h . This highlights the connection between the limit of the partial sum and the shifted normal limit experiment. Procedures such as pseudo out-of-sample model selection can be viewed as randomized estimators in the experiment of observing $Y \sim N(h, \Omega)$, depending on both Y and an additional random component $U(\cdot)$. Hirano and Wright (2017) use this result to numerically calculate and compare asymptotic risk functions for a variety of post-model selection and shrinkage estimators.

4 Treatment Assignment Rules

A natural application of statistical decision theory is to choose among different treatments or programs to assign to an individual. Recent work in economics, medical statistics, and other fields has considered treatment assignment, building on methods for causal inference and program evaluation. A classic example is assignment to a medical treatment based on patient characteristics. In economics, examples have included assignment to a job training program (Black, Smith, Berger, and Noel, 2003; Frölich, 2008), enrollment in government welfare programs (Dehejia, 2005), sentencing by judges (Bushway and Smith, 2007), environmental policy targeting (Assunção, McMillan, Murphy, and Souza-Rodrigues, 2019), and household energy use incentives (Knittel and Stolper, 2019).

4.1 Treatment Assignment as a Decision Problem

Suppose that the individuals to be assigned are drawn from some target population, and have observable background characteristics X on a space \mathcal{X} . An individual can be assigned one of two treatments, labeled 0 and 1, based on their value of X . A treatment assignment rule is a mapping from \mathcal{X} to $\{0, 1\}$. We could consider all possible such treatment rules, or some constrained set of treatment rules. In the notation of Section 2, the set of possible treatment rules is our set of actions: $\mathcal{A} \subset 2^{\mathcal{X}}$, where $2^{\mathcal{X}}$ is the power set of all binary-valued functions with domain \mathcal{X} .¹¹

We wish to select a treatment rule based on data $Z \sim P_\theta$. A *statistical treatment assignment rule* maps data into the choice of treatment rule: $\delta : \mathcal{Z} \rightarrow \mathcal{A}$. Since the object chosen by δ is itself a function of X , we simplify the notation by writing

$$\delta(x; z),$$

with the interpretation that, given a data set $Z = z$, the treatment rule $\delta(\cdot; z)$ maps covariate X to the treatment space $\{0, 1\}$. For notational simplicity we also use $\delta(x)$ to denote particular treatment rules in the action space; so $\delta(x) \in \mathcal{A}$ while $\delta(\cdot; z)$ is a mapping from \mathcal{X} to \mathcal{A} .

This framework is fairly general and encompasses a number of cases considered in the recent literature, though it could be extended to handle other settings such as those with a richer set of treatments. (We discuss some possibilities below.) An important special case occurs when the data Z represent a randomized controlled trial conducted on the same population as the target population for future treatment assignment. In general the data $Z \sim P_\theta$ do not have to come from the same population as the target population for future treatment, nor do they have to come from a randomized experiment, provided that the data are informative about a parameter θ that in turn encompasses all welfare-relevant information about the target population.

To define the welfare and associated criteria for evaluating statistical treatment rules, it is convenient to introduce potential outcomes. In the binary treatment case, let $Y(0)$ and $Y(1)$ be potential outcomes under treatments 0 and 1 respectively. Given treatment $T = 0, 1$, the realized outcome is $Y = TY(1) + (1 - T)Y(0)$. The performance of a given decision rule can then be evaluated based on the distribution of outcomes induced by the treatment assignment. Let the distributions of $Y(0)$ and $Y(1)$ conditional on characteristics be denoted by $F_0(\cdot|x, \theta)$ and $F_1(\cdot|x, \theta)$, where θ is the parameter indexing the distributions of Z . Let $F_X(x)$ denote the marginal distribution of X in the target population.

Suppose the welfare associated with a treatment rule $\delta(x)$ is given by a functional of the potential outcome distributions:

$$W(\theta, \delta(\cdot)) = \int w(F_0(\cdot|x, \theta), F_1(\cdot|x, \theta), \delta(x)) dF_X(x).$$

¹¹We could also allow the range of the treatment assignment rule to be $[0, 1]$, with the interpretation that $a(x) \in [0, 1]$ gives the probability of assigning an individual with $X = x$ to treatment 1.

Then an ideal, but generally infeasible, treatment assignment rule is the following:

$$\delta^*(x) = \begin{cases} 1 & \text{if } w(F_0(\cdot|x, \theta_0), F_1(\cdot|x, \theta_0), 1) \geq w(F_0(\cdot|x, \theta_0), F_1(\cdot|x, \theta_0), 0) \\ 0 & \text{if } w(F_0(\cdot|x, \theta_0), F_1(\cdot|x, \theta_0), 1) < w(F_0(\cdot|x, \theta_0), F_1(\cdot|x, \theta_0), 0). \end{cases} \quad (15)$$

(Here, we assign treatment 1 if the two treatments give equal welfare given X .) A special case is the utilitarian welfare criterion, which defines welfare as the expected utility associated with the outcomes or the expected outcomes themselves. The utilitarian welfare of a treatment rule δ can be expressed as

$$W(\theta, \delta) = \int \int u(y) dF_\delta(y|x, \theta) dF_X(x),$$

where u is outcome utility and $F_\delta(y|x, \theta) = \delta(x)F_1(y|x, \theta) + (1 - \delta(x))F_0(y|x, \theta)$. In this case, if the outcome is viewed as the desired utility measure, $u(y) = y$, then the welfare becomes

$$W(\theta, \delta) = \int \delta(x) E_\theta[Y(1)|X=x] + (1 - \delta(x)) E_\theta[Y(0)|X=x] dF_X(x). \quad (16)$$

An infeasible optimal treatment rule chooses, for each x , the treatment that maximizes the expected outcome:

$$\delta^*(x) = \mathbf{1}\{E_\theta[Y(1)|X=x] \geq E_\theta[Y(0)|X=x]\}. \quad (17)$$

which just specializes (15) to the case considered here.

Example 7 Suppose that the data Z come from a randomized experiment on the same population as the target population for treatment assignment. Thus $Z = \{Y_i, T_i, X_i; i = 1, \dots, n\}$, where $Y_i = T_i Y_i(1) + (1 - T_i) Y_i(0)$, the potential outcomes $Y_i(0)$ and $Y_i(1)$ are distributed $F_0(\cdot|x)$ and $F_1(\cdot|x)$ conditional on $X_i = x$, and T_i is independent of the potential outcomes (unconditionally, or conditional on X_i). In this case we identify θ with the entire joint distribution of Z_i .

If X_i is discrete, and welfare is utilitarian with Y measured in utils, the sample analog to Equation (17) is

$$\hat{\delta}(x; Z) = \mathbf{1}\left\{ \frac{\sum_i T_i \mathbf{1}(X_i = x) Y_i}{\sum_i T_i \mathbf{1}(X_i = x)} \geq \frac{\sum_i (1 - T_i) \mathbf{1}(X_i = x) Y_i}{\sum_i (1 - T_i) \mathbf{1}(X_i = x)} \right\}.$$

Manski (2004) called this the conditional empirical success rule and studied its properties.

Regardless of how we specify the welfare function $W(\theta, \delta)$, we still need to decide on a specific criterion for evaluating statistical treatment rules. As in Section 2.4, we can take the expectation of welfare with respect to the sampling distribution of $\delta(x; Z)$,

$$E_\theta[W(\theta, \delta)]$$

and use either a Bayes or maxmin criterion. In a parametric setting, for example, Dehejia (2005) considers a Bayesian approach to treatment assignment in a labor market program; see also Chamberlain

(2011), and for dynamic treatment rules, Arjas and Saarela (2010) and Zajonc (2012). However, when the underlying parameter space and the set of possible treatment rules is large, Bayes and maxmin approaches may be difficult to implement or sensitive to certain aspects of the problem setup. The literature has therefore considered other criteria, especially minmax regret, for which bounds and approximations are relatively tractable. Define the welfare regret loss as

$$L^R(\theta, \delta) = W^*(\theta) - W(\theta, \delta),$$

where $W^*(\theta) = \arg\max_{\delta} W(\theta, \delta)$ is the welfare obtained by the infeasible optimal treatment rule. The maximum welfare regret is

$$\sup_{\theta \in \Theta} E_{\theta} [L^R(\theta, \delta)].$$

This evaluates δ by its worst-case performance relative to the unknown optimal treatment rule. A minmax regret statistical treatment rule minimizes this criterion. In economics, the minmax regret approach to treatment assignment was pioneered by Manski (2004). Other choices for a loss function based on $W(\theta, \delta)$ are also possible. For example we could attach a fixed, possibly asymmetric penalty to each type of incorrect assignment, by using a loss function of the form

$$L(\theta, \delta) = \begin{cases} (1 - \delta) & \text{if } W(\theta, 1) > W(\theta, 0) \\ \delta K & \text{otherwise,} \end{cases}$$

where $K > 0$. Then we could evaluate statistical treatment rules by their maximum expected loss. This type of criterion, which can lead to decision rules similar to classical hypothesis tests, was proposed and studied by Tetenov (2012), and also studied by Hirano and Porter (2009).

The following simple example illustrates minmax regret analysis in a simple parametric case.

Example 8 Suppose the target population has no observable characteristics so that we can suppress the notational dependence on x . The potential outcome distributions are given by $Y(0) \sim N(\mu_0, \sigma_0^2)$ and $Y(1) \sim N(\mu_1, \sigma_1^2)$. Suppose σ_0^2 and σ_1^2 are known so that unknown parameters are $\theta = (\mu_1, \mu_2)$. Suppose welfare is utilitarian with Y measured in utils: $W(\theta, \delta) = \delta\mu_1 + (1 - \delta)\mu_0$. Then regret loss is

$$L^R(\theta, \delta) = \begin{cases} (1 - \delta)(\mu_1 - \mu_0) & \text{if } \mu_1 > \mu_0 \\ \delta(\mu_0 - \mu_1) & \text{otherwise} \end{cases}$$

Asymmetric “hypothesis testing” loss is

$$L^H(\theta, \delta) = \begin{cases} (1 - \delta) & \text{if } \mu_1 > \mu_0 \\ \delta K & \text{otherwise} \end{cases}$$

where $K > 0$. Let $\Delta = \mu_1 - \mu_0$ be the average treatment effect. We can write the corresponding risks of a

decision rule $\delta(Z)$ as:

$$E_\theta[L^R(\theta, \delta)] = -E_\theta[\delta(Z)]\Delta \mathbf{1}\{\Delta \leq 0\} + (1 - E_\theta[\delta(Z)])\Delta \mathbf{1}\{\Delta > 0\}$$

and

$$E_\theta[L^H(\theta, \delta)] = E_\theta[\delta(Z)]K \mathbf{1}\{\Delta \leq 0\} + (1 - E_\theta[\delta(Z)]) \mathbf{1}\{\Delta > 0\}.$$

Now, suppose we have an estimator $\hat{\Delta} = \hat{\Delta}(Z)$ of the average treatment effect based on the data Z , that satisfies $\hat{\Delta} \sim N(\Delta, \sigma^2)$. For simplicity consider threshold decision rules of the form $\delta_c = \mathbf{1}(\hat{\Delta} \geq c)$. It is straightforward to evaluate the maximum expected loss of such decision rules, and to find the minmax value of c (see Tetenov (2012) and Hirano and Porter (2009)). In the case of welfare regret loss L^R , the optimal value can be shown to be $c^* = 0$, leading to the decision rule $\delta_R(Z) = \mathbf{1}(\hat{\Delta} \geq 0)$. For L^H , the optimal value can be shown to be

$$c^* = \sigma \Phi^{-1}\left(\frac{K}{K+1}\right),$$

where $\Phi(\cdot)$ is the standard normal CDF

4.2 Welfare and Risk Analysis of Treatment Assignment Rules

Treatment assignment problems are naturally related to estimation of conditional average treatment effects. In practice, researchers often wish to avoid restrictive parametric assumptions when evaluating treatments. In our setup, this would lead to θ being high-dimensional. In such cases it can be challenging to evaluate and compare decision rules. Here we discuss some methods for analyzing decision rules. We do not consider the most general possible settings but aim to illustrate some of the techniques that can be used under relatively mild distributional assumptions.

We return to the setup of Example 7. The observed data $Z = \{Y_i, T_i, X_i; i = 1, \dots, n\}$ represent a (conditionally) randomized experiment comparing two treatments on the target population, with associated potential outcomes $Y_i(0)$ and $Y_i(1)$. Suppose $X_i \sim F_X$, $Y_i(0)|X_i = x \sim F_0(\cdot|x)$, and $Y_i(1)|X_i = x \sim F_1(\cdot|x)$, and suppose the binary treatment T_i is Bernoulli, conditionally independent of $Y_i(0)$ and $Y_i(1)$ given X_i , with

$$e(x) := \Pr(T_i = 1|X_i = x).$$

Due to the conditional independence of the treatment,

$$Y_i|T_i = t, X_i = x \sim F_t(\cdot|x),$$

and we can identify $\theta = (F_X, e, F_0, F_1)$ with the joint distribution P of (Y_i, T_i, X_i) . We will therefore treat P as the parameter in the sequel.

Under utilitarian welfare (with Y normalized to utils), the welfare $W(\theta, \delta)$ of a treatment rule δ is given

in (16). Under our assumptions we can also write welfare in terms of the observed-data distribution P :

$$\begin{aligned} W(P, \delta) &= E_P [\delta(X) E_P[Y|T=1, X=X] + (1-\delta(X)) E_P[Y|T=0, X]], \\ &= E_P \left[\delta(X) \frac{YT}{e(X)} + (1-\delta(X)) \frac{Y(1-T)}{1-e(X)} \right], \end{aligned} \quad (18)$$

where the P subscripts in (18) highlight that the expectations are taken with respect to the joint distribution of (Y_i, T_i, X_i) in the experimental data set.

The plug-in decision rule replaces P with an estimate $\hat{P}_n(Z)$, and solves:

$$\hat{\delta}_n(\cdot; Z) = \arg \max_{\delta(\cdot) \in \mathcal{A}} W(\hat{P}_n(Z), \delta(\cdot)).$$

Suppose that the support \mathcal{X} of the covariate X is finite, the action space is unrestricted so that $\mathcal{A} = 2^{\mathcal{X}}$, and we use the empirical distribution estimator $\hat{P}_n = \mathbb{P}_n$. Then

$$W(\mathbb{P}_n, \delta) = \frac{1}{n} \sum_{i=1}^n [\delta(X_i) \hat{\mu}_1(X_i) + (1-\delta(X_i)) \hat{\mu}_0(X_i)],$$

where $\hat{\mu}_0(x)$ is the conditional average of Y_i given $T_i = 0$ and $X_i = x$ and $\hat{\mu}_1(x)$ is the conditional average of Y_i given $T_i = 1$ and $X_i = x$. By straightforward calculations, if every (x, t) cell is nonempty, then we can also write the plug-in welfare as

$$W(\mathbb{P}_n, \delta) = \frac{1}{n} \sum_{i=1}^n \left[\delta(X_i) \frac{Y_i T_i}{\hat{e}(X_i)} + (1-\delta(X_i)) \frac{Y_i (1-T_i)}{1-\hat{e}(X_i)} \right],$$

where $\hat{e}(x)$ is the conditional empirical probability of $T = 1$ given $X = x$. Then the plug-in approach yields Manski's conditional empirical success rule as in Example 7.

Manski (2004) proposed to evaluate statistical treatment rules by their maximum regret risk. Under the further assumption that the outcome Y is bounded, he developed finite-sample bounds on the worst-case regret risk of the conditional empirical success rule when P is otherwise unrestricted, using Hoeffding's Large Deviations Theorem. With the same setup, Stoye (2009) used game-theoretic techniques to solve for a minmax regret statistical treatment rule. Stoye's minmax regret rule differs slightly from the conditional empirical success rule. In the case where \mathcal{X} is singleton (so we can omit X from the notation), and T is randomized with $\Pr(T=1) = \frac{1}{2}$, Stoye's decision rule is

$$\delta^*(Z) = \begin{cases} 1 & \text{if } n_1 (\hat{\mu}_1 - 1/2) - n_0 (\hat{\mu}_0 - 1/2) > 0 \\ 1/2 & \text{if } n_1 (\hat{\mu}_1 - 1/2) - n_0 (\hat{\mu}_0 - 1/2) = 0 \\ 0 & \text{if } n_1 (\hat{\mu}_1 - 1/2) - n_0 (\hat{\mu}_0 - 1/2) < 0 \end{cases}$$

where n_t denotes the number of individuals with $T_i = t$, and $\hat{\mu}_t$ denotes the sample average of Y_i in that group.

The conditional empirical success rule and Stoye’s minmax regret rules condition fully on X . If X takes on a limited number of possible values, these decision rules are natural and may work reasonably well regardless of the underlying parameter P . However, if the support of X is large (as would be the case when some elements of X are continuous, or when X is discrete but of relatively high dimension), the full action space $\mathcal{A} = 2^{\mathcal{X}}$ is very large, and it may be difficult to learn all of the conditional average treatment effects from the data Z . Yet, the results in Stoye (2009) indicate that a minmax regret rule fully conditions on X regardless of the size of \mathcal{X} . In the extreme, if there are no observations in either treatment arm for a given value of X , this can lead to “no-data” rules that employ no smoothing. Intuitively, decision rules that attempt to smooth across x can have poor worst-case risk if $F_0(\cdot|x)$ or $F_1(\cdot|x)$ vary strongly with x . This suggests that the minmax regret criterion may lead to very conservative decision rules that guard against extreme, but perhaps unrealistic, cases.

One possible way to avoid such extreme solutions is to place restrictions on the conditional distributions of potential outcomes, $F_0(\cdot|x)$ and $F_1(\cdot|x)$. For example, Stoye (2012) places bounds on the degree of variation of the conditional distributions of potential outcomes given x , and derives minmax regret rules. If the effect of a certain covariate is known to be sufficiently small, the resulting minmax regret rule will effectively ignore it. Another way to deal with a complex feature space \mathcal{X} is to restrict the set of possible treatment rules \mathcal{A} . Kitagawa and Tetenov (2018) argue that in many applications it may be reasonable to restrict the set of treatment rules under consideration. For example, there may be external constraints that prohibit the use of certain variables in X for treatment assignment, or a restriction to treatment rules based on a linear index $\beta'X$ for ease of implementation. One version of Kitagawa and Tetenov’s approach constructs the treatment assignment rule as

$$\hat{\delta}_n(x; Z) = \arg\max_{\delta \in \mathcal{A}} \hat{W}(\delta),$$

where

$$\hat{W}(\delta) = \frac{1}{n} \sum_{i=1}^n \left[\delta(X_i) \frac{Y_i T_i}{e(X_i)} + (1 - \delta(X_i)) \frac{Y_i (1 - T_i)}{1 - e(X_i)} \right]$$

and \mathcal{A} is the constrained set of treatment assignment rules. In the case of a randomized experiment, $e(x) = \Pr(T_i = 1|X_i = x)$ is known, and the Kitagawa-Tetenov approach uses these known randomization probabilities in the expression above to avoid having to estimate the propensity score nonparametrically (unlike Manski’s conditional empirical success rule, as described above). We can view $\hat{W}(\delta)$ as an estimate of $W(P, \delta)$.

If the constrained set of treatment rules \mathcal{A} is sufficiently small so that it has a finite Vapnik-Chervonenkis (VC) dimension, and some additional conditions (such as boundedness of Y) hold, then one can obtain nontrivial finite-sample probabilistic bounds on

$$\sup_P \sup_{\delta \in \mathcal{A}} |\hat{W}(\delta) - W(P, \delta)|$$

using concentration inequalities. Kitagawa and Tetenov use this approach to bound the maximum regret

welfare of $\hat{\delta}$ and from this show that the worst-case regret welfare of their decision rule converges at the fastest possible rate.

The estimator $\hat{W}(\delta)$ of $W(P, \delta)$ is consistent and asymptotically normal at a \sqrt{n} rate, but it is not efficient in the sense that its asymptotic variance does not attain the semiparametric efficiency bound derived by Hahn (1998). Athey and Wager (2019) propose the use of decision rules based on a semiparametrically efficient estimator of $W(P, \delta)$, especially when the propensity score is unknown. Athey and Wager derive alternative asymptotic characterizations of the welfare regret of statistical treatment rules under VC-type conditions on \mathcal{A} . They obtain lower bounds on risk and find that statistical treatment rules based on inefficient estimates of $W(P, \delta)$ generally do not achieve the bounds, whereas using efficient estimates of $W(P, \delta)$ can attain the bounds up to a universal constant.¹²

In the case where X_i has finite support, it is straightforward to construct asymptotically efficient estimators of $W(P, \delta)$. For example, we can use the following estimator, based on the reasoning in Hirano, Imbens, and Ridder (2003):

$$\tilde{W}(\delta) = \frac{1}{n} \sum_{i=1}^n \left[\delta(X_i) \frac{Y_i T_i}{\hat{e}(X_i)} + (1 - \delta(X_i)) \frac{Y_i (1 - T_i)}{1 - \hat{e}(X_i)} \right],$$

where $\hat{e}(\cdot)$ is a suitable nonparametric estimator of the propensity score. If $\hat{e}(x)$ is the empirical conditional probability of $T_i = 1$ given $X_i = x$, we have

$$\tilde{W}(\delta) = W(\mathbb{P}_n, \delta)$$

and maximizing this with respect to δ leads to Manski's conditional empirical success rules as discussed above.

If X_i has continuous or high-dimensional support, then additional care must be taken to construct the estimator of $W(P, \delta)$, in order to ensure that the estimator is semiparametrically efficient *uniformly* in $\delta \in \mathcal{A}$. Doubly robust estimators, such as those developed by Chernuzhukov, Chetverikov, Demirer, Duflo, Hansen, Newey, and Robins (2018) may be used. See Athey and Wager (2019) for further discussion of this approach. Other related work includes Mbakop and Tabord-Meehan (2016) and Kitagawa and Tetenov (2017).

4.3 Local Asymptotics for Treatment Rules

In the previous subsection we considered statistical treatment rules that maximize an estimate of welfare $W(P, \delta)$ within a restricted class \mathcal{A} . It is reasonable to expect that the performance of the resulting treatment assignment rule will depend on the quality of the estimator of $W(P, \delta)$, regarded as a functional of $\delta \in \mathcal{A}$. Recent work has developed semiparametrically efficient estimators for $W(P, \delta)$ and used them to

¹²See Theorem 5 of Athey and Wager (2019) and the discussion that follows their theorem. Their analysis also uses local parametrizations of the treatment effect in a manner similar to Section 4.3 below.

construct statistical treatment rules. Another approach to dealing with a complex feature space is to restrict the set of possible distributions; for example we could restrict P_θ to lie in some regular parametric or semiparametric class. Then we can use large-sample distribution theory somewhat more directly to simplify the analysis of decision rules. Here we want to explore the possibility for more refined large-sample analysis of statistical treatment assignment rules under statistical regularity conditions. We draw on Hirano and Porter (2009) for the results below.

First consider the case where θ is a finite-dimensional parameter. Suppose the data are i.i.d. with $Z_i \sim P_\theta$ and $Z^n = (Z_1, \dots, Z_n)$. The decision rule is allowed to depend on these observations, so we write it as $\delta_n(x; Z^n)$. If Assumption 1 holds, then the statistical model for Z^n is LAN. That is, in large samples, the statistical model for the data can be approximated by a simple Gaussian statistical model. A consequence of the LAN property is that any sequence of decision rules δ_n can also be approximated by a decision rule in the limiting Gaussian experiment. The following result is drawn from Proposition 3.1 in Hirano and Porter (2009).

Proposition 6 *Suppose $\theta_0 \in \Theta$ and assume Θ is an open subset of \mathbb{R}^k . Consider a sequence of statistical models $P_{\theta_{n,h}}$ where $\theta_{n,h} = \theta_0 + h/\sqrt{n}$, and suppose that Assumption 1 holds. If $E_{\theta_{n,h}}[\delta_n(x; Z^n)]$ has a well defined limit for every h and x , then there exists a decision rule $\delta : \mathcal{X} \times \mathbb{R}^k \rightarrow [0, 1]$ such that*

$$\lim_{n \rightarrow \infty} E_{\theta_{n,h}}[\delta_n(x; Z^n)] = \int \delta(x; \xi) dN(\xi|h, J_0^{-1}) \quad (19)$$

for all $h \in \mathbb{R}^k$ and each $x \in \mathcal{X}$, where $N(\xi|h, J_0^{-1})$ is the Gaussian distribution with mean h and variance J_0^{-1} with J_0 defined in Assumption 1.

This asymptotic representation result states that every converging sequence of decision rules is matched by some decision rule in the Gaussian shift statistical model. One can therefore analyze the decision problem in this relatively simple limiting model to characterize the set of attainable expected welfare and risk functions under various criteria. For example, suppose that there is no covariate, so we can drop X_i from the notation, and suppose that the model is parametrized so that the difference in welfare between the two treatments is

$$g(\theta) = W(\theta, 1) - W(\theta, 0),$$

where the “1” and “0” in the expression above indicate treatment rules that assign all individuals to treatment arms 1 and 0 respectively. Since we will be taking limits as the sample size increases, we focus on θ_0 such that $g(\theta_0) = 0$. Then the local parameter sequences $\theta_0 + h/\sqrt{n}$ for $h \in \mathbb{R}^k$ correspond to cases where the “better” treatment cannot be learned perfectly, even as the sample size increases.

Welfare regret loss is

$$L^R(\theta, \delta) = g(\theta) [\mathbf{1}(g(\theta) > 0) - \delta],$$

Under the local parameter sequence, this will converge to zero, and if the function g is smooth, a scaled

version of the loss has the following limiting form:

$$\sqrt{n}L^R(\theta_0 + h/\sqrt{n}, \delta) \rightarrow L_\infty^R(h, \delta) := \dot{g}'h[\mathbf{1}(\dot{g}'h - \delta)], \quad (20)$$

where \dot{g} is the gradient of g at θ_0 .

The asymptotic representation of decision rules in (19) and the limiting representation of the loss function in (20) yield a local asymptotic version of the statistical decision problem. In the limiting problem, one observes a single draw from the shifted Gaussian model:

$$Z \sim N(h, J_0^{-1}).$$

A decision rule δ maps the observation Z into $[0, 1]$, with the interpretation that future individuals will be assigned treatment 1 with probability δ . The loss of treatment rule δ is $L_\infty^R(h, \delta)$, which is based on a known linear function $\dot{g}'h$ of the parameter h . This is a mild extension of the decision problem considered in Example 8. It can be shown (see Hirano and Porter (2009), Theorem 3.4 and Lemma 5) that the simple cutoff decision rule $\delta^*(Z) = \mathbf{1}(\dot{g}'Z > 0)$ is minmax for regret loss. Let

$$R^* = \sup_{h \in \mathbb{R}^k} E_h[L_\infty^R(h, \delta^*)]$$

be the regret risk of the minmax decision rule in the limiting problem. In the original problem, any converging sequence of decision rules δ_n can be represented by some decision rule in the limiting problem. As a result, for any finite set H of possible values of $h \in \mathbb{R}^k$, it follows that

$$\liminf_{n \rightarrow \infty} \sup_{h \in H} \sqrt{n} E_{\theta_0 + h/\sqrt{n}}[L^R(\theta_0 + h/\sqrt{n}, \delta_n(Z^n))] \geq R^*.$$

Furthermore, under suitable conditions, the feasible decision rule

$$\delta_n(Z^n) = \mathbf{1}(g(\hat{\theta}) > 0)$$

will achieve the lower bound on expected welfare regret, provided that $\hat{\theta}$ is an asymptotically efficient estimator of θ in the sense that $\sqrt{n}(\hat{\theta} - \theta_0 - h/\sqrt{n}) \overset{\theta_0 + h/\sqrt{n}}{\rightsquigarrow} N(0, J_0^{-1})$ for every h .

This argument can be extended to semiparametric settings where the distribution $P \in \mathcal{M}$ of the data Z is not constrained to lie in a finite-dimensional parametric family, but the relevant welfare contrast is a smooth functional of P . In analogy with the $g(\theta)$ notation used above, let $g(P)$ denote the difference in welfare between the two treatments. For example, under utilitarian welfare the welfare contrast is the average treatment effect. If the average treatment effect is point-identified from the observed data distribution P , then it can be written as a functional $g(P)$.

We can analyze decision rules in a similar way to the parametric case. We pick a centering value P_0 of the unknown distribution such that $g(P_0) = 0$. We then need to define local sequences of distributions

around P_0 . The *tangent space* contains one-dimensional submodels of \mathcal{M} that pass through P_0 .¹³ We can view each submodel as a “direction” h away from P_0 , and consider sequences of measures $P_{n,h}$ that satisfy a quadratic mean differentiability condition, as in Assumption 1(a), where h is the score function. If the tangent space is a separable linear space, then h can be expressed as an element $h = (h_1, h_2, \dots)$ of ℓ^2 , the space of square-summable sequences, where the h_j are Fourier coefficients associated with an orthonormal basis of the tangent space.

Following van der Vaart (1991a), an asymptotic representation result similar to (19) holds in this semi-parametric setting and allows us to examine the large sample properties of a sequence of statistical treatment rules through the corresponding problem in the limit experiment. The following result is drawn from Hirano and Porter (2009), Proposition 4.1, which includes an explicit statement of the required differentiability in quadratic mean condition.

Proposition 7 *Suppose that a sequence of decision rules δ_n has limits under every local sequence of measures $P_{n,h}$, in the sense that*

$$\lim_{n \rightarrow \infty} E_{P_{n,h}} [\delta_n(x; Z^n)] \text{ exists.}$$

Then, there exists a decision rule δ such that

$$\lim_{n \rightarrow \infty} E_{P_{n,h}} [\delta_n(x; Z^n)] = E_h [\delta(x; \Delta_1, \Delta_2, \dots)],$$

where $(\Delta_1, \Delta_2, \dots)$ is a sequence of independent random variables with $\Delta_j \stackrel{h}{\sim} N(h_j, 1)$.

Based on this result, we can regard the Gaussian sequence model $(\Delta_1, \Delta_2, \dots)$ as the limit experiment for the analysis of decision rules δ .

For example, in the case where X_i is singleton and we evaluate decision rules by welfare regret loss, Hirano and Porter (2009) show that a simple cutoff decision rule is minmax in the limiting problem. We can then seek a decision rule in the original problem that matches the minmax risk of this decision rule asymptotically. One such decision rule based on the data Z^n is

$$\delta_n(Z^n) = \mathbf{1}(\hat{g}_n(Z^n) > 0)$$

where \hat{g}_n is a semiparametrically efficient estimator of $g(P)$.

4.4 Other Treatment Assignment Problems

We have focused on the case of a binary treatment, but similar analysis is possible in other related assignment or allocation problems. For example, there may be more than two treatment arms available.

¹³See, for example, Bickel, Klaassen, Ritov, and Wellner (1993), van der Vaart (1991a), and van der Vaart (1998) for discussions of tangent spaces in semiparametric statistics.

Then results from the literature on estimating multi-valued treatment effects (e.g., Imbens (2000), Lechner (2001), Cattaneo (2010), Imai and van Dyk (2004), Hirano and Imbens (2004)) provide a starting point for developing treatment assignment rules mapping some feature space \mathcal{X} into the set of available treatments, as in recent work by Kallus and Zhou (2018) and Demirer, Syrgkanis, Lewis, and Chernuzhukov (2019). Other work has considered assignment of binary treatments under different types of budget constraints; see for example Bhattacharya and Dupas (2012) and Adusumilli, Geiecke, and Schilter (2020).

The literature on estimating *dynamic treatment regimes*, beginning with Murphy (2003) and Robins (2004), considers assignment rules that choose a sequence of treatments dynamically based on past outcomes of the individual. See Chakraborty and Moodie (2013) and Chakraborty and Murphy (2014) for surveys of the literature in statistics and biostatistics on dynamic treatment regimes. A related decision problem is to target coupons or other marketing interventions to individual characteristics and past purchase history; see for example Rossi, McCulloch, and Allenby (1996) and Dubé, Fang, Fong, and Luo (2017). In adaptive treatment assignment problems, which can be viewed as combining aspects of treatment assignment in the sense we have described above, and dynamic experimental design, raise additional issues. They have been extensively studied in the form of bandit problems; recent work includes Kock and Thyrgaard (2017).

When there are peer (social interaction) effects, the problem of allocating individuals to different peer groups (for example, assigning students to classrooms), and the problem of assigning treatments to individuals taking into account the potential effect on peers, raise additional conceptual and practical challenges. Recent work on allocation rules under peer effects includes Graham, Imbens, and Ridder (2014) and Bhattacharya (2009), among others.

5 Other Topics

In this section we consider other applications where viewing procedures as statistical decision rules can provide useful insights.

5.1 Nonstandard Functionals

In some economic applications, the underlying model satisfies conventional regularity conditions, but the nature of the decision problem leads to nonstandard distributional theory. Suppose interest centers on some function of the parameter $\kappa(\theta)$ where $\kappa : \Theta \rightarrow \mathbb{R}$. Under conventional smoothness conditions on the sequence of experiments $\mathcal{E}^n = \{P_\theta^n : \theta \in \Theta\}$, the MLE $\hat{\theta}_{ML}$ and other estimators such as the Bayes estimator are asymptotically efficient. However, the limit distributions of derived estimators of $\kappa(\theta)$, and more generally the feasible limit distributions of any estimators of $\kappa(\theta)$, will depend crucially on the smoothness in κ at any θ_0 that is the centering point of possible localized sequences. This may lead to nonstandard limiting distributions for estimators and other decision rules targeting κ .

In some important economic applications, the parameter of interest κ is a directionally differentiable, but not fully differentiable, functional of the distribution of the data. For example, suppose that $\theta = (\theta_1, \theta_2)$ is a two-dimensional parameter and the statistical model is $Z^n = (Z_1, \dots, Z_n) \sim P_\theta^n$. Suppose the parametrization by θ is such that the model satisfies Assumption 1 so that local asymptotic normality holds, and let $\hat{\theta}$ be a regular estimator of θ as discussed in Section 3. However, suppose we are interested in the following function of θ :

$$\kappa(\theta) := \min\{\theta_1, \theta_2\}.$$

This type of estimand can arise in partially identified econometric models based on moment inequalities, as we discuss further in Section 5.2 below. It can also arise in other applications such as the bounds on valuation distributions, as in Haile and Tamer (2003), or inference on the best treatment as related to treatment assignment in Section 4, see Hirano and Porter (2012) for additional examples. In many of these examples, the parameter θ is a reduced form parameter, and κ is a structural or latent quantity that is of primary interest.

The function $\kappa(\cdot)$ is well behaved in many respects, being continuous, homogeneous of degree one, and directionally differentiable, but it is not fully differentiable. Consider the natural analog estimator

$$\hat{\kappa} = \min\{\hat{\theta}_1, \hat{\theta}_2\}.$$

If the true value of (θ_1, θ_2) satisfies $\theta_1 \neq \theta_2$, then by application of the Delta Method, $\hat{\kappa}$ will be asymptotically normal and centered at the true value of κ . However, if θ_1 is close to θ_2 , then the estimator $\hat{\kappa}$ will be downward biased and the normal distribution will not provide an accurate approximation to its sampling distribution.

We can capture this situation formally by considering a localization where the two parameters are within $O(1/\sqrt{n})$ of each other:

$$\begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} = \begin{pmatrix} \theta_0 + h_1/\sqrt{n} \\ \theta_0 + h_2/\sqrt{n} \end{pmatrix}.$$

Under this localization we will obtain a nonnormal limit distribution for $\hat{\kappa}$, but the implications are deeper. In particular, results in van der Vaart (1991b) imply that there exists *no* regular (locally asymptotically translation equivariant) estimator for κ , Hirano and Porter (2012) show that there exist no estimators that are locally asymptotically unbiased or quantile-unbiased. These results can be derived by considering the estimation problem in the limiting normal shift model. Fang and Santos (2019) develop inference methods for directionally differentiable parameters, and Fang (2016) develops locally asymptotically minmax estimators within a class of plug-in estimators.

Another set of econometric models exhibiting nonstandard distributional theory are models exhibiting

weak identification. Consider a classical linear instrumental variables model:

$$\begin{aligned}x_i &= \theta_1' w_i + v_{1i}, \\y_i &= \kappa x_i + \epsilon_i,\end{aligned}$$

where y_i is the outcome of primary interest, x_i is a scalar endogenous explanatory variable, and w_i is a vector of exogenous instruments. We wish to estimate the structural parameter κ . The model can be written in reduced form as

$$\begin{aligned}x_i &= \theta_1' w_i + v_{1i}, \\y_i &= \theta_2' w_i + v_{2i},\end{aligned}$$

with $\theta_2 = \kappa \theta_1$. This restricts the possible values of $\theta = (\theta_1, \theta_2)$ in general. Subject to this restriction, we can view κ as a function of the reduced form parameters: $\kappa = \kappa(\theta)$.

If the instrument is weakly correlated with the endogenous regressor, then it is well known that conventional asymptotic distribution theory provides a poor approximation to the distributions of estimators of the structural parameters. Staiger and Stock (1997) proposed to consider a local parameter sequence in which the coefficient on the instrument in the first stage is local to zero:

$$\theta_{1n} = \frac{h}{\sqrt{n}}.$$

Under this parameter sequence, one obtains non-standard limit distributions for estimators such as 2SLS and LIML that, in some cases, well approximate their finite-sample distributions. A large literature including Stock and Wright (2000), Kleibergen (2002), Moreira (2003), and Andrews, Moreira, and Stock (2006) have used this type of local asymptotic parametrization to study estimation and inference procedures for weakly identified models.

The Staiger-Stock local parametrization chooses a particular centering for θ_{1n} in order to approximate the situation where the correlation of the instrument with the endogenous regressor is close to zero. Subject to this particular choice for the centering, however, the localization is the standard one. Cattaneo, Crump, and Jansson (2012) point out that as a result, conventional LAN theory applies to the reduced form parameter θ , and use the limits of experiments framework to study optimal inference. Nonstandard behavior of the estimator of κ arises because κ depends on θ in a non-smooth fashion when θ_1 is close to zero. Hirano and Porter (2015) show that due to the non-smoothness of $\kappa(\theta)$, no locally asymptotically unbiased estimators for β exist. Andrews and Armstrong (2017) show that unbiased estimation of β is possible if the sign of the first-stage coefficient is known.

For both types of problems considered in this subsection, we can regard the problem as being regular in its reduced form, where we use “reduced form” in the classic econometric sense to mean describing the distribution of the observable data. In our notation, this means that the statistical model \mathcal{M} has a

parametrization in terms of some θ such that the model is locally asymptotically normal. The irregularity arises because the “structural” parameter of interest κ is related to θ in a non-smooth way.

While the impossibility results cited above suggest that standard approaches to estimation and inference may not be available in these cases, recent work has proposed creative new approaches. Kaji (2017) develops an alternative criterion for efficiency of weakly identified structural parameters based on minimal sufficiency in the reduced form. And as mentioned in section 2.4.3, Müller and Wang (2019) suggest a numerically attractive approach for proceeding with efficient estimation under constraints in the presence of non-regularity. Chen, Christensen, and Tamer (2018) consider the construction of confidence sets for structural parameters, and use level sets of the quasi-posterior of the quasi-likelihood ratio to obtain correct coverage even when the model exhibits singularities as above.

5.2 Partial Identification

Recently, a large literature has developed dealing with estimation and inference in partially identified models (see e.g. Manski (1995, 2003, 2007), Tamer (2012)). The chapter by Molinari (2019) provides an extensive survey of partial identification in econometrics. In partially identified models, we observe $Z \sim P$, but knowledge of P does not pin down the value of the parameters of interest. Formally, let $\theta \in \Theta$ be the parameter of some underlying latent variable or structural model. The value of θ determines the distribution of the observable data: $Z \sim P = P_\theta$, but different values of the parameter may be observationally equivalent in the sense that

$$P_{\theta'} = P_{\theta''} \text{ for some } \theta', \theta'' \in \Theta \text{ with } \theta' \neq \theta''.$$

Let $\Theta(P) \subset \Theta$ be the set of values for θ consistent with $Z \sim P$:

$$\Theta(P) = \{\theta \in \Theta : P_\theta = P\}.$$

This is the set-valued inverse of the mapping $\theta \mapsto P_\theta$. If interest centers on a subparameter, say $\kappa := \kappa(\theta)$, let $K(P) = \kappa(\Theta(P))$ be the set of values of κ consistent with $Z \sim P$. We say the model is point identified if $\theta' \neq \theta'' \Rightarrow P_{\theta'} \neq P_{\theta''}$, in other words if the mapping $\theta \mapsto P_\theta$ is injective. That is, point identification implies that $\Theta(P_\theta) = \{\theta\}$ for all θ . We say that the model is partially identified if $\Theta(P)$ is not singleton, but is a strict subset of Θ for at least some values of P .

Whether or not the statistical model is point identified, the general statistical decision framework of Section 2 can still be applied. However, lack of point identification raises additional technical and conceptual issues. To illustrate some of these issues in a simple setting, consider the following simple moment inequality problem. Suppose we are interested in a scalar subparameter $\kappa \in \mathbb{R}$, and its identified set can be represented by a set of moment inequalities

$$K(P) = \{\kappa : E_P[m(Z, \kappa)] \geq 0\},$$

where $m(Z, \kappa)$ is a given (vector-valued) moment function.

The standard point estimation problem, where the action space is $\mathcal{A} = \mathbb{R}$, is still well defined. Given a point estimator $\hat{\kappa} : \mathcal{Z} \rightarrow \mathbb{R}$, we could apply one of the loss functions considered in Section 3, such as squared error loss $(\kappa - \hat{\kappa})^2$. However, because κ cannot be perfectly learned even in large samples, the risk of the point estimator will not converge to 0 in general, leading to different considerations for the asymptotic analysis. See Aryal and Kim (2013) and Song (2014) for examples of point decisions in partially identified settings.

Alternatively, we can consider “point” estimation of the set $K(P)$ by a set \hat{K} . Here the action space \mathcal{A} is a collection of subsets of \mathbb{R} . In this case we need to posit a loss function compatible with the action space. For example, we could set $L(P, \hat{K}) = d(K(P), \hat{K})$, where $d(\cdot)$ is some measure of distance between the identified set $K(P)$ and the set estimate \hat{K} , such as the Hausdorff metric.

As a concrete example, suppose that (Y_i, V_i, W_i) are i.i.d. for $i = 1, \dots, n$, and we are interested in $\kappa := E[Y_i]$. Both V_i and W_i have nonnegative means, but otherwise we put no restrictions on the joint distribution of (Y_i, V_i, W_i) . We can think of the parameter θ as indexing this joint distribution, and κ as a subparameter of θ . Suppose we only observe $Z_i = (Z_{i1}, Z_{i2})$, where

$$\begin{aligned} Z_{i1} &= Y_i - V_i, \\ Z_{i2} &= Y_i + W_i. \end{aligned}$$

In other words, Z_{i1} and Z_{i2} are downward and upward biased measures of κ , respectively. The observed data is $Z^n = (Z_1, \dots, Z_n)$ where Z_i are i.i.d. P . Then a natural moment function for κ is

$$m(Z_i, \kappa) = \begin{pmatrix} \kappa - Z_{i1} \\ Z_{i2} - \kappa \end{pmatrix}.$$

Let $\mu_1 = \mu_1(P) = E_P[Z_{i1}]$ and $\mu_2 = \mu_2(P) = E_P[Z_{i2}]$. Then the identified set is

$$K(P) = [\mu_1, \mu_2].$$

Suppose the action space \mathcal{A} consists of closed intervals in \mathbb{R} , so that estimators \hat{K} have the form

$$\hat{K}(Z^n) = [\delta_1(Z^n), \delta_2(Z^n)].$$

In this case a simple loss function that could be used is

$$L(P, \hat{K}) = (\delta_1 - \mu_1)^2 + (\delta_2 - \mu_2)^2.$$

This is the sum of the squared error losses of the estimates of the endpoints of the identified set. In this case, the decision problem is equivalent to a two-dimensional point estimation problem, and an analysis

similar to that in Section 3 can be applied.¹⁴

On the other hand, suppose $Z_{i3} = Y_i + V_i$ with $\mu_3 = \mu_3(P) = E_P[Z_{i3}]$, and we only observe $\tilde{Z}_i = (Z_{i2}, Z_{i3})$. Then, Z_{i2} and Z_{i3} provide two upward biased measures of κ , leading to the moment function:

$$\tilde{m}(\tilde{Z}_i, \kappa) = \begin{pmatrix} Z_{i2} - \kappa \\ Z_{i3} - \kappa \end{pmatrix}.$$

The resulting identified set is

$$K(P) = (-\infty, \min\{\mu_2, \mu_3\}].$$

Suppose the action space consists of sets of the form $\hat{K} = (-\infty, \delta]$ for $\delta \in \mathbb{R}$. We could use the loss function $L(P, \hat{K}) = (\min\{\mu_2, \mu_3\} - \delta)^2$, which measures the closeness of the estimated upper bound on κ to the true upper bound. Thus the decision problem is equivalent to point estimation of $\min\{\mu_2, \mu_3\}$, but since the function $h(P) = \min\{\mu_2(P), \mu_3(P)\}$ is not differentiable in P , the estimation problem is nonstandard, as we discussed in Section 5.1. For example, no regular estimators of $h(P)$ exist, and locally unbiased estimation of $h(P)$ is not possible.

These simple examples can be generalized to many other decision problems with partially identified parameters. If the subparameter of interest κ is a d -dimensional vector, then its identified set $K(P)$ will be a subset of \mathbb{R}^d . If $K = K(P)$ is convex, then it can be characterized by its *support function*

$$h_K(x) = \sup_{v \in K} x'v, \quad x \in \mathbb{R}^d,$$

which encodes its supporting hyperplanes in all directions. The support function approach has been used by Beresteanu and Molinari (2008), Bontemps, Magnac, and Maurin (2012), and Kaido (2016), among others. Then, estimation of $K(P)$ can be viewed as a functional point estimation problem. If the support function of $K(P)$ is smooth (as a function of P), then classical results on efficiency bounds for estimation of functions, such as those in Bickel, Klaassen, Ritov, and Wellner (1993), can be applied. Kaido and Santos (2014) use this approach to obtain asymptotically efficient estimators for models defined by convex moment inequalities, under sufficient smoothness. On the other hand, in many econometric applications $K(P)$ is not differentiable in P . This can arise, for example, in moment inequality problems when the number of binding (or nearly binding) inequalities at a given point on the boundary of the identified set is greater than d . Then $K(P)$ will not be differentiable in P , and standard asymptotic optimality theory does not apply.

¹⁴However, there is an additional implicit inequality: $\mu_1 \leq \mu_2$. If μ_2 is sufficiently greater than μ_1 , then in large samples this constraint will have a negligible impact on the analysis. But if μ_1 is arbitrarily close to μ_2 in the collection of measures P , then the analysis is more subtle.

5.3 Confidence Intervals and Sets

Confidence intervals and confidence sets can be regarded as rules of the form $\delta : \mathcal{Z} \rightarrow \mathcal{A}$, where \mathcal{A} is some collection of subsets of the Euclidean parameter space Θ . For example, if θ is a scalar parameter we could consider rules that take the data Z and produce a closed interval. In the classical Neyman-Pearson approach to constructing confidence intervals, we require that, for any $\theta \in \Theta$, if $Z \sim P_\theta$ then $\delta(Z)$ contains θ with probability equal to or greater to some prespecified value, such as 0.95. This coverage condition amounts to a restriction on the set of decision rules, as discussed in Section 2.4.3. Subject to this restriction, we could seek to find a rule that generates small sets, by using a loss function that depends on the volume of the set or some other criterion.

An alternative decision theoretic approach could entertain a trade-off between coverage and length or volume of the set. Again suppose $\delta : \mathcal{Z} \rightarrow \mathcal{A}$ where \mathcal{A} is a collection of subsets of Θ . A generic form of the loss function for confidence sets would be

$$L(\theta, a) = \ell(\text{vol}(a), \chi(\theta, a)),$$

where $\text{vol}(a)$ denotes the length or volume of the set a , and $\chi(\theta, a)$ measures lack of coverage or precision of the set. The most common choice for χ is

$$\chi(\theta, a) = \mathbf{1}\{\theta \notin a\}.$$

When loss is additively separable in the arguments $\text{vol}(a)$ and $\chi(\theta, a)$, the part of risk that comes from the $\chi(\theta, a)$ term directly measures coverage. A modification of χ to penalize non-coverage depending on the distance between the parameter θ and the set δ would set

$$\chi(\theta, a) = \inf_{t \in a} \|\theta - t\|$$

where $\|\cdot\|$ is used to denote a Euclidean distance.

There is a long history of work adopting a decision theoretic perspective on the construction of confidence sets. In a setting with a scalar parameter, Winkler (1972) and Cohen and Strawderman (1973) consider specific choices for loss and analyze Bayes rules and admissability. Hwang and Casella (1982) consider restricting the space of confidence set rules to satisfy a minimum coverage requirement and then minimize a loss that depends only on volume. They show that when $\dim(\theta) \geq 4$, confidence sets centered at James-Stein estimators can be constructed that dominate classical confidence intervals. Evans, Hansen, and Stark (2005) consider minimax expected volume (length) subject to a coverage condition for a scalar parameter on a bounded space. Casella, Hwang, and Robert (1993) note that using a loss function that is linear in lack of coverage and volume can lead to unappealing decision rules. Casella, Hwang, and Robert (1994) and Rice, Lumley, and Szpiro (2008) consider various loss functions that avoid the problem pointed out in Casella, Hwang, and Robert (1993) and provide a decision theoretic analysis

of confidence set rules for these loss functions.

Given the prevalence of non-regular models and partial identification problems in econometrics, there is potential for further work on decision theoretic analysis of confidence set rules under these conditions. For example, Chamberlain (2007) considers a loss function that depends on volume and coverage and uses invariance properties to develop a Bayesian confidence procedure in an instrumental variables setting with possibly weak instruments. Müller and Norets (2016) consider an alternative bet-proofness criterion to construct confidence sets in nonstandard problems.

5.4 Experimental and Data Collection Design

Experimental design, and the more general problem of choosing a data collection scheme (such as choosing sample weights in a stratified survey) is a classical problem in statistics and there is a vast literature on this topic. Some standard texts on experimental design include Fisher (1935), Cochran and Cox (1957), Cox (1958), and Pukelsheim (2006); discussions of survey and sampling design include Kish (1965), Manski and McFadden (1981), and Thompson (2012). There has been renewed interest in data collection design in recent years, for various reasons including: the increasing use of randomized controlled trials in microeconomics; the development of new surveys with design input from social scientists; and technological advances that facilitate new data collection, such as online experiments and richer observational data sources. As a result, the design of experiments and surveys remains an active area of research, often with immediate applications in economics and other fields.

The decision theoretic framework we have outlined in this chapter can provide a useful set of organizing concepts for a range of data collection problems. However, there are also some limitations of the framework that we shall discuss below. To apply our framework, we must specify an action space \mathcal{A} , a parameter space Θ , a random variable $Z \sim P_\theta$, and an evaluation criterion in the form of a welfare function $W(\theta, a)$ or a loss function $L(\theta, a)$.

The action space embodies the choices available to the designer of the experiment or survey. For example, in a simple randomized control trial, the experimenter could choose the probability of assigning individuals to treatment $p \in (0, 1)$. Then we could set $\mathcal{A} = (0, 1)$. In a stratified survey, there is a characteristic X , usually discrete taking values in some set $\{\xi_1, \dots, \xi_K\}$. A stratification scheme consists of a weighting $p = (p_1, \dots, p_K)$ in the simplex Δ^{K-1} , with the interpretation that individuals with $X = \xi_k$ are included in the sample with probability p_k . Then the action space \mathcal{A} is the simplex Δ^{K-1} or some pre-specified subset of the simplex. Other design problems involve different choices for the action space. Economically motivated mechanisms for data design can be considered, as in Philipson (1997), Narita (2018), and others.

The parameter $\theta \in \Theta$ characterizes all relevant features of the population; usually at least some components of θ are unknown. Hence there is a role for considering the different outcomes that may result

from a given design under different values of θ .

Recall that $Z \sim P_\theta$ represents the data available *before* the action is chosen. Here, Z would correspond to data available in advance of specification of the survey or experiment. In some cases, no such data may be available. In other cases, there may be data from prior surveys or experiments that are informative about θ . For example, Sukhatme (1935) and Solomon and Zacks (1970) consider the use of prior samples to select a stratified sampling scheme. Hahn, Hirano, and Karlan (2011) consider using data from the initial wave of a two-wave experiment to choose conditional treatment probabilities in the second wave; see also Chambaz, van der Laan, and Zheng (2015) and Tabord-Meehan (2018). More generally, treatment assignment could be sequential or in multiple waves. Hu and Rosenberger (2006) surveys the theory of response-adaptive randomization; for recent work on multi-stage dynamic experimental design, see among others Kasy and Sautmann (2019) and Xiong, Athey, Bayati, and Imbens (2019).

To apply the framework of this chapter we need to specify a welfare function $W(\theta, a)$, or, equivalently, a loss function $L(\theta, a)$. The choice of the evaluation criterion may be motivated by the intended use of the experiment or survey. In practice, there may not be a single obvious choice for W or L . One way to specify the evaluation criterion is to use a general measure of informativeness, such as Fisher information or some generalization of the information as in Lin, Martin, and Yang (2019). An alternative is to focus on some particular estimator (and estimand) that will be applied to the data collected from the experiment or survey, and evaluate its expected loss under different design choices a .

Example 9 *Suppose the decision maker will run a randomized controlled trial, assigning individuals to treatments 0 and 1. Suppose potential outcomes $Y(0)$ and $Y(1)$ are distributed as*

$$Y(0) \sim N(\mu_0, \sigma_0^2)$$

$$Y(1) \sim N(\mu_1, \sigma_1^2)$$

The experiment will produce observations on Y_i, T_i . Suppose T_i is Bernoulli with probability $a \in (0, 1)$ and Y_i is the observed outcome, satisfying $Y_i = T_i Y_i(1) + (1 - T_i) Y_i(0)$. The decision maker plans to estimate the average treatment effect $\mu_1 - \mu_0$ with the simple difference in means estimator

$$\hat{\tau} = \frac{1}{n_1} \sum_{i=1}^n T_i Y_i - \frac{1}{n_0} \sum_{i=1}^n (1 - T_i) Y_i,$$

where $n_1 = \sum_{i=1}^n T_i$ and $n_0 = n - n_1$. The estimator is unbiased (conditional on both treatment cells being nonempty), and its large-sample normalized variance is

$$\frac{\sigma_1^2}{na} + \frac{\sigma_0^2}{n(1-a)}.$$

Then we could take $\theta = (\mu_0, \mu_1, \sigma_0, \sigma_1)$ and the loss function to be

$$L(\theta, a) = \frac{\sigma_1^2}{na} + \frac{\sigma_0^2}{n(1-a)}. \quad (21)$$

If σ_0 and σ_1 are known (i.e., Θ is degenerate along the relevant dimensions), then it is straightforward to minimize (21) to obtain the optimal randomization probability

$$a^*(\theta) = \frac{\sigma_1}{\sigma_0 + \sigma_1}.$$

In Example 9, there is an optimal action given knowledge of σ_0^2 and σ_1^2 . If the variances are not known, then this optimal action is not feasible. Following our framework, we can adopt a Bayesian approach by specifying a prior distribution over the parameters, or a minmax or minmax regret approach. For discussions of Bayesian experimental design see Spiegelhalter, Freedman, and Parmar (1994) and Chaloner and Verdinelli (1995). For minmax-regret approaches to experiment and data design, see for example Schlag (2007), Stoye (2012), Manski and Tetenov (2016), and Dominitz and Manski (2017).

An important practical issue in experimental design is stratified randomization. If a binary treatment is simply randomized over a sample of individuals, this can lead to *ex post* imbalance of covariates between treated and control groups. As a consequence, estimators of treatment effects can have higher variance than is possible if the treatment is assigned in a way that reduces covariate imbalance. Bai (2019) shows that certain matched pair designs minimize the variance of the estimator $\hat{\tau}$ (defined in Example 9) among randomized stratification schemes. Another approach to improve covariate balance is to re-randomize. Morgan and Rubin (2012) propose to establish rules for re-randomization and account for these rules in inference. In Example 9, for instance, one could use an assignment rule that re-randomizes if either treatment cell is empty. See Athey and Imbens (2017) for a recent survey of approaches to stratified randomization.

When peer effects (or “interference” in the classical experimental design terminology) are of interest, one may wish to design experiments that are informative about the magnitude of such effects. Economic studies that use randomized experiments to measure peer effects include Duflo and Saez (2003), Miguel and Kremer (2003), and Angelucci and De Giorgi (2009). The design of experiments to measure peer effects has been studied in Hirano and Hahn (2010) and Baird, Bohren, McIntosh, and Özler (2017).

The Wald framework we have emphasized in this chapter is a single-agent theory of statistical decision-making. As such, it does not easily handle strategic considerations that could arise from the behavior of the subjects of the experiment or survey, unless they can be subsumed into the loss or welfare function. Some recent work has explicitly modeled design and other statistical decision problems as games between some designer and her human subjects. Some examples include Chassang, Padró i Miquel, and Snowberg (2012), Tetenov (2016), and Spiess (2018).

6 Conclusion

In this chapter we have discussed the construction and analysis of decision rules in econometrics, focusing on obtaining useful approximations to their risk (or expected welfare) properties. The framework is quite general and provides useful insights not only into conventional statistical activities such as estimation, but also into policy decision-making (through empirical auction design and treatment assignment rules for example), design of experiments and surveys, point forecasting, and other empirical problems. There are many other potential applications of the framework which we have not covered here, such as the construction and comparison of forecast intervals (e.g. Christoffersen (1998) and Askanazi, Diebold, Schorfheide, and Shin (2018)) and forecast densities (e.g. Diebold, Gunther, and Tay (1998), Tay and Wallis (2000), and Hall and Mitchell (2004)).

We focused on local asymptotic methods as a tool for obtaining relatively simple characterizations of the risk/welfare properties of decision rules. Local asymptotics can be used to study not only smooth settings where normality emerges naturally, but also other settings involving parameters that do not change smoothly with the underlying distribution of the data, rare events, and nonstationary time series. Of course, depending on the problem, other approaches such as the use of concentration inequalities to obtain risk bounds, and the use of global nonparametric approximations (as in Brown and Low (1996), Nussbaum (1996), and Armstrong and Kolesár (2018)) may provide more useful characterizations of complex decision rules.

References

- ABADIE, A., AND M. KASY (2019): “Choosing among Regularized Estimators in Empirical Economics: the Risk of Machine Learning,” forthcoming, *Review of Economics and Statistics*.
- ADUSUMILLI, K., F. GEIECKE, AND C. SCHILTER (2020): “Dynamically Optimal Treatment Allocation using Reinforcement Learning,” working paper.
- ANDREWS, D. W. K., M. J. MOREIRA, AND J. H. STOCK (2006): “Optimal Two-Sided Invariant Similar Tests for Instrumental Variables Regression,” *Econometrica*, 74(3), 715–752.
- ANDREWS, I., AND T. B. ARMSTRONG (2017): “Unbiased Instrumental Variables Estimation under Known First-Stage Sign,” *Quantitative Economics*, 8(2), 479–503.
- ANGELUCCI, M., AND G. DE GIORGI (2009): “Indirect Effects of an Aid Program: How do Cash Transfers Affect Ineligibles’ Consumption?,” *American Economic Review*, 99(1), 486–508.
- ANSCOMBE, F. J., AND R. J. AUMANN (1963): “A Definition of Subjective Probability,” *Annals of Mathematical Statistics*, 34, 199–205.
- ARJAS, E., AND O. SAARELA (2010): “Optimal Dynamic Regimes: Presenting a Case for Predictive Inference,” *The International Journal of Biostatistics*, 6(2).
- ARMSTRONG, T. B., AND M. KOLESÁR (2018): “Optimal Inference in a Class of Regression Models,” *Econometrica*, 86(2), 655–683.
- ARYAL, G., AND D.-H. KIM (2013): “A Point Decision for Partially Identified Auction Models,” *Journal of Business and Economic Statistics*, 31(4), 384–397.
- ASKANAZI, R., F. X. DIEBOLD, F. SCHORFHEIDE, AND M. SHIN (2018): “On the Comparison of Interval Forecasts,” *Journal of Time Series Analysis*, 39(6), 953–965.
- ASSUNÇÃO, J., R. MCMILLAN, J. MURPHY, AND E. SOUZA-RODRIGUES (2019): “Optimal Environmental Targeting in the Amazon Rainforest,” working paper.
- ATHEY, S., AND G. W. IMBENS (2017): “The Econometrics of Randomized Experiments,” in *Handbook of Economic Field Experiments*, ed. by E. Duflo, and A. Banerjee, vol. 1, pp. 73–140. Elsevier.
- ATHEY, S., AND S. WAGER (2019): “Efficient Policy Learning,” working paper.
- AVRAMOV, D., AND G. ZHOU (2010): “Bayesian Portfolio Analysis,” *Annual Review of Financial Economics*, 2(1), 25–47.
- BAI, Y. (2019): “Optimality of Matched-Pair Designs in Randomized Controlled Trials,” working paper.
- BAIRD, S., J. A. BOHREN, C. MCINTOSH, AND B. ÖZLER (2017): “Optimal Design of Experiments in the Presence of Interference,” working paper.

- BARANCHIK, A. J. (1964): "Multiple Regression and Estimation of the Mean of a Multivariate Normal Distribution," Stanford University Department of Statistics Technical Report 51.
- BARBERIS, N. (2000): "Investing for the Long Run when Returns Are Predictable," *The Journal of Finance*, 55(1), 225–264.
- BERESTEANU, A., AND F. MOLINARI (2008): "Asymptotic Properties for a Class of Partially Identified Models," *Econometrica*, 76(4), 763–814.
- BERGER, J. O. (1993): *Statistical Decision Theory and Bayesian Analysis*. Springer, New York, 2nd edn.
- BERNARDO, J. M., AND A. F. M. SMITH (1994): *Bayesian Theory*. John Wiley & Sons, New York.
- BHATTACHARYA, D. (2009): "Inferring Optimal Peer Assignment from Experimental Data," *Journal of the American Statistical Association*, 104, 486–500.
- BHATTACHARYA, D., AND P. DUPAS (2012): "Inferring Welfare Maximizing Treatment Assignment under Budget Constraints," *Journal of Econometrics*, 167, 168–196.
- BICKEL, P. J., C. A. KLAASEN, Y. RITOV, AND J. A. WELLNER (1993): *Efficient and Adaptive Estimation for Semiparametric Models*. Springer-Verlag, New York.
- BLACK, D., J. SMITH, M. BERGER, AND B. NOEL (2003): "Is the Threat of Training More Effective than Training Itself? Experimental Evidence from the UI System," *American Economic Review*, 93(4), 1313–1327.
- BLACKWELL, D. A., AND M. A. GIRSHICK (1979): *Theory of Games and Statistical Decisions*. Dover, New York.
- BONTEMPS, C., T. MAGNAC, AND E. MAURIN (2012): "Set Identified Linear Models," *Econometrica*, 80(3), 1129–1155.
- BOUSQUET, O., S. BOUCHERON, AND G. LUGOSI (2004): "Introduction to Statistical Learning Theory," in *Advanced Lectures on Machine Learning*, pp. 169–207. Springer-Verlag.
- BREIMAN, L. (1996): "Bagging Predictors," *Machine Learning*, 36, 105–139.
- BROWN, L. D., AND M. G. LOW (1996): "Asymptotic Equivalence of Nonparametric Regression and White Noise," *The Annals of Statistics*, 24(6), 2384–2398.
- BÜHLMANN, P., AND B. YU (2002): "Analyzing Bagging," *Annals of Statistics*, 30, 927–961.
- BUSHWAY, S., AND J. SMITH (2007): "Sentencing Using Statistical Treatment Rules: What We Don't Know Can Hurt Us," *Journal of Quantitative Criminology*, 23(4), 377–387.
- CASELLA, G., J. T. G. HWANG, AND C. ROBERT (1993): "A Paradox in Decision-Theoretic Interval Estimation," *Statistica Sinica*, 3, 141–155.

- (1994): “Loss Functions for Set Estimation,” in *Statistical Decision Theory and Related Topics V*, ed. by S. S. Gupta, and J. O. Berger. Springer, New York.
- CATTANEO, M. D. (2010): “Efficient Semiparametric Estimation of Multi-Valued Treatment Effects under Ignorability,” *Journal of Econometrics*, 155, 138–154.
- CATTANEO, M. D., R. K. CRUMP, AND M. JANSSON (2012): “Optimal Inference for Instrumental Variables Regression with Non-Gaussian Errors,” *Journal of Econometrics*, 167(1), 1–15.
- CHAKRABORTY, B., AND E. E. M. MOODIE (2013): *Statistical Methods for Dynamic Treatment Regimes*. Springer, New York.
- CHAKRABORTY, B., AND S. A. MURPHY (2014): “Dynamic Treatment Regimes,” *Annual Reviews of Statistics and Its Application*, 1, 447–464.
- CHALONER, K., AND I. VERDINELLI (1995): “Bayesian Experimental Design: A Review,” *Statistical Science*, 10(3), 273–304.
- CHAMBAZ, A., M. VAN DER LAAN, AND W. ZHENG (2015): “Targeted Covariate-Adjusted Response-Adaptive LASSO-Based Randomized Controlled Trials,” in *Modern Adaptive Randomized Clinical Trials: Statistical, Operational, and Regulatory Aspects*, ed. by O. Sverdlov. CRC Press, New York.
- CHAMBERLAIN, G. (2007): “Decision Theory Applied to an Instrumental Variables Model,” *Econometrica*, 75(3), 609–652.
- (2011): “Bayesian Aspects of Treatment Choice,” in *The Oxford Handbook of Bayesian Econometrics*, ed. by J. Geweke, G. Koop, and H. van Dijk, pp. 11–39. Oxford University Press.
- CHASSANG, S., G. PADRÓ I MIQUEL, AND E. SNOWBERG (2012): “Selective Trials: A Principal-Agent Approach to Randomized Control Trials,” *American Economic Review*, 102(4), 1279–1309.
- CHAWLA, S., J. D. HARTLINE, AND D. NEKIPELOV (2017): “Mechanism Redesign,” working paper.
- CHEN, X., T. M. CHRISTENSEN, AND E. TAMER (2018): “Monte Carlo Confidence Sets for Identified Sets,” *Econometrica*, 86(6), 1965–2018.
- CHERNUZHUKOV, V., D. CHETVERIKOV, M. DEMIRER, E. DUFLO, C. HANSEN, W. NEWEY, AND J. ROBINS (2018): “Double/Debiased Machine Learning for Treatment and Structural Parameters,” *The Econometrics Journal*, 21, C1–C68.
- CHRISTOFFERSEN, P. F. (1998): “Evaluating Interval Forecasts,” *International Economic Review*, 39, 841–862.
- CHRISTOFFERSEN, P. F., AND F. X. DIEBOLD (1997): “Optimal Prediction Under Asymmetric Loss,” *Econometric Theory*, 13(6), 808–817.

- CLAESKENS, G., AND N. L. HJORT (2008): *Model Selection and Model Averaging*. Cambridge University Press, Cambridge.
- COCHRANE, W. G., AND G. M. COX (1957): *Experimental Designs*. Wiley, New York.
- COHEN, A., AND W. E. STRAWDERMAN (1973): “Admissible Confidence Interval and Point Estimation for Translation or Scale Parameters,” *The Annals of Statistics*, 1(3), 545–550.
- COX, D. R. (1958): *Planning of Experiments*. Wiley, New York.
- DAVIES, R. B. (1985): “Asymptotic Inference when the Amount of Information is Random,” in *Proceedings of the Berkeley Conference in Honor of Jerzy Neyman and Jack Kiefer*, ed. by L. M. Le Cam, and R. A. Olshen, vol. II. Wadsworth, Belmont, CA.
- DEHEJIA, R. H. (2005): “Program Evaluation as a Decision Problem,” *Journal of Econometrics*, 125, 141–173.
- DEL NEGRO, M., AND F. SCHORFHEIDE (2004): “Priors from General Equilibrium Models for VARs,” *International Economic Review*, 45(2), 643–673.
- DEMIRER, M., V. SYRGKANIS, G. LEWIS, AND V. CHERNUZHUKOV (2019): “Semi-Parametric Efficient Policy Learning with Continuous Actions,” working paper.
- DIEBOLD, F. X., T. A. GUNTHER, AND A. S. TAY (1998): “Evaluating Density Forecasts with Applications to Financial Risk Management,” *International Economic Review*, 39(4), 863–883.
- DOAN, T., R. LITTERMAN, AND C. SIMS (1984): “Forecasting and Conditional Projection using Realistic Prior Distributions,” *Econometric Reviews*, 3(1), 1–100.
- DOMINITZ, J., AND C. F. MANSKI (2017): “More Data or Better Data? A Statistical Decision Problem,” *Review of Economic Studies*, 84, 1583–1605.
- DUBÉ, J.-P., Z. FANG, N. FONG, AND X. LUO (2017): “Competitive Price Targeting with Smartphone Coupons,” *Marketing Science*, 36(6), 944–975.
- DUFLO, E., AND E. SAEZ (2003): “The Role of Information and Social Interactions in Retirement Plan Decisions: Evidence from a Randomized Experiment,” *Quarterly Journal of Economics*, 118, 815–842.
- EATON, M. L. (1989): *Group Invariance Applications in Statistics*, Regional Conference Series in Probability and Statistics. IMS-ASA.
- EGERTON, M. F., AND P. J. LAYCOCK (1982): “An Explicit Formula for the Risk of James-Stein Estimators,” *The Canadian Journal of Statistics*, 10(3), 199–205.
- EVANS, S., B. HANSEN, AND P. STARK (2005): “Minimax Expected Measure Confidence Sets for Restricted Location Parameters,” *Bernoulli*, 11(4), 571–590.

- FANG, Z. (2016): “Optimal Plug-in Estimators of Directionally Differentiable Functionals,” working paper.
- FANG, Z., AND A. SANTOS (2019): “Inference on Directionally Differentiable Functionals,” *The Review of Economic Studies*, 86(1), 377–412.
- FERGUSON, T. S. (1967): *Mathematical Statistics: A Decision Theoretic Approach*. Academic Press, New York.
- FESSLER, P., AND M. KASY (2019): “How to Use Economic Theory to Improve Estimators: Shrinking toward Theoretical Restrictions,” *The Review of Economics and Statistics*, 101(4), 681–698.
- FISHER, R. A. (1935): *The Design of Experiments*. Oliver and Boyd, Edinburgh.
- FRÖLICH, M. (2008): “Statistical Treatment Choice: An Application to Active Labor Market Programs,” *Journal of the American Statistical Association*, 103, 547–558.
- GELMAN, A., J. B. CARLIN, H. S. STERN, D. B. DUNSON, A. VEHTARI, AND D. B. RUBIN (2013): *Bayesian Data Analysis*. Chapman & Hall/CRC Press, New York, third edn.
- GEWEKE, J., G. KOOP, AND H. VAN DIJK (eds.) (2011): *The Oxford Handbook of Bayesian Econometrics*. Oxford University Press, New York.
- GILBOA, I., AND D. SCHMEIDLER (1989): “Maxmin Expected Utility with Non-unique Prior,” *Journal of Mathematical Economics*, 18, 141–153.
- GRAHAM, B. S., AND K. HIRANO (2011): “Robustness to Parametric Assumptions in Missing Data Models,” *American Economic Review: Papers & Proceedings*, 101(3), 538–543.
- GRAHAM, B. S., G. W. IMBENS, AND G. RIDDER (2014): “Complementarity and Aggregate Implications of Assortative Matching: A Nonparametric Analysis,” *Quantitative Economics*, 5, 29–66.
- HAHN, J. (1998): “On the Role of the Propensity Score in Efficient Semiparametric Estimation of Average Treatment Effects,” *Econometrica*, 66(2), 315–331.
- HAHN, J., K. HIRANO, AND D. KARLAN (2011): “Adaptive Experimental Design using the Propensity Score,” *Journal of Business and Economic Statistics*, 29(1), 96–108.
- HAILE, P., AND E. TAMER (2003): “Inference with an Incomplete Model of English Auctions,” *Journal of Political Economy*, 111(1), 1–51.
- HALL, S. G., AND J. MITCHELL (2004): “‘Optimal’ Combination of Density Forecasts,” working paper.
- HANSEN, B. E. (2015): “Shrinkage Efficiency Bounds,” *Econometric Theory*, 31, 860–879.
- (2016a): “Efficient Shrinkage in Parametric Models,” *Journal of Econometrics*, 190, 115–132.

- (2016b): “The Risk of James-Stein and Lasso Shrinkage,” *Econometric Reviews*, 35(8-10), 1456–1470.
- HAYASHI, T. (2008): “Regret Aversion and Opportunity Dependence,” *Journal of Economic Theory*, 139(1), 242–268.
- HIRANO, K., AND J. HAHN (2010): “Design of Randomized Experiments to Measure Social Interaction Effects,” *Economics Letters*, 106, 51–53.
- HIRANO, K., AND G. W. IMBENS (2004): “The Propensity Score with Continuous Treatments,” in *Applied Bayesian Modeling and Causal Inference from Incomplete Data Perspectives*, ed. by A. Gelman, and X.-L. Meng. Wiley, New York.
- HIRANO, K., G. W. IMBENS, AND G. RIDDER (2003): “Efficient Estimation of Average Treatment Effects Using the Estimated Propensity Score,” *Econometrica*, 71(4), 1161–1189.
- HIRANO, K., AND J. R. PORTER (2003a): “Asymptotic Efficiency in Parametric Structural Models with Parameter-Dependent Support,” *Econometrica*, 71(5), 1307–1338.
- (2003b): “Efficiency in Asymptotic Shift Experiments,” working paper.
- (2009): “Asymptotics for Statistical Treatment Rules,” *Econometrica*, 77(5), 1683–1701.
- (2012): “Impossibility Results for Nondifferentiable Functionals,” *Econometrica*, 80(4), 1769–1790.
- (2015): “Location Properties of Point Estimators in Linear Instrumental Variables and Related Models,” *Econometric Reviews*, 34(6-10), 720–733.
- HIRANO, K., AND J. H. WRIGHT (2017): “Forecasting with Model Uncertainty: Representations and Risk Reduction,” *Econometrica*, 85(2), 617–643.
- HOERL, A. E., AND R. W. KENNARD (1970): “Ridge Regression: Biased Estimation for Nonorthogonal Problems,” *Technometrics*, 12(1), 55–67.
- HU, F., AND W. F. ROSENBERGER (2006): *The Theory of Response-Adaptive Randomization in Clinical Trials*. Wiley, New York.
- HWANG, J. T. G., AND G. CASELLA (1982): “Minimax Confidence Sets for the Mean of a Multivariate Normal Distribution,” *The Annals of Statistics*, 10(3), 868–881.
- IBRAGIMOV, I., AND R. HASMINSKII (1981): *Statistical Estimation: Asymptotic Theory*. Springer-Verlag, New York.
- IMAI, K., AND D. A. VAN DYK (2004): “Causal Inference with General Treatment Regimes: Generalizing the Propensity Score,” *Journal of the American Statistical Association*, 99, 854–866.

- IMBENS, G. W. (2000): "The Role of the Propensity Score in Estimating Dose-Response Functions," *Biometrika*, 87, 706–710.
- INGRAM, B. F., AND C. H. WHITEMAN (1994): "Supplanting the 'Minnesota' Prior: Forecasting Macroeconomic Time Series using Real Business Cycle Model Priors," *Journal of Monetary Economics*, 34, 497–510.
- INOUE, A., AND L. KILIAN (2006): "On the Selection of Forecasting Models," *Journal of Econometrics*, 130, 273–306.
- JAMES, W., AND C. STEIN (1961): "Estimation with quadratic loss," *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability*, pp. 361–379.
- JEGANATHAN, P. (1982): "On the Asymptotic Theory of Estimation when the Limit of the Log-Likelihood Ratios is Mixed Normal," *Sankhya: A*, 44, 173–212.
- (1995): "Some Aspects of Asymptotic Theory with Applications to Time Series Models," *Econometric Theory*, 11, 818–887.
- KAIDO, H. (2016): "A Dual Approach to Inference for Partially Identified Econometric Models," *Journal of Econometrics*, 192(1), 269–290.
- KAIDO, H., AND A. SANTOS (2014): "Asymptotically Efficient Estimation of Models Defined by Convex Moment Inequalities," *Econometrica*, 82(1), 387–413.
- KAJI, T. (2017): "Theory of Weak Identification in Semiparametric Models," working paper.
- KALLUS, N., AND A. ZHOU (2018): "Policy Evaluation and Optimization with Continuous Treatments," working paper.
- KAN, R., AND G. ZHOU (2007): "Optimal Portfolio Choice with Parameter Uncertainty," *The Journal of Financial and Quantitative Analysis*, 42(3), 621–656.
- KANDEL, S., AND R. F. STAMBAUGH (1996): "On the Predictability of Stock Returns: An Asset-Allocation Perspective," *The Journal of Finance*, 51(2), 385–424.
- KASY, M., AND A. SAUTMANN (2019): "Adaptive Treatment Assignment in Experiments for Policy Choice," working paper.
- KIM, D.-H. (2010): "Bayesian Econometrics for Auction Models," Ph.D. thesis, University of Arizona.
- (2013): "Optimal Choice of a Reserve Price Under Uncertainty," *International Journal of Industrial Organization*, 31(5), 5887–602.
- (2015): "Flexible Bayesian Analysis of First Price Auctions Using a Simulated Likelihood," *Quantitative Economics*, 6(2), 429–461.

- KISH, L. (1965): *Survey Sampling*. Wiley, New York.
- KITAGAWA, T., AND A. TETENOV (2017): “Equality-Minded Treatment Choice,” working paper.
- (2018): “Who Should be Treated? Empirical Welfare Maximization Methods for Treatment Choice,” *Econometrica*, 86(2), 591–616.
- KLEIBERGEN, F. (2002): “Pivotal Statistics for Testing Structural Parameters in Instrumental Variables Regression,” *Econometrica*, 70(5), 1781–1803.
- KNIGHT, K., AND W. FU (2000): “Asymptotics for LASSO-Type Estimators,” *Annals of Statistics*, 28, 1356–1378.
- KNITTEL, C. R., AND S. STOLPER (2019): “Using Machine Learning to Target Treatment: The Case of Household Energy Use,” working paper.
- KOCK, A. B., AND M. THYRSGAARD (2017): “Optimal Sequential Treatment Allocation,” working paper.
- LANCASTER, T. (2004): *An Introduction to Modern Bayesian Econometrics*. Blackwell, Malden, MA.
- LE CAM, L. (1970): “On the Assumptions Used to Prove Asymptotic Normality of Maximum Likelihood Estimates,” *The Annals of Mathematical Statistics*, 41(3), 802–828.
- LE CAM, L. M. (1972): “Limits of Experiments,” in *Proceedings of the Sixth Berkeley Symposium of Mathematical Statistics*, vol. 1, pp. 245–261.
- (1986): *Asymptotic Methods in Statistical Theory*. Springer-Verlag, New York.
- LE CAM, L. M., AND G. L. YANG (2000): *Asymptotics in Statistics: Some Basic Concepts*. Springer-Verlag, 2nd edn.
- LECHNER, M. (2001): “Identification and Estimation of Causal Effects of Multiple Treatments under the Conditional Independence Assumption,” in *Econometric Evaluations of Active Labor Market Policies in Europe*, ed. by M. Lechner, and F. Pfeiffer. Physica, Heidelberg.
- LI, T., I. PERRIGNE, AND Q. VUONG (2002): “Structural Estimation of the Affiliated Private Value Model,” *The RAND Journal of Economics*, 33(2), 171–193.
- LIESE, F., AND K.-J. MIESCKE (2008): *Statistical Decision Theory: Estimation, Testing, and Selection*. Springer, New York.
- LIN, Y., R. MARTIN, AND M. YANG (2019): “On Optimal Designs for Nonregular Models,” *The Annals of Statistics*, 47(6), 3335–3359.
- MANSKI, C. F. (1995): *Identification Problems in the Social Sciences*. Harvard University Press.
- (2003): *Partial Identification of Probability Distributions*. Springer-Verlag.

- (2004): “Statistical Treatment Rules for Heterogeneous Populations,” *Econometrica*, 72(4), 1221–1246.
- (2007): *Identification for Prediction and Decision*. Harvard University Press.
- MANSKI, C. F. (2019): “Econometrics for Decision Making: Building Foundations Sketched by Haavelmo and Wald,” working paper.
- MANSKI, C. F., AND D. L. MCFADDEN (1981): “Alternative Estimators and Sampling Designs for Discrete Choice Analysis,” in *Structural Analysis of Discrete Data and Econometric Applications*, ed. by C. F. Manski, and D. L. McFadden. MIT Press, Cambridge, MA.
- MANSKI, C. F., AND A. TETENOV (2016): “Sufficient Trial Size to Inform Clinical Practice,” *Proceedings of the National Academy of Sciences*, 113(38), 10518–10523.
- MTAKOP, E., AND M. TABORD-MEEHAN (2016): “Model Selection for Treatment Choice: Penalized Welfare Maximization,” working paper.
- MIGUEL, E., AND M. KREMER (2003): “Worms: Identifying Impacts on Education and Health in the Presence of Treatment Externalities,” *Econometrica*, 72(1), 159–217.
- MILNOR, J. (1954): “Decisions Against Nature,” in *Decision Processes*, ed. by R. M. Thrall, C. H. Coombs, and R. L. Davis. Wiley, New York.
- MOLINARI, F. (2019): “Econometrics with Partial Identification,” this volume.
- MOREIRA, M. J. (2003): “A Conditional Likelihood Ratio Test for Structural Models,” *Econometrica*, 71(4), 1027–1048.
- MORGAN, K. L., AND D. RUBIN (2012): “Rerandomization to Improve Covariate Balance in Experiments,” *The Annals of Statistics*, 40(2), 1263–1282.
- MORGENSTERN, J., AND T. ROUGHGARDEN (2015): “The Pseudo-Dimension of Near-Optimal Auctions,” working paper.
- MORI, H. (2001): “Asymptotic Inadmissibility of Maximum Likelihood Estimator in a Quadratic Programming Problem,” *Gakushuin Economic Paper*, 38(1), 35–49.
- (2004): “Finite Sample Properties of Estimators for the Optimal Portfolio Weight,” *Journal of the Japan Statistical Society*, 34(1), 27–46.
- MORRIS, C. N. (1983): “Parametric Empirical Bayes Inference: Theory and Applications,” *Journal of the American Statistical Association*, 78(381), 47–55.
- MÜLLER, U. K. (2011): “Efficient Tests under a Weak Convergence Assumption,” *Econometrica*, 79(2), 395–435.

- MÜLLER, U. K., AND A. NORETS (2016): “Credibility of Confidence Sets in Nonstandard Econometric Problems,” *Econometrica*, 84(6), 2183–2213.
- MÜLLER, U. K., AND Y. WANG (2019): “Nearly Weighted Risk Minimal Unbiased Estimation,” *Journal of Econometrics*, 209, 18–34.
- MURPHY, S. A. (2003): “Optimal Dynamic Treatment Regimes,” *Journal of the Royal Statistical Society, Series B*, 65, 331–366.
- NARITA, Y. (2018): “Toward an Ethical Experiment,” working paper.
- NUSSBAUM, M. (1996): “Asymptotic Equivalence of Density Estimation and Gaussian White Noise,” *The Annals of Statistics*, 24(6), 2399–2430.
- PAARSCH, H. J. (1992): “Deciding between the common and private value paradigms in empirical models of auctions,” *Journal of Econometrics*, 51(1), 191–215.
- PFANZAGL, J. (1994): *Parametric Statistical Theory*. de Gruyter, New York.
- PHILIPSON, T. (1997): “Data Markets and the Production of Surveys,” *Review of Economic Studies*, 64, 47–72.
- PUKELSHEIM, F. (2006): *Optimal Design of Experiments*. Society for Industrial and Applied Mathematics, Philadelphia.
- RAMSEY, F. P. (1926): “Truth and Probability,” in *The Foundations of Mathematics and other Logical Essays*.
- RICE, K. M., T. LUMLEY, AND A. A. SZPIRO (2008): “Trading Bias for Precision: Decision Theory for Intervals and Sets,” UW Biostatistics Working Paper.
- ROBERT, C. (1988): “An Explicit Formula for the Risk of the Positive-Part James-Stein Estimator,” *The Canadian Journal of Statistics*, 16(2), 161–168.
- ROBINS, J. M. (2004): “Optimal Structural Nested Models for Optimal Decisions,” in *Proceedings of the Second Seattle Symposium on Biostatistics*, ed. by D. Y. Lin, and P. Heagerty, New York. Springer-Verlag.
- ROSSI, P. E., R. E. MCCULLOCH, AND G. M. ALLENBY (1996): “The Value of Purchase History in Target Marketing,” *Marketing Science*, 15(4), 321–340.
- SAVAGE, L. J. (1951): “The Theory of Statistical Decision,” *Journal of the American Statistical Association*, 46(253), 55–67.
- (1972): *The Foundations of Statistics*. Dover, New York, 2nd edn.
- SCHLAG, K. H. (2007): “Eleven - Designing Randomized Experiments under Minimax Regret,” working paper.

- SOLOMON, H., AND S. ZACKS (1970): "Optimal Design of Sampling from Finite Populations: A Critical Review and Indication of New Research Areas," *Journal of the American Statistical Association*, 65(330), 653–677.
- SONG, K. (2014): "Point Decisions for Interval-Identified Parameters," *Econometric Theory*, 30(2), 334–356.
- SPIEGELHALTER, D. J., L. S. FREEDMAN, AND M. K. B. PARMAR (1994): "Bayesian Approaches to Randomized Trials," *Journal of the Royal Statistical Society, Series A*, 157(3), 357–416.
- SPIESS, J. (2018): "Optimal Estimation when Researcher and Social Preferences are Misaligned," working paper.
- STAIGER, D., AND J. H. STOCK (1997): "Instrumental Variables Regression with Weak Instruments," *Econometrica*, 65(3), 557–586.
- STOCK, J. H., AND J. H. WRIGHT (2000): "GMM with Weak Identification," *Econometrica*, 68(5), 1055–1096.
- STOYE, J. (2009): "Minimax Regret Treatment Choice with Finite Samples," *Journal of Econometrics*, 151, 70–81.
- (2011): "Axioms for Minimax Regret Choice Correspondences," *Journal of Economic Theory*, 146(6), 2226–2251.
- (2012): "Minimax Regret Treatment Choice with Covariates or with Limited Validity of Experiments," *Journal of Econometrics*, 166(1), 138–156.
- STRASSER, H. (1985): *Mathematical Theory of Statistics*. Walter de Gruyter & Co., New York.
- SUKHATME, P. V. (1935): "Contribution to the Theory of the Representative Method," *Supplement to the Journal of the Royal Statistical Society*, 2(2), 253–268.
- TABORD-MEEHAN, M. (2018): "Stratification Trees for Adaptive Randomization in Randomized Controlled Trials," working paper.
- TAMER, E. (2012): "Partial Identification in Econometrics," *Annual Review of Economics*, 2, 167–195.
- TAY, A. S., AND K. F. WALLIS (2000): "Density Forecasting: A Survey," *Journal of Forecasting*, 19, 235–254.
- TETENOV, A. (2012): "Statistical Treatment Choice Based on Asymmetric Minmax Regret Criteria," *Journal of Econometrics*, 166, 157–165.
- (2016): "An Economic Theory of Statistical Testing," working paper.
- THOMPSON, S. K. (2012): *Sampling*. Wiley, New York, 3rd edn.

- TIBSHIRANI, R. (1996): “Regression Shrinkage and Selection via the Lasso,” *Journal of the Royal Statistical Society, Series B*, 58(1), 267–288.
- VAN DER VAART, A. W. (1991a): “An Asymptotic Representation Theorem,” *International Statistical Review*, 59, 99–121.
- (1991b): “On Differentiable Functionals,” *The Annals of Statistics*, 19, 178–204.
- (1998): *Asymptotic Statistics*. Cambridge University Press, New York.
- VAN DER VAART, A. W., AND J. A. WELLNER (1996): *Weak Convergence and Empirical Processes*. Springer-Verlag, New York.
- VARIAN, H. (1974): “A Bayesian Approach to Real Estate Assessment,” in *Studies in Bayesian Econometrics and Statistics in Honor of L. J. Savage*, ed. by S. E. Feiberg, and A. Zellner. North-Holland.
- WALD, A. (1950): *Statistical Decision Functions*. Wiley, New York.
- WESLER, O. (1959): “Invariance Theory and a Modified Minimax Principle,” *The Annals of Mathematical Statistics*, 30(1), 1–20.
- WINKLER, R. L. (1972): “A Decision-Theoretic Approach to Interval Estimation,” *Journal of the American Statistical Association*, 67(337), 187–191.
- XIE, X., S. C. KOU, AND L. D. BROWN (2012): “SURE Estimates for a Heteroscedastic Hierarchical Model,” *Journal of the American Statistical Association*, 107(500), 1465–1479.
- XIONG, R., S. ATHEY, M. BAYATI, AND G. IMBENS (2019): “Optimal Experimental Design for Staggered Rollouts,” working paper.
- ZAJONC, T. (2012): “Bayesian Inference for Dynamic Treatment Regimes: Mobility, Equity, and Efficiency in Student Tracking,” *Journal of the American Statistical Association*, 107(497), 80–92.
- ZELLNER, A. (1986): “Bayesian Estimation and Prediction Using Asymmetric Loss Functions,” *Journal of the American Statistical Association*, 81(394), 446–451.
- ZELLNER, A., AND V. K. CHETTY (1965): “Prediction and Decision Problems in Regression Models from the Bayesian Point of View,” *Journal of the American Statistical Association*, 60(310), 608–616.