

Wald's Statistical Decision Theory for Policy Analysis and Adaptive Experiments*

Keisuke Hirano[†]

October 31, 2025

Abstract

This paper reviews Wald's framework for statistical decision theory and discusses how it can be applied to empirical policy choice and adaptive experimental design in econometrics. The Wald framework provides an elegant and comprehensive approach to a wide array of empirical problems, but it can be difficult to solve statistical decision problems in realistic empirical settings. I provide some strategies for simplifying statistical decision problems to make them more tractable.

1 Introduction

Wald introduced his frequentist framework for statistical decision theory in a series of papers and monographs (Wald, 1939, 1945b, 1947a,b, 1949, 1950). It has greatly influenced the development of modern statistical and econometric theory (see, for example, Brown (2000)). Yet its direct impact on empirical economics was exceedingly modest in the decades after Wald's original contributions. Modern graduate-level textbooks in econometrics do not cover Wald's framework, or discuss it only briefly. So it may be surprising that it has seen renewed interest and attention in recent years, especially since Manski (2004) and Dehejia (2005) called attention to the potential for statistical decision theory to address empirical policy analysis. Further developments, including Stoye (2009), Tetenov (2012), Hirano and Porter (2009), Kitagawa and Tetenov (2018), and Athey and Wager (2021), developed new analytical results on statistical policy choice. The methodological literature on statistical policy choice is now quite active and intersects with the literatures in causal inference, identification analysis, and experimental design.

I argue that the renewed interest in statistical decision theory in econometrics arises in part from its potential to speak directly to a wide range of policy questions, and from an accumulation of technical innovations that have made frequentist analysis of statistical decision problems more tractable. However,

*This paper is based on a presentation in the semi-plenary session “Causal Inference and Statistical Decisions” at the 2025 World Congress of the Econometric Society. I thank Michael Gechter, Jack Porter, Jörg Stoye, and Han Xu for helpful comments and discussions. This research was supported by the U.S. National Science Foundation under grant SES-2117260. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author and do not necessarily reflect the views of the U.S. National Science Foundation.

[†]Department of Economics, Pennsylvania State University. Email: kuh237@psu.edu

many of these technical developments require relaxing or dropping some of the stringent requirements of Wald's framework in its pure form, with some inevitable tradeoffs in conceptual elegance and generality. By now, many variations of the Wald paradigm have been developed, often tailored to specific settings. How these alternative frameworks relate to each other is not always clear.

In this paper, I review Wald's framework and some of its variations and descendants, and provide a set of strategies for applying the approach to empirical policy choice problems. A key organizing principle is the role of dimension-reduction strategies that make statistical decision problems more tractable. These include the use of sufficiency and asymptotic sufficiency to reduce the dimensionality of the data, isolation of decision-relevant subparameters (or statistical functionals), and restriction of the action space to low-dimensional parametric classes.

To keep the discussion focused, I limit this review to methods that are closely related to the classic Wald framework. The subjective Bayesian approach (Savage, 1972), and other decision-theoretic frameworks that have emerged since Savage, provide an important alternative. The Wald framework can be connected to Bayesian approaches through average (Bayes) risk, as discussed in Section 2.4, but the emphasis here will be on frequentist evaluation of statistical procedures and connections to classical inference and estimation theory. For the sake of brevity, I also forgo discussing important recent developments in the empirical Bayes approach of Robbins (1956, 1985). Recent work on empirical Bayes methods for compound decision problems, including Gu and Koenker (2017), Armstrong, Kolesár, and Plagborg-Møller (2022), Montiel Olea, O'Flaherty, and Sethi (2021), Chen (2023), and Koenker and Gu (2024), has demonstrated its relevance for economics, but the topic deserves a more complete treatment than this article can provide.

Section 2 briefly reviews the classic Wald framework in static settings. The material will be familiar to readers already acquainted with statistical decision theory, but I emphasize aspects of the approach that make it relevant for policy analysis, and highlight some potential points of friction in its interpretation and application to empirical problems. Section 3 discusses dimension-reduction strategies, connecting concepts such as statistical sufficiency, local asymptotic limit experiments, and semiparametric efficiency for point estimation. Section 4 discusses Wald's approach for adaptive problems, especially the design and analysis of adaptive experiments.

2 The Basic Wald Framework

In its basic form, Wald's framework for statistical decision theory is remarkably simple and elegant, yet it encompasses a wide range of data-based decision-making problems. It consists of: (1) a statistical model for the data; (2) an action space and associated concept of a statistical decision rule; (3) a loss function and associated risk function; and (4) a menu of criteria for evaluating and solving the problem, such as minimaxity and average risk optimality. I briefly discuss each in turn, emphasizing some potential stumbling blocks in their interpretation, and identifying some practical challenges that arise in applying Wald's approach to concrete empirical problems.

2.1 Statistical Model

The data Z are assumed to follow $Z \sim P \in \mathcal{P}$. Here the notation Z indicates the entire data set. For example, if we have a sample of size n on (X_i, Y_i) , then $Z = \{(X_1, Y_1), \dots, (X_n, Y_n)\}$. This formulation is quite general, covering observational data sets as well as experimental studies, and encompassing cross-section, time-series and panel settings. When we need to be explicit about sample size, we can write Z^n for the data and P^n for its (joint) distribution. The set \mathcal{P} is the *statistical model*, the set of possible distributions for the data.

The notation $P \in \mathcal{P}$ encompasses both nonparametric and parametric settings. For example, we could take \mathcal{P} to be the set of all distributions on \mathbb{R} that have finite mean and variance. Or the statistical model \mathcal{P} could be a parametric class, for example $\mathcal{P} = \{N(\mu, \sigma^2)^n; \mu \in \mathbb{R}, \sigma^2 \in \mathbb{R}_+\}$. Then we could equivalently work with a parameter space Θ , such that $Z \sim P_\theta$ for some $\theta \in \Theta$. We will freely switch between these two notations for the statistical model in the sequel. It is important to note that the parameter θ needs to fully specify the distribution of the data. For example, in the linear regression specification $E[Y_i | X_i = x] = x'\beta$, the “parameter” β does not fully specify the joint distribution of (X_i, Y_i) or even the conditional distribution of Y_i given X_i . It is better regarded as a “subparameter” $\beta(\theta)$ of a more comprehensive parameter θ , or as a functional $\beta(P)$ of the joint distribution P of (X_i, Y_i) . Subparameters will be discussed further in Section 3.3.

Nothing in this setup requires that the statistical model satisfy any specific regularity or smoothness conditions, nor that the parameter θ , or decision-relevant subparameters β , be point-identified (see e.g. Hirano and Porter (2020) and Manski (2024)). Such “nonstandard” statistical setups can lead to substantial technical complications in solving Wald’s program.

2.2 Action Space and Decision Rules

The *action space* \mathcal{A} is the set of choices available to the decision-maker. A *statistical decision rule* $\delta(\cdot)$ maps data $Z = z$ to an action $a = \delta(z)$, or a probability distribution over actions (in which case we call it a randomized statistical decision rule). In many problems, there is no gain from allowing randomization (i.e., optimal rules are nonrandomized), but working with the larger class of randomized rules convexifies the set of decision rules, making the problem easier to analyze. However, some recent work including Manski (2024) and Montiel Olea, Qiu, and Stoye (2025) has shown that in certain statistical treatment choice problems involving partial identification, minimax regret rules are randomized. A randomized rule can usually be specified as a function $\delta(Z, U)$ of Z and a randomization variable $U \sim \text{Unif}[0, 1]$ which is independent of Z .¹

With this component, the Wald framework encompasses and generalizes the Neyman-Pearson paradigm. If the goal is to produce a point estimate of the parameter $\theta \in \Theta$, we can take $\mathcal{A} = \Theta$, whereas in hypothesis testing we might take $\mathcal{A} = \{0, 1\}$ with the action 1 corresponding to rejecting the null hypothesis.

¹If the action space \mathcal{A} is a subset of finite-dimensional Euclidean space, or more generally of a Borel space, any rule δ that maps Z into probability distributions over \mathcal{A} can be represented in this way.

We are not limited to point estimation and statistical inference. Since the action space can be discrete, continuous, or even a space of functions, this framework enables rigorous evaluations of empirically driven policy choices in a variety of settings.

Example 1 (Binary Policy Choice). In the case where the set of possible policies (or treatments) is binary, the action space can be taken as $\{0, 1\}$, the same action space as in hypothesis testing. The actions $a = 0, 1$ are interpreted as possible policies or treatments that will be assigned to all future individuals in some population. The loss function with which we evaluate this action, to be discussed below, represents the negative welfare associated with an action for some particular population. A statistical treatment rule (like a hypothesis test) specifies how this binary choice will be determined on the basis of the data, which might arise from a randomized controlled trial or from observational data.

□

Example 2 (Conditional Policy Rules). Manski (2004) and Dehejia (2005) imagined a policy-maker choosing policy rules $a : \mathcal{X} \rightarrow \mathcal{T}$ that assign (future) individuals to treatment arms $t \in \mathcal{T}$ on the basis of individual characteristics $X \in \mathcal{X}$. In medicine, such policies are often called “personalized” or “precision” medicine. Absent further restrictions, the action space \mathcal{A} consists of all such functions. Kitagawa and Tetenov (2018) argued that in practice, policymakers would restrict the set of possible policies to a relatively simple class. We discuss this further in Section 3.4.

□

Example 3 (Counterfactual Policy Analysis). In structural econometric analyses that estimate models of economic primitives and agent behavior, it has become common to report policy counterfactuals that predict the effects of changes in economic policies. Here we can think of the action space as the menu of policies under consideration. These counterfactual policy predictions are usually constructed by evaluating or simulating the economic model at its estimated parameter values, and rarely take explicit account of estimation error and the implications of parameter uncertainty for policy conclusions. Some exceptions include Aryal and Kim (2013) and Kim (2013), who study auction design based on empirical data using decision-theoretic methods.

□

2.3 Loss and Risk Functions

The real-valued *loss function* $L(P, a)$ quantifies the loss from choosing action a when the true distribution is P . Alternatively, if we evaluate the action using a utility or welfare function $W(P, a)$, we can take loss to be its negative. The *regret loss* considers loss or welfare relative to the best action, and can be defined as $L_R(P, a) = L(P, a) - \min_{\tilde{a}} L(P, \tilde{a})$ or $L_R(P, a) = \max_{\tilde{a}} W(P, \tilde{a}) - W(P, a)$.

Recall that a statistical decision rule δ specifies, *ex ante*, how data will be turned into actions. The *risk function* evaluates the expected loss when actions are chosen by a statistical decision rule $\delta(Z)$:

$$R(P, \delta) = \mathbb{E}_P [L(P, \delta(Z))] = \int L(P, \delta(z)) dP(z),$$

with a straightforward extension to the case of a randomized decision rule. Since risk is based on the sampling distribution of the data, it implicitly accounts for estimation uncertainty.

The framework requires the decision-maker to commit to a loss or welfare measure. Solutions may be sensitive to the choice of loss or welfare function, though in well-established problems there may be some reasonable default choices.

Example 4 (Loss Functions for Point Estimation). For estimation problems it is common to take $L(\theta, a) = \|\theta - a\|^2$, or some other convex, symmetric function. Note that this loss function is already in regret form: $L_R(\theta, a) = L(\theta, a)$.

□

Example 5 (Utilitarian Welfare). Suppose there is a binary treatment. The target population P consists of individuals with potential outcomes $Y(0), Y(1)$ and observable characteristics $X \in \mathcal{X} \subset \mathbb{R}^{d_x}$. A treatment policy $a : \mathcal{X} \rightarrow \{0, 1\}$ assigns treatments based on characteristics, as in Example 2. The utilitarian welfare of a policy a is

$$W(P, a) = \mathbb{E}_P [Y(1) \cdot a(X) + Y(0) \cdot (1 - a(X))].$$

For a class of feasible rules \mathcal{A} , let $a^*(P) = \arg \max_{a \in \mathcal{A}} W(P, a)$. Then the welfare regret loss can be defined as

$$L_R(P, a) = W(P, a^*(P)) - W(P, a).$$

Alternative social welfare functions for evaluating treatment policies have been considered by Kasy (2016), Kitagawa and Tetenov (2021), Kock, Preinerstorfer, and Veliyev (2023), Escanciano and Terschuur (2023), and Terschuur (2025), among others.

□

Predictive Loss

The terminology of loss and risk is used throughout econometrics and statistics, but with different meanings that are easily conflated. Another common usage in forecasting and machine learning defines loss in terms of a future *outcome* rather than in terms of P (or θ), e.g. $L(y_{n+1}, \hat{f}(x_{n+1}))$ where $\hat{f}(x_{n+1})$ is the predictor of y_{n+1} based on x_{n+1} and an estimated function \hat{f} . This predictive loss concept is used, for example, in empirical risk minimization methods (Vapnik, 1991).

Defining utilities in terms of outcomes is also typical in axiomatic decision theory. One could take Waldean loss to be the expected value (under P) of predictive loss, or of the negative of outcome-based utility, but this extra step has substantive implications.

Suppose that a policy-maker has a utility function $u(y)$ over outcomes y . In the setup of Example 5 we could take $W(P, a)$ to be the expected utility

$$W(P, a) = \mathbb{E}_P [u(Y(1))a(X) + u(Y(0))(1 - a(X))] = \mathbb{E}_P [u(Y(a(X)))].$$

If the policy-maker's preferences over outcomes are not of the expected utility form, however, the definition of $W(P, a)$ should be modified. For example, with quantile preferences as in Manski (1988) and Rostek (2009), one can define $W(P, a)$ to be the appropriate quantile of the distribution of Y induced by P and a ; see for example Guggenberger, Mehta, and Pavlov (2024). In general, quantile or inequality-averse preferences should be incorporated into the definition of the Waldean loss or welfare function, even before accounting for statistical uncertainty through a risk function. Other applications and discussion of ordinal and quantile utility include Manski and Tetenov (2023) and Kitagawa, Lee, and Qiu (2025).

2.4 Evaluation and Solution Criteria

Having specified the statistical model, the set of decision rules under consideration, and their risk functions, it remains to formulate a criterion with which to compare rules, and then solve the problem by finding a rule that optimizes the criterion.

One criterion for evaluating statistical decision rules is admissibility. A decision rule is admissible if no other decision rule dominates it in terms of risk $R(P, \delta)$ over all $P \in \mathcal{P}$. While admissibility is intuitively appealing, it is not always a useful concept. Montiel Olea, Qiu, and Stoye (2025) show that, in a Gaussian policy choice setting with partial identification, *all* decision rules are admissible. Meanwhile, in the classical normal linear regression model the ordinary least squares estimator is generally inadmissible, a result that dates back to Stein (1956).

A popular solution concept is minimaxity. Minimax rules solve the functional optimization problem

$$\inf_{\delta} \sup_{P \in \mathcal{P}} R(P, \delta).$$

Often the loss function is taken to be regret loss, leading to the minimax regret optimality criterion. Savage (1951) is often viewed as the originator of the minimax regret concept, but he himself credits Wald (Savage (1972), p. 170.) Savage (1972) also gives a group decision-making interpretation to minimax regret; see Armstrong, Kitagawa, and Tetenov (2024) for a lucid discussion of Savage's interpretation and connections to recent work on interpersonal extensions of statistical decision theory.

Another common criterion is weighted average (or Bayes) risk. Here it is notationally convenient to assume that \mathcal{P} can be written as a parametric family $\mathcal{P} = \{P_\theta : \theta \in \Theta\}$. We do not require that θ be finite-dimensional. Let $\Pi(\theta)$ be a probability distribution, or more generally a measure, on Θ . The Bayes risk averages frequentist risk over θ :

$$B(\Pi, \delta) = \int R(P_\theta, \delta) d\Pi(\theta).$$

For a given Π , the Bayes-risk optimal rule minimizes $B(\Pi, \delta)$. This criterion is also used in the hypothesis testing literature in the guise of (local) weighted average power optimality. Wald considered the Bayes risk criterion as a valid criterion in its own right, and as a technical tool to reduce the set of rules under consideration. Complete-class theorems link the class of Bayes-risk optimal rules to the class of admissible rules, and one can sometimes solve for minimax rules by finding a “least favorable” prior Π .

The Bayes risk criterion requires a choice for the prior Π . Since a different prior could lead to different prescriptions for the decision rule, this raises a burden in producing empirical analyses that are broadly convincing. The minimax criterion is attractive because it frees the user from having to choose a prior. But in some applications it can be quite pessimistic, focusing attention on extreme or peculiar points in the parameter space. Adopting the minimax (or minimax regret) criterion does not entirely circumvent the problem.

Wald viewed the minimax and minimax regret criteria as being on less firm ground than other aspects of his approach. Wolfowitz (1952), in an article written shortly after Wald's untimely death in 1950, wrote: "Wald was searching for other criteria, and his last joint work with this writer concerned this problem. He was dissatisfied with known results on the problem and had no great faith in the necessity for the minimax criterion." See Brown (1994) for additional historical discussion of the minimax criterion.

3 Dimension Reduction Strategies

Fully solving statistical decision problems in the Wald paradigm of Section 2 can be daunting. Even if the action space is simple (e.g. binary as in the simple treatment choice case), the set of possible statistical decision rules consists of *functions* from the sample space of Z to \mathcal{A} , and it is over this set that one optimizes to find minimax and average risk optimal rules. For this reason, it is noteworthy that Stoye (2009) and Schlag (2006) were able to find the exact minimax statistical treatment rule in an experimental setting where the underlying statistical model is nonparametric.

This section discusses strategies to simplify statistical decision problems to make them more tractable, focusing on the special role of *dimension reduction*. I consider four such strategies: reduction by sufficiency in Section 3.1; approximate (asymptotic) sufficiency in Section 3.2; subparameters in Section 3.3; and parametrized policy classes in 3.4. Applying these strategies may require relaxing or moving outside of the strict finite-sample Wald framework. I discuss some of these tradeoffs, with the goal of maintaining the conceptual clarity of Wald's framework while enabling practical applications.

A number of computational strategies for finding minimax and minimax regret rules have been explored in the literature, including Chamberlain (2000), Manski and Tabord-Meehan (2017), Litvin and Manski (2021), Masten (2023), Guggenberger and Huang (2025), and Aradillas Fernandez, Blanchet, Montiel Olea, Qiu, Stoye, and Tan (2024). There is also a well-developed literature on numerical methods for calculating Bayesian posteriors and finding Bayes-optimal rules. Many of these methods are more effective when first employing some of the dimension-reduction strategies discussed below to reduce the problem complexity.

3.1 Sufficiency

Recall that a statistic $S = S(Z)$ is sufficient² for \mathcal{P} in the classical (Fisher-Neyman) sense if the conditional distribution of Z given S does not depend on P . This notion of sufficiency can be checked using the Neyman factorization criterion. For our purposes an alternative notion of sufficiency, Blackwell sufficiency (Blackwell, 1951), is especially useful. In general, a statistic that is sufficient in the classical sense need not be Blackwell sufficient (see Bahadur (1954)), but under a domination condition that is commonly satisfied in empirical applications, the two concepts are equivalent. The following result gives a useful decision-theoretic implication of sufficiency. It was first shown by Bahadur (1954), building on Halmos and Savage (1949). In the form stated below, it can be deduced from Proposition 4.58 in Liese and Miescke (2008).

Theorem 1 (Bahadur, 1954). *Suppose that the model \mathcal{P} is dominated by a σ -finite measure and that the sample space \mathcal{Z} is a Borel space. Let $S = S(Z)$ be a sufficient statistic for the model \mathcal{P} . Consider any decision rule $\delta(Z)$ from \mathcal{Z} to a subset of a Euclidean space. Then there exists a randomized decision rule $\tilde{\delta}(S, U)$ based on S and an independent randomization $U \sim \text{Unif}[0, 1]$ such that*

$$\delta(Z) \sim \tilde{\delta}(S, U) \quad \forall P \in \mathcal{P}.$$

This result states that given any statistical decision rule, we can achieve the same distributions over actions, and hence the same risk function, by using some randomized rule based on a sufficient statistic. The condition that the model is dominated is satisfied if the model can be expressed as a collection of probability density functions with respect to Lebesgue measure or some other well-behaved measure, so it is often met in empirical applications.

The class of statistical decision rules based on sufficient statistics can be much simpler than the class of all decision rules, as the following example illustrates.

Example 6 (Multivariate Normal Location Model). Let $Z = (Z_1, \dots, Z_n)$ where the Z_i are k -dimensional vectors distributed iid $N(\theta, \Sigma)$. The mean $\theta \in \mathbb{R}^k$ is unknown. Assume that the variance matrix Σ is known. Then $\hat{\theta}_n = \frac{1}{n} \sum_{i=1}^n Z_i$ is a k -dimensional sufficient statistic. For any decision rule $\delta(Z_1, \dots, Z_n)$ taking values in a subset of a Euclidean space, there is a rule of the form $\tilde{\delta}(\hat{\theta}_n, U)$, where $U \sim \text{Unif}[0, 1]$, such that

$$\delta(Z_1, \dots, Z_n) \sim \tilde{\delta}(\hat{\theta}_n, U) \quad \forall \theta.$$

□

In the setup of Example 6, suppose we have a simple binary action space (as in Example 1 with a binary policy choice). Then the set of all decision rules consists of all functions from $\mathbb{R}^{k \times n}$ to $\{0, 1\}$. Even if we ignore the ordering of the observations Z_1, \dots, Z_n , this is a huge space over which to search for the optimal rule. Theorem 1 lets us reduce the search to a much simpler class of rules, those depending only on $\hat{\theta}_n$ (and possibly some independent randomization). Since $\hat{\theta}_n$ itself has a normal distribution with

²See Liese and Miescke (2008), Ch. 4, and Pfanzagl (2017), Ch. 2, for more extensive discussions of sufficiency.

known variance, further simplifications are often possible, depending on the form of the loss function and the solution criterion.

3.2 Approximate Sufficiency

Unfortunately, useful low-dimensional sufficient statistics are available only for a small number of statistical models (most prominently for exponential family models). However, a type of asymptotic or approximate sufficiency holds for a much broader class of statistical models.

Consider a parametric model with parameter $\theta \in \Theta \subset \mathbb{R}^k$, where k is finite. Under conventional regularity conditions, the maximum likelihood estimator $\hat{\theta}_{\text{ML}}$ is approximately normally distributed,

$$\hat{\theta}_{\text{ML}} \xrightarrow{\text{approx}} N(\theta, \Sigma/n),$$

where Σ equals the inverse of the Fisher information matrix and can be consistently estimated. The similarity with the sufficient statistic appearing in Example 6 seems promising. But the MLE is *not* a sufficient statistic in general, and the fact that its variance shrinks to zero as the sample size increases means that, in the limit, the statistical uncertainty about θ vanishes as sample size increases.

An approximate sufficiency concept can be developed, however, by employing the limits of experiments technique.³ The most common form uses local sequences of parameters. Fix some $\theta_0 \in \Theta$, and consider sequences of parameter values of the form $\theta_0 + h/\sqrt{n}$, where $h \in \mathbb{R}^k$ and n is the sample size. This replaces the original k -dimensional parameter space Θ by a local parameter space of vectors h that define local deviations from a point θ_0 .⁴ The idea is to consider the performance of statistical decision rules under alternative “nearby” parameters such that the statistical uncertainty about θ does not vanish in the limit.

If the parametric model is smooth and has a nonsingular Fisher information matrix whose inverse is Σ , we have the following result.⁵

Theorem 2. *Consider any sequence of statistical decision rules $\delta_n(Z^n)$ that possess limits under every local parameter sequence. Then there exists a randomized decision rule $\tilde{\delta}(Z, U)$ based on the normal model $Z \sim N(h, \Sigma)$ with $U \sim \text{Unif}[0, 1]$ independent of Z , such that*

$$\delta_n(Z^n) \xrightarrow{d} \tilde{\delta}(Z, U) \quad \text{under } \theta_0 + h/\sqrt{n}, \text{ for all } h \in \mathbb{R}^k. \quad (1)$$

Note the similarity with Theorem 1 and Example 6. The maximum likelihood estimator $\hat{\theta}_{\text{ML}}$ typically satisfies

$$\sqrt{n}(\hat{\theta}_{\text{ML}} - \theta_0) \xrightarrow{d} Z \sim N(h, \Sigma)$$

³The “experiment” in “limits of experiments” follows the textbook usage in statistics, referring to a situation where data are drawn from some well-defined probability distribution. This concept applies to any well-specified statistical models, and is not restricted to randomized controlled trials.

⁴Local parameter sequences are sometimes referred to as Pitman drifts, but were first proposed by Neyman (1937); see McManus (1991).

⁵For a more precise statement of this result, see for example van der Vaart (1998), Ch. 7.

under every local parameter sequence $\theta_0 + h/\sqrt{n}$, which can be interpreted as saying that the (normalized) MLE is approximately sufficient for decision-theoretic purposes. This approximate sufficiency of the MLE was first shown by Wald, in an analysis of local asymptotic optimality in testing problems (Wald, 1943). In the decades following, the idea of approximating a statistical decision problem by one involving a simpler “limit experiment” was developed in great generality by Le Cam (1972, 1986), Hájek (1970), Ibragimov and Hasminskii (1981), van der Vaart (1991), and others. The locally asymptotically normal (LAN) case of Theorem 2 is encountered most frequently in applications to parametric statistical models, but in some other models that do not satisfy the conditions required for this case, alternative asymptotic sufficiency results have been obtained. An especially important extension to semiparametric models is discussed further in Section 3.3.

The dimension reduction achieved by local asymptotic sufficiency is remarkable, but takes the analysis outside the pure finite-sample Wald framework outlined in Section 2. Not only are the results only approximate, but local reparametrization may fundamentally change the interpretation of the resulting analysis. For example, the local asymptotic minimax criterion considers the worst-case risk of a procedure over a local (i.e., shrinking) neighborhood, which can lead to different conclusions than the global minimax criterion in Section 2.4. Moreover, to obtain a nontrivial analysis, the risk functions must be normalized so that they do not diverge as $n \rightarrow \infty$. In practice, this often requires that the loss function itself be centered at zero, which forces the use of regret loss or welfare regret. On the other hand, it could be argued that considering worst-case performance over a local neighborhood where the statistical uncertainty is nonvanishing, and working with regret loss, is no more unreasonable than the global minimax criterion.

Other challenges in interpretation can arise as well, such as in settings where treatment effects are partially identified. Manski (2009, 2010, 2011) initiated the study of treatment assignment under partial identification, and it has been an active area of methodological research in recent years. Following Yata (2023), suppose there is a structural parameter of interest $\theta \in \Theta$, where Θ is a convex set in some Euclidean space. There is a mapping $m : \Theta \rightarrow \mathbb{R}^k$ that yields the reduced-form parameter $\mu = m(\theta)$. The mapping is not assumed to be injective, so for a given value of the reduced-form parameter μ , the identified set of values for θ is

$$\Theta(\mu) := \{\theta \in \Theta : m(\theta) = \mu\}.$$

We observe data Z that are informative about the reduced-form parameter μ . Yata (2023), Montiel Olea, Qiu, and Stoye (2025), and other recent work consider the situation where

$$Z \sim N(m(\theta), \Sigma)$$

where Σ is known, and derive novel finite-sample optimality results for certain decision problems.

Given the shifted normal form for the reduced-form model, it is tempting to think of this statistical setup as corresponding to the limit experiment in the locally asymptotically normal case, so that these results can be interpreted as simultaneously providing asymptotic optimality results when the reduced-form

model is regular and parametric. However, making this connection is somewhat complicated. Applying the normal limit experiment theory would require centering the reduced-form parameter μ at some point μ_0 , and considering local sequences of the form $\mu_0 + h/\sqrt{n}$ for $h \in \mathbb{R}^k$. The corresponding sequences of identified sets are

$$\tilde{\Theta}_n(\mu_0; h) = \{\theta \in \Theta : m(\theta) = \mu_0 + h/\sqrt{n}\}.$$

As $n \rightarrow \infty$ and for any fixed μ_0 , the structure of $\tilde{\Theta}_n(\mu_0; h)$ as h ranges over \mathbb{R}^k may become quite different from the structure of $\Theta(\mu)$ as μ ranges over its parameter space. Thus, the localization inherent in large-sample distributional approximations is not innocuous in this class of problems. One way around this is through profiling as in Christensen, Moon, and Schorfheide (2023); see also Kido (2023) and Xu (2025).

An alternative way to reduce the dimensionality of the data is to impose a restriction at the outset, that the statistical decisions are made on the basis of some function of the data such as the sample mean or the maximum likelihood estimator. Müller (2011) develops a notion of asymptotic optimality of decision rules based on a user-specified converging sequence of statistics. Like the local asymptotic sufficiency theory outlined above, this approach typically requires the use of local parameter sequences, to ensure that the statistical decision problem remains nondegenerate in the limit.

3.3 Subparameters

Recall that the parameter θ must fully characterize the distribution of the data. In many empirical problems, this requires that the statistical model \mathcal{P} be moderately to extremely high dimensional. Even a problem as basic as estimating the linear regression model $Y_i = \beta' X_i + \epsilon_i$, without parametric restrictions on the distribution of ϵ_i , may involve an infinite-dimensional model space. But our interest may focus on a low-dimensional subparameter or functional of the unknown distribution, such as β in this case.

Example 7 (Subparameter of a Multivariate Normal Model). Consider the multivariate normal case of Example 6 where $Z \sim N(\theta, \Sigma)$. Assume that Σ is full rank. Suppose we are interested in a d -dimensional linear transform $\beta := G\theta$ where $G \in \mathbb{R}^{d \times k}$ is known and $d < k$. There exists a full-rank $k \times k$ matrix M with G as its upper block, such that

$$MZ \sim N\left(\begin{pmatrix} \beta \\ \gamma \end{pmatrix}, \begin{pmatrix} G\Sigma G' & 0 \\ 0 & \tilde{\Sigma}_\gamma \end{pmatrix}\right). \quad (2)$$

□

In this example, the statistical experiment of observing MZ is equivalent⁶ to observing Z , and the first d -dimensional subvector of MZ is distributed as $[MZ]_1 = GZ \sim N(\beta, G\Sigma G')$, independently of the other components. A natural question is: when can one restrict attention to decision rules based on the d -dimensional statistic GZ ?

⁶More precisely, the mapping $Z \mapsto MZ$ is one-to-one and measurable, so MZ is sufficient for Z .

Suppose that the loss function depends on θ only through $G\theta$. For Bayes risk optimality, there is an optimal rule that depends on Z only through GZ if the prior Π over $M\theta = (\beta, \gamma)$ can be written in a product form:

$$\Pi(\beta, \gamma) = \Pi_\beta(\beta) \times \Pi_\gamma(\gamma).$$

In other words, the prior for β is independent of the prior for the other components of the parameter. If this is not the case, however, it is not necessarily justified to restrict attention to rules based on GZ .

For minimax optimality, it can often be shown that there is a minimax rule that depends only on GZ , though this should be verified for the problem at hand. One strategy that can be useful is to carve the full parameter space into one-dimensional “slices” within which the orthogonalized nuisance component γ is fixed. If the loss function is invariant to the nuisance component and if we can find a minimax decision rule based on GZ , then it is also minimax over the full parameter space. The following lemma gives a simple result that facilitates this approach.

Lemma 1 (Minimax via Orthogonal Slicing). *Let $Z \sim P_\theta$ for $\theta \in \Theta$, where $\theta = (\beta, \gamma) \in B \times \Gamma$. Suppose Z can be partitioned as (Z_1, Z_2) where*

$$Z_1 \sim P_{1,\beta}, \quad Z_2 \sim P_{2,\gamma}, \quad Z_1 \perp Z_2.$$

Suppose the loss function $L(\theta, a)$ depends on θ only through β : $L(\theta, a) = L((\beta, \gamma), a) = \tilde{L}(\beta, a)$. Suppose there exists a decision rule $d^(Z_1)$ based on the first component of Z that is minimax for some fixed γ^* , i.e. for any $\delta: \mathcal{Z} \rightarrow \mathcal{A}$,*

$$\sup_{\beta} \mathbb{E}_{\beta, \gamma^*} [\tilde{L}(\beta, d^*(Z_1))] \leq \sup_{\beta} \mathbb{E}_{\beta, \gamma^*} [\tilde{L}(\beta, \delta(Z))].$$

Then d^ is a minimax statistical decision rule when both β and γ are unknown.*

Proof:

By construction, the risk of d^* does not depend on γ , hence

$$\sup_{\beta, \gamma} \mathbb{E}_{\beta, \gamma} [\tilde{L}(\beta, d^*(Z_1))] = \sup_{\gamma} \sup_{\beta} \mathbb{E}_{\beta, \gamma} [\tilde{L}(\beta, d^*(Z_1))] = \sup_{\beta} \mathbb{E}_{\beta, \gamma^*} [\tilde{L}(\beta, d^*(Z_1))].$$

If d^* is not minimax over (β, γ) , then there would exist another decision rule $\delta(Z)$ such that

$$\sup_{\beta} \mathbb{E}_{\beta, \gamma^*} [\tilde{L}(\beta, d^*(Z_1))] = \sup_{\beta, \gamma} \mathbb{E}_{\beta, \gamma} [\tilde{L}(\beta, d^*(Z_1))] > \sup_{\beta, \gamma} \mathbb{E}_{\beta, \gamma} [\tilde{L}(\beta, \delta(Z))] \geq \sup_{\beta} \mathbb{E}_{\beta, \gamma^*} [\tilde{L}(\beta, \delta(Z))],$$

which contradicts the assumption.

□

Although fairly straightforward, the lemma can be useful to solve minimax problems in subparameter settings after an orthogonal reparametrization, especially when loss depends on the full parameter vector through a scalar or low-dimensional subparameter. It generalizes the argument used in the proof of Hirano and Porter (2009), Theorem 3.4. Its proof is adapted from a result for point estimators given in Lehmann and Casella (1998), Chapter 5, Lemma 1.15, which was used in Xu (2025) to justify subpa-

rameter restrictions. Lemma 1 can also be regarded as a special case of the more general concept of “invariant sufficiency” developed in Hall, Wijsman, and Ghosh (1965), but the particular form given here is adequate for our purposes. See Choi, Hall, and Schick (1996) for a related argument for the case of hypothesis testing.

The result requires that there exists a minimax decision rule in the sliced model that only uses Z_1 . Note that in the sliced model, for any fixed γ , the statistic Z_1 is sufficient for β . Thus the condition always holds if the statistical model is dominated (as it is in Example 7) and we allow for randomization, by Theorem 1. Moreover, if the loss function is convex, then by the Rao-Blackwell theorem, any decision rule depending on the full vector Z is weakly dominated by its conditional expectation given the sufficient statistic. Then there exists a nonrandomized minimax decision rule depending only on Z_1 in the sliced model.

Focusing on a subparameter combines well with local asymptotic approximations of the type discussed in Section 3.2, when the subparameter is smooth.

Example 8 (Asymptotic Subparameter of a Parametric Model). Suppose that the parametric model $\{P_\theta : \theta \in \Theta\}$ with k -dimensional parameter θ satisfies the regularity conditions for Theorem 2 for local parameter $\theta_0 + h/\sqrt{n}$. Let $\beta : \Theta \rightarrow \mathbb{R}^d$, with $d < k$, be continuously differentiable at θ_0 , so that

$$\beta(\theta_0 + h/\sqrt{n}) = \beta(\theta_0) + \frac{1}{\sqrt{n}}Gh + o\left(\frac{1}{n}\right)$$

where $G = \frac{\partial}{\partial \theta} \beta(\theta_0)$ is the $d \times k$ matrix of partial derivatives of $\beta(\cdot)$ at θ_0 . Suppose the maximum likelihood estimator $\hat{\theta}_{ML}$ satisfies

$$\sqrt{n}(\hat{\theta}_{ML} - \theta_0) \xrightarrow{d} N(h, \Sigma) \quad \forall h,$$

where Σ denotes the inverse of the Fisher information matrix at θ_0 . Then by the delta method, the plug-in estimator of $\beta(\theta)$ will satisfy

$$\sqrt{n}(\beta(\hat{\theta}_{ML}) - \beta(\theta_0)) \xrightarrow{d} N(Gh, G\Sigma G').$$

□

In this example, the local asymptotic sufficiency result in Theorem 2 gives a limiting multivariate normal location model with local parameter vector $h \in \mathbb{R}^k$. We can then apply the subparameter orthogonalization of Example 7. Put simply, we can restrict attention to decision rules based on the natural plug-in estimator $\hat{\beta} = \beta(\hat{\theta}_{ML})$ with its approximate normal distribution, and seek to find a rule $\delta(\hat{\beta})$ that has good risk properties.

This strategy further extends to semiparametric models, where the parameter of interest can be estimated efficiently in the semiparametric sense.

Example 9 (Semiparametric Functionals). Suppose the model \mathcal{P} is nonparametric, for example consisting of all distributions on some sample space that possess a certain number of moments. Let $\beta(P)$ be a (scalar) functional, such as the population expectation $\beta(P) = \mathbb{E}_P[Z]$, or the average treatment effect under unconfoundedness $\beta(P) = \mathbb{E}_P[\tau(X)]$ where $\tau(x) := \mathbb{E}_P[Y|D=1, X=x] - \mathbb{E}_P[Y|D=0, X=x]$.

If $\beta(P)$ is sufficiently smooth as a function of P , and if \mathcal{P} is sufficiently well-behaved, then $\beta(P)$ can be estimated at a parametric rate, and a similar analysis to Example 8 applies. In particular, one can define a local parametrization by fixing $P_0 \in \mathcal{P}$ and considering local deviations in the form of smooth one-dimensional parametric submodels. The score functions of such parametric submodels form a local parameter space, called the tangent space at P_0 . The efficient influence function represents the derivative of $\beta(P)$ with respect to P at P_0 , analogously to the derivative G in Example 8.

An asymptotic sufficiency result for smooth semiparametric models was shown by van der Vaart (1991). Analogously to Theorem 2, any converging sequence of decision rules δ_n can be represented by some rule $\tilde{\delta}(Z, U)$, where Z is a Gaussian sequence model. This can be orthogonalized (see Hirano and Porter (2009)), such that $Z = (Z_1, Z_2)$, where

$$Z_1 \sim N(h_\beta, \sigma_\beta^2),$$

and Z_2 is a Gaussian process independent of Z_1 . Here h_β represents the local (sub)parameter $h_\beta = \sqrt{n}(\beta(P_n) - \beta(P_0))$ for $\beta(P)$ corresponding to local sequences of distributions $P_n \rightarrow P_0$, and σ_β^2 equals the semiparametric efficiency bound for estimating $\beta(P)$. Semiparametrically efficient point estimators $\hat{\beta}$ of $\beta(P)$ satisfy

$$\sqrt{n}(\hat{\beta} - \beta(P_0)) \xrightarrow{d} Z_1$$

under local sequences P_n . Thus, if the loss function depends only on $\beta(P)$, the arguments given above can be used to justify restricting attention to decision rules based on $\hat{\beta}$.

□

3.4 Parametrized Policies

When the policies under consideration are complex, for example the conditional treatment rules in Example 2, the action space \mathcal{A} can be unwieldy. A key insight of Kitagawa and Tetenov (2018) is that it is often more practical (both in terms of applicability and technical analysis) to limit attention to a restricted class of policies. In this subsection I consider this idea at a general level and point out that it can often be recast as a subparameter estimation problem.

Suppose we have a statistical model $P \in \mathcal{P}$ and a welfare function $W(P, a)$ defined over \mathcal{P} and a large action space \mathcal{A} . Here the action space could consist of conditional treatment rules as in Example 2, but in other applications it could contain other types of policies. We consider a finite-dimensional, parametrized class $\mathcal{A}_1 \subset \mathcal{A}$ with

$$\mathcal{A}_1 = \left\{ a_\beta : \beta \in B \subset \mathbb{R}^k \right\}.$$

For example, one class considered by Kitagawa and Tetenov (2018) consists of rules of the form

$$\alpha_\beta(x) = \mathbf{1}\{\beta_0 + \beta_1 x \geq 0\},$$

which select individuals into treatment based on a linear eligibility score.

Let $\beta^*(P) = \arg \max_{\beta \in B} W(P, a_\beta)$ be the optimal policy parameter. It is a k -dimensional functional of P . A statistical decision rule will select some $\hat{\beta} \in B$ on the basis of the empirical data. Formally, this turns the policy choice problem into a subparameter estimation problem. The novelty is in how we evaluate the policy-parameter estimator. In the spirit of Wald's framework, we could work with welfare regret

$$L_R(P, \beta) = W(P, a_{\beta^*(P)}) - W(P, a_\beta).$$

As with the case of subparameters considered in Section 3.3, this problem can sometimes be further simplified via large-sample approximations, although it does not always admit the kind of sharp results available for regular estimation problems.

Note that we have defined regret *relative* to the best choice $\beta^*(P)$ in the parametric class, not with respect to the unconstrained best policy. Alternatively, we could define regret relative to the unconstrained best policy $a^* = \arg \max_{a \in \mathcal{A}} W(P, a)$, but this would make large-sample analysis difficult because the regret (scaled by sample size) would diverge as the sample size increases. In line with the discussion in Section 3.2, working with regret relative to $\beta^*(P)$ centers regret loss at zero, making large-sample analysis feasible. Even after scaling and centering, however, this loss $L_R(P, \beta)$ may differ substantially from conventional loss functions for point estimation.

In the binary treatment assignment case, the welfare function can be written as

$$W(P, \beta) = \mathbb{E}_P[m_0(X)] + \mathbb{E}_P[\tau(X)a_\beta(X)],$$

where $m_0(x) = \mathbb{E}_P[Y(0) | X = x]$ and $\tau(x) = \mathbb{E}_P[Y(1) - Y(0) | X = x]$. For linear eligibility score rules, the target policy parameter is

$$\beta^*(P) = \arg \max_{\beta \in B} \mathbb{E}_P[\tau(X)\mathbf{1}\{\beta_0 + \beta_1 X \geq 0\}]. \quad (3)$$

This has a similar form to Manski's maximum score problem (Manski, 1975), as was noted by Crippa (2025). Such problems are not smooth enough to admit \sqrt{n} -consistent estimation of $\beta^*(P)$, nor do they permit application of standard semiparametric efficiency theory for point estimators.

Kitagawa and Tetenov (2018) deal with these technical complications by working with the maximizer of an empirical analog of welfare,

$$\hat{\beta} = \arg \max_{\beta \in B} \widehat{W}(a_\beta),$$

and studying the difference

$$W(P, a_\beta) - \widehat{W}(a_\beta)$$

between the true (but unknown) welfare function and its empirical analog. By controlling this difference uniformly in β using modern concentration results, they can ensure that the welfare regret from using $\hat{\beta}$ instead of $\beta^*(P)$ is small with high probability. Athey and Wager (2021) refine this approach by drawing on results from semiparametric efficiency theory to construct improved empirical welfare estimators $\widehat{W}(a_\beta)$. Further work in this line includes Mbakop and Tabord-Meehan (2021), Zhou, Athey, and Wager

(2023), Nie, Brunskill, and Wager (2020), D'Adamo (2023), and Terschuur (2025).

The currently available results in the empirical welfare maximization literature certainly have a Waldean flavor, for example showing that certain procedures have welfare regret converging to zero at the minimax *rate*. While minimax rate optimality is a desirable property, it may be much less informative than local asymptotic minimax regret optimality. As a comparison, in standard randomized controlled trials, or in observational studies satisfying the unconfoundedness and overlap conditions, the minimax rate for estimating the average treatment effect is \sqrt{n} . There are estimators that achieve this rate but with a much larger asymptotic variance than the semiparametric variance bound for estimation of the ATE. For example, in an RCT one could randomly discard half the observations, and calculate the sample difference in means between the remaining treated and controls. The resulting estimator will have the same rate of convergence as the full-sample difference in means estimator, but with an inflated variance. One would not normally regard such an estimator as optimal simply by virtue of it being \sqrt{n} -consistent.

In other cases, however, the problem may be sufficiently smooth to enable much stronger large-sample optimality results. The following example is motivated by Kasy (2018).

Let $t \in \mathcal{T} = [0, 1]$ be a continuous policy variable, such as a tax rate or price. Suppose that individuals in the target population have observable characteristics X and latent characteristics ϵ , and that outcomes (in utils) are given by $Y = u(t, X, \epsilon)$. A statistical analysis identifies the average structural function

$$m(t, x) = \mathbb{E}_P [u(t, X, \epsilon) | X = x]$$

for all $x \in \mathcal{X}$ and $t \in [0, 1]$. Here P indicates the joint distribution of (ϵ, X) .

For a given policy $a : \mathcal{X} \rightarrow \mathcal{T}$, let welfare be

$$W(P, a) = \mathbb{E}_P [u(a(X), X, \epsilon)] = \mathbb{E}_P [m(a(X), X)].$$

Suppose we parametrize the policies as a_β for $\beta \in B \subset \mathbb{R}^k$, and define

$$\beta^*(P) = \arg \max_{\beta \in B} W(P, a_\beta).$$

Under sufficient smoothness of the parametrization a_β and of $m(t, x)$, and assuming that the solution $\beta^*(P)$ is interior to B , we would expect that $\beta^*(P)$ is sufficiently smooth in P to allow us to apply semiparametric efficiency theory. Moreover, the assumed smoothness would imply that $L_R(P, \beta)$ can be locally approximated by a quadratic function in a neighborhood of $P_0 \in \mathcal{P}$. Thus the loss is asymptotically equivalent to squared error loss. If we can find an estimator of $\beta^*(P)$ that is semiparametrically efficient as a point estimator, it will be locally asymptotic minimax for welfare regret.

4 Adaptive Policies and Experimental Design

In sequential statistical decision problems, data arrive sequentially over time (individually or in batches). The decision-maker can make decisions at different points in time based on the available data, and these decisions may in turn alter the data observed afterwards. Although we summarized Wald's framework for statistical decisions in its static, non-sequential form in Section 2, in fact Wald developed a sequential version of his framework early on, starting with his early work on the sequential probability ratio test (Wald, 1945a) and developed more fully in Wald (1947b), Wald (1949), and Wald (1950). This helped spawn the field of sequential statistical analysis, reviewed in Siegmund (1994) and Lai (2001). It also played an important role in the literatures on multi-armed bandits and reinforcement learning (see e.g. Lattimore and Szepesvári (2020) and Sutton and Barto (2018)).

While sequential statistical decision theory has had industrial and biostatistical applications, it has not had much impact on applied economic research until fairly recently. Empirical studies in economics often work with a fixed data set, for which the static setup of Section 2 is adequate. Recently, however, sequential decision problems have become more relevant for empirical economics, driven by both the convergence of econometrics, statistics, and computer science methodologies in dynamic e-commerce settings, and the increasing interest in designing more complex, adaptive economic experiments.

Sequential and adaptive statistical methods offer potential gains to statistical inference and decision-making. For example, experimenters may be able to test a hypothesis at lower cost by using a sequential stopping rule that terminates the experiment when sufficient evidence for or against a hypothesis has been obtained, or they may use a sequential treatment assignment rule that adaptively allocates treatments to target certain objectives. However, sequential decision rules may be much more difficult to analyze, and can exhibit complex sampling properties not seen in static settings. Even classic concepts such as sufficiency need to be considered carefully (Bahadur, 1954; Greenshtein, 1996). The following simple sequential model, loosely based on an example in Hadad, Hirshberg, Zhan, Wager, and Athey (2021), illustrates some of the subtleties of sufficiency concepts in sequential decision problems.

Example 10 (A Simple Sequential Sampling Problem). Suppose that observations Z_i are distributed iid $N(\theta, 1)$, where the only unknown parameter is $\theta \in \mathbb{R}$. Initially, in “batch 1,” we observe n draws for Z_i . Based on these observations, a statistical decision rule $A = A(Z_1, \dots, Z_n)$ selects a value from $\mathcal{A}_1 \subset (0, 1)$. Then, a further $m = [a \cdot n]$ samples are drawn for Z_i , where $a = A(Z_1, \dots, Z_n)$ and $[\cdot]$ denotes the integer ceiling (so that m is a positive integer). The observations Z_{n+1}, \dots, Z_{n+m} constitute “batch 2.” A terminal statistical decision rule $T = T(Z_1, \dots, Z_{n+m})$ maps the data from both batches into some action space \mathcal{T} . Note that T is not an ordinary statistic, because one must specify the terminal decision for every possible value of $m = [a \cdot n]$ for $a \in \mathcal{A}_1$.

Considering only batch 1, the likelihood function $\ell_1(\theta)$ can be written as

$$\ell_1(\theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}(Z_i - \theta)^2\right) \propto \exp\left(-\frac{n(\bar{Z}_1 - \theta)^2}{2}\right), \quad \text{where } \bar{Z}_1 = \frac{1}{n} \sum_{i=1}^n Z_i.$$

The sample average \bar{Z}_1 is distributed $N(\theta, 1/n)$, and is a sufficient statistic for the first batch. By Theorem 1, any statistical decision rule $A(Z_1, \dots, Z_n)$ can be represented by a rule $\bar{A}(\bar{Z}_1, U)$, where U is an independent randomization, such that A and \bar{A} have the same distribution for all θ .

The likelihood function for the batch 2 data can be written as

$$\ell_2(\theta) \propto \exp\left(-\frac{m(\bar{Z}_2 - \theta)^2}{2}\right), \quad \text{where } \bar{Z}_2 = \frac{1}{m} \sum_{i=n+1}^{n+m} Z_i,$$

while the joint likelihood function based on both batches is

$$\ell_{12}(\theta) \propto \exp\left(-\frac{(n+m)(\bar{Z}_{12} - \theta)^2}{2}\right), \quad \text{where } \bar{Z}_{12} = \frac{1}{n+m} \sum_{i=1}^{n+m} Z_i.$$

\bar{Z}_{12} is not unconditionally normally distributed in general, because the sample size $(n+m)$ depends on the data Z_1, \dots, Z_n through $A(\cdot)$. As a result, conventional statistical inference procedures based on \bar{Z}_{12} for the static normal sampling model can fail to have correct frequentist size or coverage, an issue highlighted by Hadad, Hirshberg, Zhan, Wager, and Athey (2021) and Zhang, Janson, and Murphy (2021).

Terminal decision rules T that are based on the likelihood function (such as Bayes procedures) can be written as functions of the one-dimensional statistic \bar{Z}_{12} , which aggregates the information from both batches. On the other hand, there exist pairs of rules (A, T) that cannot be represented by rules of the form $(\bar{A}(\bar{Z}_1, U), \bar{T}(\bar{Z}_{12}, U))$. Thus \bar{Z}_{12} is “likelihood-sufficient” for the terminal decision, but is not a sufficient statistic in a wider sense. What is true is that the two-dimensional statistic (\bar{Z}_1, \bar{Z}_2) is sufficient: for any pair (A, T) , there are rules $\bar{A} = \bar{A}(\bar{Z}_1, U)$ and $\bar{T} = \bar{T}(\bar{Z}_1, \bar{Z}_2, U)$ such that

$$(\bar{A}, \bar{T}) \sim (A, T) \quad \forall \theta.$$

□

The preceding analysis was very stylized but already points to some of the complications that arise when analyzing sequential decision rules. Given the additional technical challenges encountered in sequential problems, there is perhaps even more to be gained from employing dimension-reduction and other techniques to make the analysis more tractable. Local asymptotic approximations of the kind discussed in Section 3.2 appear initially in Le Cam (1986), Ch. 13, for the case of sequential stopping problems. More recently, they have been used to characterize bandit algorithms and adaptive experimental methods by a number of authors, including Adusumilli (2024, 2025), Armstrong (2025), Chen and Andrews (2023), Hadad, Hirshberg, Zhan, Wager, and Athey (2021), Higbee (2025), Hirano and Porter (2025), Kalvit and Zeevi (2021), Kuang and Wager (2024), Niu and Ren (2025), Xu and Zhou (2025), and Zhang, Janson, and Murphy (2021). The following example illustrates how large sample approximations can simplify the analysis of an adaptive statistical decision problem.

Example 11 (Local Asymptotic Approximation of a Two Stage Adaptive Experiment). Consider a two-batch setup as in Example 10, but now with a general parametric model instead of normally distributed

observations. In each batch we observe

$$Z_i \stackrel{iid}{\sim} P_\theta, \quad \theta \in \Theta \subset \mathbb{R}^k,$$

where the parametric model for Z_i satisfies the same regularity conditions as in Theorem 2.

Suppose that (A_n, T_n) converges jointly in distribution under every $\theta_0 + h/\sqrt{n}$. Then, by Theorem 2 in Hirano and Porter (2025), the limits of (A_n, T_n) can be represented by

$$\bar{A}(\bar{Z}_1, U), \quad \bar{T}(\bar{Z}_1, \bar{Z}_2, U),$$

where

$$\bar{Z}_1 \sim N(h, \Sigma), \quad \bar{Z}_2 | \bar{Z}_1, U \sim N(h, \Sigma / \bar{A}).$$

These representations extend to multiple treatment arms and more than 2 batches.

□

This example shows that the normal case of Example 10 (or more precisely, a multivariate normal version of it) can serve as an approximate representation of the statistical decision problem when the underlying data are not normal, but follow a smooth parametric model. While the normal model may still present challenges, it emerges as a canonical case for which to seek solutions.

5 Conclusion

This paper has surveyed some recent work that uses or draws inspiration from Wald's statistical decision theory to study empirical policy choice and experimental design. While the framework is general and powerful, solving statistical decision problems in realistic settings can be challenging. Exact (finite-sample) solutions may not always be feasible, but problems may be recast in more tractable forms, while preserving some of the conceptual elegance of Wald's original theory. Strategies that involve reducing the dimensionality of the underlying optimization problems, combined with approximations and computational tools, can be helpful toward these goals.

References

- ADUSUMILLI, K. (2024): “Optimal Tests Following Sequential Experiments,” .
- (2025): “Risk and Optimal Policies in Bandit Experiments,” *Econometrica*, 93(3), 1003–1029.
- ARADILLAS FERNANDEZ, A., J. BLANCHET, J. L. MONTIEL OLEA, C. QIU, J. STOYE, AND L. TAN (2024): “ ϵ -Minimax Solutions of Statistical Decision Problems via the Hedge Algorithm,” .
- ARMSTRONG, T. B. (2025): “Asymptotic Efficiency Bounds for a Class of Experimental Designs,” .
- ARMSTRONG, T. B., T. KITAGAWA, AND A. TETENOV (2024): “Statistical Decision Theory and Empirical Practice,” .
- ARMSTRONG, T. B., M. KOLESÁR, AND M. PLAGBORG-MØLLER (2022): “Robust Empirical Bayes Confidence Intervals,” *Econometrica*, 90(6), 2567–2602.
- 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.
- ATHEY, S., AND S. WAGER (2021): “Policy Learning With Observational Data,” *Econometrica*, 89(1), 133–161.
- BAHADUR, R. R. (1954): “Sufficiency and Statistical Decision Functions,” *The Annals of Mathematical Statistics*, 25(3), 423–462.
- BLACKWELL, D. (1951): “Comparison of Experiments,” in *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability*, ed. by J. Neyman, pp. 93–102. University of California Press.
- BROWN, L. D. (1994): “Minimaxity, More or Less,” in *Statistical Decision Theory and Related Topics V*, pp. 1–18. Springer-Verlag, New York.
- (2000): “An Essay on Statistical Decision Theory,” *Journal of the American Statistical Association*, 95(452), 1277–1281.
- CHAMBERLAIN, G. (2000): “Econometric Applications of Maxmin Expected Utility,” *Journal of Applied Econometrics*, 15(6), 625–644.
- CHEN, J. (2023): “Empirical Bayes When Estimation Precision Predicts Parameters,” .
- CHEN, J., AND I. ANDREWS (2023): “Optimal Conditional Inference in Adaptive Experiments,” .
- CHOI, S., W. J. HALL, AND A. SCHICK (1996): “Asymptotically Uniformly Most Powerful Tests in Parametric and Semiparametric Models,” *The Annals of Statistics*, 24(2).
- CHRISTENSEN, T., H. R. MOON, AND F. SCHORFHEIDE (2023): “Optimal Decision Rules When Payoffs Are Partially Identified,” .

- CRIPPA, F. (2025): “Regret Analysis in Threshold Policy Design,” *Journal of Econometrics*, 249, 105998.
- D’ADAMO, R. (2023): “Orthogonal Policy Learning Under Ambiguity,” .
- DEHEJIA, R. H. (2005): “Program Evaluation as a Decision Problem,” *Journal of Econometrics*, 125, 141–173.
- ESCANCIANO, J. C., AND J. R. TERSCHUUR (2023): “Machine Learning Inference on Inequality of Opportunity,” .
- GREENSHTAIN, E. (1996): “Comparison of Sequential Experiments,” *The Annals of Statistics*, 24(1).
- GU, J., AND R. KOENKER (2017): “Unobserved Heterogeneity in Income Dynamics: An Empirical Bayes Perspective,” *Journal of Business & Economic Statistics*, 35, 1–16.
- GUGGENBERGER, P., AND J. HUANG (2025): “On the Numerical Approximation of Minimax Regret Rules Via Fictitious Play,” .
- GUGGENBERGER, P., N. MEHTA, AND N. PAVLOV (2024): “Minimax Regret Treatment Rules with Finite Samples When a Quantile Is the Object of Interest,” .
- HADAD, V., D. A. HIRSHBERG, R. ZHAN, S. WAGER, AND S. ATHEY (2021): “Confidence Intervals for Policy Evaluation in Adaptive Experiments,” *Proceedings of the National Academy of Sciences*, 118(15), e2014602118.
- HÁJEK, J. (1970): “A Characterization of Limiting Distributions of Regular Estimates,” *Zeitschrift für Wahrscheinlichkeitstheorie und Verwandte Gebiete*, 14(4), 323–330.
- HALL, W. J., R. A. WIJSMAN, AND J. K. GHOSH (1965): “The Relationship Between Sufficiency and Invariance with Applications in Sequential Analysis,” *The Annals of Mathematical Statistics*, 36(2), 575–614.
- HALMOS, P. R., AND L. J. SAVAGE (1949): “Application of the Radon-Nikodym Theorem to the Theory of Sufficient Statistics,” *The Annals of Mathematical Statistics*, 20(2), 225–241.
- HIGBEE, S. D. (2025): “Experimental Design for Policy Choice,” .
- HIRANO, K., AND J. R. PORTER (2009): “Asymptotics for Statistical Treatment Rules,” *Econometrica*, 77(5), 1683–1701.
- (2020): “Asymptotic Analysis of Statistical Decision Rules in Econometrics,” in *Handbook of Econometrics*, vol. 7, pp. 283–354. Elsevier.
- (2025): “Asymptotic Representations for Sequential Decisions, Adaptive Experiments, and Batched Bandits,” .
- IBRAGIMOV, I., AND R. HASMINSKII (1981): *Statistical Estimation: Asymptotic Theory*. Springer-Verlag, New York.

- KALVIT, A., AND A. ZEEVI (2021): “A Closer Look at the Worst-case Behavior of Multi-armed Bandit Algorithms,” *35th Conference on Neural Information Processing Systems (NeurIPS 2021)*.
- KASY, M. (2016): “Partial Identification, Distributional Preferences, and the Welfare Ranking of Policies,” *Review of Economics and Statistics*, 98(1), 111–131.
- (2018): “Optimal Taxation and Insurance Using Machine Learning — Sufficient Statistics and Beyond,” *Journal of Public Economics*, 167, 205–219.
- KIDO, D. (2023): “Locally Asymptotically Minimax Statistical Treatment Rules Under Partial Identification,” .
- KIM, D.-H. (2013): “Optimal Choice of a Reserve Price Under Uncertainty,” *International Journal of Industrial Organization*, 31(5), 5887–602.
- KITAGAWA, T., S. LEE, AND C. QIU (2025): “Leave No One Undermined: Policy Targeting with Regret Aversion,” .
- KITAGAWA, T., AND A. TETENOV (2018): “Who Should Be Treated? Empirical Welfare Maximization Methods for Treatment Choice,” *Econometrica*, 86(2), 591–616.
- (2021): “Equality-Minded Treatment Choice,” *Journal of Business & Economic Statistics*, 39(2), 561–574.
- KOCK, A. B., D. PREINERSTORFER, AND B. VELIYEV (2023): “Treatment Recommendation with Distributional Targets,” *Journal of Econometrics*, 234(2), 624–646.
- KOENKER, R., AND J. GU (2024): “Empirical Bayes for the Reluctant Frequentist,” .
- KUANG, X., AND S. WAGER (2024): “Weak Signal Asymptotics for Sequentially Randomized Experiments,” *Management Science*, 70(10), 7024–7041.
- LAI, T. L. (2001): “Sequential Analysis: Some Classical Problems and New Challenges,” *Statistica Sinica*, 11, 303–408.
- LATTIMORE, T., AND C. SZEPESVÁRI (2020): *Bandit Algorithms*. Cambridge University Press, 1 edn.
- 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.
- LEHMANN, E. L., AND G. CASELLA (1998): *Theory of Point Estimation*. Springer, New York, second edn.
- LIESE, F., AND K.-J. MIESCKE (2008): *Statistical Decision Theory: Estimation, Testing, and Selection*. Springer, New York.

- LITVIN, V., AND C. F. MANSKI (2021): “Evaluating the Maximum Regret of Statistical Treatment Rules with Sample Data on Treatment Response,” *The Stata Journal*, 21(1), 97–122.
- MANSKI, C. F. (1975): “Maximum Score Estimation of the Stochastic Utility Model of Choice,” *Journal of Econometrics*, 3(3), 205–228.
- (1988): “Ordinal Utility Models of Decision Making Under Uncertainty,” *Theory and Decision*, 25(1), 79–104.
- (2004): “Statistical Treatment Rules for Heterogeneous Populations,” *Econometrica*, 72(4), 1221–1246.
- (2009): “The 2009 Lawrence R. Klein Lecture: Diversified Treatment Under Ambiguity,” *International Economic Review*, 50(4), 1013–1041.
- (2010): “Vaccination with Partial Knowledge of External Effectiveness,” *Proceedings of the National Academy of Sciences*, 107(9), 3953–3960.
- (2011): “Choosing Treatment Policies Under Ambiguity,” *Annual Review of Economics*, 3, 25–49.
- (2024): “Identification and Statistical Decision Theory,” *Econometric Theory*, pp. 1–17.
- MANSKI, C. F., AND M. TABORD-MEEHAN (2017): “Evaluating the Maximum MSE of Mean Estimators with Missing Data,” *The Stata Journal*, 17(3), 723–735.
- MANSKI, C. F., AND A. TETENOV (2023): “Statistical Decision Theory Respecting Stochastic Dominance,” *The Japanese Economic Review*, 74(4), 447–469.
- MASTEN, M. A. (2023): “Minimax-Regret Treatment Rules with Many Treatments,” *The Japanese Economic Review*, 74(4), 501–537.
- MBAKOP, E., AND M. TABORD-MEEHAN (2021): “Model Selection for Treatment Choice: Penalized Welfare Maximization,” *Econometrica*, 89(2), 825–848.
- MCMANUS, D. A. (1991): “Who Invented Local Power Analysis?,” *Econometric Theory*, 7(2), 265–268.
- MONTIEL OLEA, J. L., B. O’FLAHERTY, AND R. SETHI (2021): “Empirical Bayes Counterfactuals in Poisson Regression with an Application to Police Use of Deadly Force.”
- MONTIEL OLEA, J. L., C. QIU, AND J. STOYE (2025): “Decision Theory for Treatment Choice Problems with Partial Identification.”
- MÜLLER, U. K. (2011): “Efficient Tests Under a Weak Convergence Assumption,” *Econometrica*, 79(2), 395–435.
- NEYMAN, J. (1937): “Smooth Test for Goodness of Fit,” *Scandinavian Actuarial Journal*, 1937(3-4), 149–199.

- NIE, X., E. BRUNSKILL, AND S. WAGER (2020): “Learning When-to-Treat Policies,” .
- NIU, Z., AND Z. REN (2025): “Assumption-Less Weak Limits and Tests for Two-Stage Adaptive Experiments,” .
- PFANZAGL, J. (2017): *Mathematical Statistics*, Springer Series in Statistics. Springer Berlin Heidelberg, Berlin, Heidelberg.
- ROBBINS, H. (1956): “An Empirical Bayes Approach to Statistics,” in *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*.
- (1985): “Asymptotically Subminimax Solutions of Compound Statistical Decision Problems,” in *Herbert Robbins Selected Papers*, ed. by T. L. Lai, and D. Siegmund, pp. 7–24. Springer New York, New York, NY.
- ROSTEK, M. (2009): “Quantile Maximization in Decision Theory,” *Review of Economic Studies*, 77(1), 339–371.
- SAVAGE, L. J. (1951): “The Theory of Statistical Decision,” *Journal of the American Statistical Association*, 46(253), 55–67.
- SAVAGE, L. J. (1972): *The Foundations of Statistics*. Dover, New York, 2nd edn.
- SCHLAG, K. H. (2006): “ELEVEN - Tests Needed for a Recommendation,” .
- SIEGMUND, D. (1994): “A Retrospective of Wald’s Sequential Analysis—Its Relation to Change-point Detection and Sequential Clinical Trials,” in *Statistical Decision Theory and Related Topics V*. Springer-Verlag, New York.
- STEIN, C. (1956): “Inadmissibility of the Usual Estimator for the Mean of a Multivariate Normal Distribution,” in *Contribution to the Theory of Statistics*, pp. 197–206. University of California Press.
- STOYE, J. (2009): “Minimax Regret Treatment Choice with Finite Samples,” *Journal of Econometrics*, 151, 70–81.
- SUTTON, R. S., AND A. G. BARTO (2018): *Reinforcement Learning: An Introduction*, Adaptive Computation and Machine Learning Series. The MIT Press, Cambridge, Massachusetts, second edition edn.
- TERSCHUUR, J. (2025): “Locally Robust Policy Learning: Inequality, Inequality of Opportunity and Intergenerational Mobility,” .
- TETENOV, A. (2012): “Statistical Treatment Choice Based on Asymmetric Minmax Regret Criteria,” *Journal of Econometrics*, 166, 157–165.
- VAN DER VAART, A. (1991): “An Asymptotic Representation Theorem,” *International Statistical Review / Revue Internationale de Statistique*, 59(1), 97.

- VAN DER VAART, A. W. (1998): *Asymptotic Statistics*. Cambridge University Press, New York.
- VAPNIK, V. (1991): “Principles of Risk Minimization for Learning Theory,” in *Advances in Neural Information Processing Systems*, vol. 4. Morgan-Kaufmann.
- WALD, A. (1939): “Contributions to the Theory of Statistical Estimation and Testing Hypotheses,” *The Annals of Mathematical Statistics*, 10(4), 299–326.
- (1943): “Tests of Statistical Hypotheses Concerning Several Parameters When the Number of Observations Is Large,” *Transactions of the American Mathematical Society*, 54(3), 426–482.
- (1945a): “Sequential Tests of Statistical Hypotheses,” *The Annals of Mathematical Statistics*, 16(2), 117–186.
- (1945b): “Statistical Decision Functions Which Minimize the Maximum Risk,” *Annals of Mathematics*, 46(2), 265–280.
- (1947a): “An Essentially Complete Class of Admissible Decision Functions,” *The Annals of Mathematical Statistics*, 18(4), 549–555.
- (1947b): “Foundations of a General Theory of Sequential Decision Functions,” *Econometrica*, 15(4), 279–313.
- (1949): “Statistical Decision Functions,” *The Annals of Mathematical Statistics*, 20(2), 165–205.
- (1950): *Statistical Decision Functions*. Wiley, New York.
- WOLFOWITZ, J. (1952): “Abraham Wald, 1902-1950,” *The Annals of Mathematical Statistics*, 23(1), 1–13.
- XU, H. (2025): “Asymptotic Analysis of Point Decisions with General Loss Functions,” .
- XU, Y., AND B. ZHOU (2025): “Batched Adaptive Network Formation,” .
- YATA, K. (2023): “Optimal Decision Rules Under Partial Identification,” .
- ZHANG, K. W., L. JANSON, AND S. A. MURPHY (2021): “Inference for Batched Bandits,” .
- ZHOU, Z., S. ATHEY, AND S. WAGER (2023): “Offline Multi-Action Policy Learning: Generalization and Optimization,” *Operations Research*, 71(1), 148–183.