

The Economic Case for Probability-Based Sentencing

Ron Siegel and Bruno Strulovici*

September 6, 2019

Abstract

Evidence in criminal trials determines whether the defendant is found guilty, but is usually not one of the factors formally considered during sentencing. A number of legal scholars have advocated the use of sentences that reflect the strength of evidence. This paper proposes an economic model that unifies the arguments put forward in this literature and addresses three of the remaining objections to the use of evidence-based sentencing: i) political legitimacy (the impact on the coercive power of the state), ii) robustness to details of the environment, and iii) incentives to acquire evidence.

*We thank Robert Burns, Andy Daughety, Eddie Dekel, Louis Kaplow, Fuhito Kojima, Adi Leibovitz, Aki Matsui, Paul Milgrom, Jennifer Reiganum, Ariel Rubinstein, Kathy Spier, Jean Tirole, and Leeat Yariv for their comments. The paper benefited from the reactions of seminar participants at UC Berkeley, Seoul National University, the NBER, the World Congress of the Econometric Society, the Harvard/MIT Theory workshop, Caltech's NDAM conference, Duke, Penn State, Johns Hopkins, the Pennsylvania Economic Theory Conference, Bocconi University, Oxford University, Kyoto University, Tokyo University, the Toulouse School of Economics, the Harris School of Public Policy, and the Summer School of the Econometric Society 2019. David Rodina provided excellent research assistance. Previous versions of the paper were circulated under the name "Multiverdict Systems." Strulovici acknowledges financial support from an NSF CAREER Award (Grant No. 1151410) and a fellowship from the Alfred P. Sloan Foundation. Siegel: Department of Economics, The Pennsylvania State University, University Park, PA 16802, rus41@psu.edu. Strulovici: Department of Economics, Northwestern University, Evanston, IL 60208, b-strulovici@northwestern.edu.

1 Introduction

Criminal cases traditionally end with one of two outcomes: acquittal or conviction. An acquittal entails no punishment, regardless of the strength of evidence in the case, whereas a conviction requires that guilt be established “beyond a reasonable doubt” (or some other demanding criterion) and entails a sentence that does not (officially, at least) depend on any residual doubt concerning the defendant’s guilt.

Recently, a number of scholars (Bray (2005), Lando (2005), Fisher (2011, 2012), Teichman (2017), and Spottswood (2019)) have advocated the use of a *probabilistic sentencing* model that would reflect the strength of evidence and residual doubt more gradually than the *threshold model* that underlies the binary verdicts used in practice.

These authors put forth several arguments in favor of probabilistic sentencing, which range from achieving more effective deterrence to improving outcome equality for similar trials and to reducing the effect of biases on judicial outcomes. This recent development in legal thought parallels an earlier literature concerning probability-based damages, pioneered by Coons (1966) and further explored by Kaye (1982), Shavell (1985), Fischer (1993), and Davis (2003), among others.¹ Legal scholars have also noted that binary verdicts may owe their persistence to idiosyncratic historical roots rather than to deep philosophical principles and refuted a number of philosophical objections to the use of probabilistic sentencing.²

¹Similar ideas have recently found their way into actual legal practice, such as the “loss of chance” doctrine in medical malpractice lawsuits (Fischer (2001)) and the concept of market share liability in tort law (Rostron (2004)), a theory under which a plaintiff unable to identify the manufacturer of the product that caused his injury can recover on a proportional basis from each manufacturer that might have made the product.

²Fisher (2012) provides an in-depth account and defense of the philosophical implications of probabilistic sentencing, which includes the expressive role of verdicts (i.e., the message sent to society about a particular defendant and a particular crime) and the retributive role of verdicts, which requires that, as a moral principle, guilty defendants must be punished and innocent defendants must be acquitted regardless of any utilitarian consideration. In particular, Fisher observes (p. 875) that “Retributivists would criticize the probabilistic model for eroding the evidence threshold and thereby weakening the moral legitimacy of inflicting criminal punishment. Imposing punishment based on evidence with a probative weight falling below the beyond-a-reasonable-doubt standard (i.e., below the standard of moral certainty) results in the defendant’s objectification, causing moral harm that is incommensurable with deterrence utilities,” Fisher answers (p. 876) that “retributivist theories cannot, in fact, justify any type of decision regime, the present threshold model included.” She quotes Reiman and van den Haag (1990), who assert that “the retributivist theories cannot justify the beyond-a-reasonable-

The recent advocacy of probabilistic sentencing has been weaker along three dimensions. The first one concerns *political legitimacy*: the fact that probabilistic sentencing would give the state the right to incarcerate citizens of “less than certain” guilt. Fisher (2012, p. 882) describes this concern as follows: “the threshold model, and its pivotal requirement of proof beyond a reasonable doubt, are, thus, a matter of political morality and serve as one of the main limitations on state power in modern systems,” and acknowledges (p. 883) that “the challenge posed by probabilistic decision making to the liberal state model is, indeed, a strong case against its implementation.”³

Second, probabilistic sentencing may have an adverse effect on the incentives to acquire evidence. In particular, Lando (2005, p. 285) expresses a concern that “graduating sanctions may affect the strategies and incentives to search for evidence of both defense lawyers and prosecutors” but notes that “a it is not obvious what the effect will be.”

Third, the robustness of improvements offered by probabilistic sentencing have yet to be explored. Beyond qualitative arguments and Bayesian computations requiring minute knowledge of all institutional details concerning a case, are there robust ways of improving welfare with probabilistic sentencing? One reason that little is known about robustness is that most legal theories of probabilistic sentencing have, despite the quantitative nature of probabilities, been discussed qualitatively, with the exception of Lando (2005), and to some extent Fisher (2011), who anchors her analysis with some numerical examples. This has led to some confusion as to which insights of this literature were actually correct. In particular, Fisher (2012, p. 856) suggests that Lando’s analysis may in fact *justify*, rather than discredit, the threshold model, and argues that defendants’ risk attitudes, as well as the adverse effect of wrongful acquittals on deterrence, must be taken into account in order to justify the probabilistic sentencing model.⁴

doubt threshold, as it is inconsistent with the retributivist commitment to punishing the guilty: ”Should we try to convict fewer innocents and risk letting more of the guilty escape, or try to convict more of the guilty, and, unavoidably, more of the innocent? Retributivism (although not necessarily retributivists) is mute on how high standards of proof ought to be...”

³Fisher does propose arguments in favor of probabilistic voting even in this case, the most prominent of which echoes some of Lando’s (2005) analysis. As we explain below, however, these arguments are problematic because they entail raising the sentence to some defendants who are found guilty.

⁴Another confusion, which underlies this literature, concerns the view that deterrence is achieved more effectively by threshold rules than by probabilistic ones. This view can be traced to Kaye (1982) who, in the context of tort law, argues that the preponderance of evidence criterion (which is a threshold rule) provides more effective

This paper offers an economic analysis of probabilistic sentencing that addresses these three issues. To simplify the exposition, we focus our analysis on the addition of a third, intermediate verdict to reflect more finely the strength evidence against a defendant, and later discuss an extension to more verdicts.

In our framework, citizens benefit heterogeneously from committing a crime and abstain from its commission if the expected disutility from punishment exceeds their benefit from the crime. We consider a welfare objective function that captures deterrence, type I/type II error considerations, and the cost of generating evidence.

Our first set of results shows how an intermediate verdict can be used to improve welfare while addressing fully the political legitimacy concern. We show that, starting from the best possible binary verdict (i.e., one for which the sentence is a 2-step function of the strength of evidence), it is always possible to construct a three-verdict system that i) increases expected welfare, ii) does not increase the probability of conviction of any defendant, and iii) does not increase the sentence length of any defendant. We show that such an improvement exists even if deterrence plays a major role in the social objective. Part iii) guarantees that probabilistic sentencing improves welfare, even under the constraint that the state cannot be more coercive than under the binary verdict system.

Moreover, we show that the form of the three-verdict improvement does not depend on the minute details of the welfare and utility functions used for the analysis, on the particular threshold used by the binary verdict, or on the particular technology used to gather evidence. Our construction thus addresses both the political legitimacy and the robustness concerns.

Next, we exploit the expressive function of verdicts to provide a second answer to the political legitimacy concern: when criminal trials carry some stigma for the defendant, we show that acquittals too can be split into two distinct verdicts that entail no direct punishment but carry different stigmas. Our result provides an explicit illustration of Fisher's (2012) argument that the expressive function of trials can in fact be enhanced by finer evidence-based verdicts and a formal

deterrence than probabilistic damages. This view has been developed further by Shavell (1987) and is invoked in criminal law to suggest that the social costs arising from type I and type II errors (wrongful convictions and wrongful acquittals) would be better addressed by a threshold rule than by probabilistic sentencing. Thus, Fisher (2012, p. 861) is concerned that "Incorporating the ex post error-costs into the equation may, indeed, constrict the probabilistic regime's scope of operation." These confusions call for a formal framework to investigate the validity and robustness of earlier insights.

foundation for Bray’s (2005) advocacy of “not proven” verdicts. Several countries, including Israel, Italy, and Scotland distinguish among acquitted defendants based on the residual doubt regarding their guilt. In Scotland, for example, a conviction in a criminal trial leads to a “guilty” verdict, but an acquittal leads to either a verdict of “not guilty” or “not proven.” Neither of the two acquittal verdicts carries any jail time, but the latter indicates a higher likelihood that the defendant is in fact guilty.⁵ The likelihood is, however, insufficiently high for conviction.⁶

To address the incentives for acquiring evidence, we explore the consequence of adding a third verdict for the value of acquiring evidence. As noted by Lando (2005), one may a priori be concerned that probabilistic sentencing reduces incentives for acquiring evidence. We show that introducing an intermediary verdict can *increase* the value of evidence acquisition, and in fact *systematically* for some regions of the belief space. To study the value of information acquisition, we begin with a simple model of one-shot information acquisition, and then embed the analysis into a modern, tractable model of continuous information acquisition, for which the shape of the value function is easy to characterize and the systematic positive impact of a third verdict on the value of evidence is more easily presented.

In our framework, a defendant can have an arbitrary risk attitude with respect to his sentence and an arbitrary prior probability of guilt,⁷ the social costs due type I and type II errors follow arbitrary specifications, the risk of wrongful acquittal is fully taken into account, and the relative weights of deterrence and error considerations are also arbitrary. We show that probabilistic sentencing is valuable irrespective of these dimensions, and thus unify and clarify earlier insights. To focus on the trial aspect of criminal justice, and in line with the aforementioned literature, our model does not allow for plea bargaining. While plea bargaining is of paramount importance in the United States, its use in many other countries is much more limited, especially for serious

⁵The introduction of a not-proven verdict is considered by Daughety and Reinganum (2015a), who study how the effect of informal sanctions on defendants and prosecutors affect the plea bargaining process and its acceptance rate, and consider the effect of a not-proven verdict in this context. Daughety and Reinganum (2015b) consider two implementations of a not-proven verdict. In the first one, the defendant can choose between the standard binary verdict system and the system with a not-proven verdict. In equilibrium, all defendants choose the latter system. The authors also analyze an alternative implementation in which some defendants who are found not guilty are compensated.

⁶This may happen, for example, if an eye-witness testimony exists, but the testimony cannot be corroborated.

⁷Both features distinguish our work from Lando (2005), who assumes risk neutrality and implicitly imposes a uniform prior.

crimes. In Siegel and Strulovici (2019) we consider a different question of criminal trial design in a model that allows for plea bargains.

The appendix contains the continuous-time evidence gathering model and a micro-foundation for the Bayesian formulation used in Section 4.1. This micro-foundation establishes that trial technology conceptualized as a mapping from accumulated evidence to a verdict can always be reformulated in Bayesian fashion: accumulated evidence is a signal that turns the prior probability that the defendant is guilty into a posterior probability, on which the verdict is based. Moreover, this transformation establishes a relationship between two notions of ‘incriminating’ and ‘exculpatory’ evidence. One notion is based on decisions and the other on beliefs. What makes a piece of evidence ‘incriminating’ is the fact that it increases the likelihood of guilt of a defendant and, hence, results in a longer expected sentence. In particular, there is no loss of generality when one says that a guilty defendant is more likely than innocent defendant to generate incriminating evidence.

2 Model

We consider a trial whose objective is to determine whether a defendant is guilty of committing a certain crime and to deliver the corresponding sentence. In our baseline model the trial is summarized by two numbers: the probability π_g that the defendant is found guilty if he is actually guilty, and the probability π_i that the defendant is found guilty if he is actually innocent.⁸ Corresponding to a guilty verdict is a sentence $s > 0$, interpreted as jail time (so a higher value of s corresponds to a harsher punishment).⁹

Society wishes to avoid punishing the defendant if he is innocent, and adequately punish him if he is guilty. This dual goal is modeled by a differentiable welfare function W . Jailing an innocent defendant for s years leads to a welfare of $W(s, i)$, with $W(0, i) = 0$ and W decreasing in s . Jailing a guilty defendant leads to a welfare of $W(s, g)$, which has a single peak at $\hat{s} > 0$. Thus, \hat{s} is the punishment deemed optimal by society if it is certain that the defendant is guilty. The assumption that $W(s, g)$ increases up to \hat{s} and then decreases is in line with US sentencing guidelines, which state that “The court shall impose a sentence sufficient, but not greater than

⁸It is natural to assume that $\pi_g > \pi_i$, i.e., a defendant is more likely to be found guilty if he is actually guilty than if he is innocent. This assumption is, however, not required for this section.

⁹We leave aside such issues as mitigating circumstances, which are tangential to the focus of the paper.

necessary, to...reflect the seriousness of the offense... and to provide just punishment for the offense.”¹⁰

The relative importance of punishing the defendant if he is guilty and not punishing him if he is innocent depends on the prior probability $\lambda \in (0, 1)$ that the defendant is guilty. This is captured by the interim social welfare from the defendant going to trial when the punishment of being found guilty is s :

$$\mathcal{W}_2(s) = \lambda [\pi_g W(s, g) + (1 - \pi_g)W(0, g)] + (1 - \lambda) [\pi_i W(s, i) + (1 - \pi_i)W(0, i)]. \quad (1)$$

Since $W(\cdot, i)$ is decreasing and $W(\cdot, g)$ peaks at \hat{s} , it is never interim optimal to choose $s > \hat{s}$.

Society’s ex-ante welfare also depends on whether the crime is committed in the first place. To model this, we consider an individual’s decision whether to commit the crime, and assume that at most one individual is prosecuted for the crime if it is committed.¹¹ In a large society, the probability that any particular innocent individual is prosecuted for the crime is infinitesimal, so for expositional convenience we assume that an innocent individual treats this probability as 0.¹² If the individual commits the crime, he obtains a benefit b (in utility terms), but faces a probability η_g of being arrested and prosecuted.¹³ Thus, the individual commits the crime if

$$b + \eta_g (\pi_g u(s) + (1 - \pi_g)u(0)) > 0, \quad (2)$$

where $u(s) \leq 0$ is the defendant’s differentiable utility from a sentence s , and the utility from not being prosecuted is normalized to 0. Denote by $H(s)$ the fraction of individuals for whom (2) holds. The benefit b is distributed in the population according to an absolutely continuous cdf B , so by (2) we have

$$H(s) = 1 - B(-\eta_g (\pi_g u(s) + (1 - \pi_g)u(0))). \quad (3)$$

¹⁰See 18 U.S.C § 3553. These guidelines also state that another goal is “to protect the public from further crimes of the defendant.” This incapacitation reasonably increases at a rate that decreases in the sentence, whereas the disutility a prisoner experiences increases with his sentence, which together may also give rise to single-peaked social welfare.

¹¹This allows us to abstract from interdependencies between multiple defendants, an issue that is tangential to the focus of this paper. See Silva (2018) for an analysis of such issues.

¹²The probability $1 - \lambda > 0$ that a prosecuted individual is innocent is, however, not infinitesimal.

¹³This probability can be endogenized by including the amount of costly law enforcement as a decision variable without changing any of the results.

By normalizing the welfare from no crime to 0, we obtain that the ex-ante social welfare is

$$H(s) (\eta_g (\pi_g W(s, g) + (1 - \pi_g) W(0, g)) + \eta_i (\pi_i W(s, i) + (1 - \pi_i) W(0, i)) - h), \quad (4)$$

where η_i is the probability that an innocent defendant is prosecuted and h is the social harm from the crime.¹⁴ In particular, since that sentence s determines which individuals commit the crime and which are deterred,¹⁵ the sentence affects social welfare beyond its direct effect on the ex-post welfare W . Finally, when an individual is prosecuted the crime has already been committed so the social harm h from the crime is “sunk” and the prior that the defendant is guilty is $\lambda = \eta_g / (\eta_g + \eta_i)$, so we recover (1).

Because we will later consider more than two verdicts, we rewrite the interim social welfare (1) more generally as

$$\lambda E_g (W(\tilde{s}, g)) + (1 - \lambda) E_i (W(\tilde{s}, i)), \quad (5)$$

where the sentence \tilde{s} is a random variable whose distribution depends on whether the defendant committed the crime. Similarly, we rewrite (2) as

$$b + \eta_g E_g (u(\tilde{s})) > 0, \quad (6)$$

and rewrite (3) as

$$H(\tilde{s}) = 1 - B(-\eta_g E_g (u(\tilde{s}))), \quad (7)$$

where $H(\tilde{s})$ is the fraction of individuals for whom (6) holds. We rewrite (4) as

$$H(\tilde{s}) (\eta_g E_g (W(\tilde{s}, g)) + \eta_i E_i (W(\tilde{s}, i)) - h). \quad (8)$$

Throughout the analysis we assume that all sentences are interior, in the sense that they can be made more severe.¹⁶ We also assume that the harm caused by the crime exceeds the social welfare from punishing its perpetrator, i.e.,

$$W(\hat{s}, g) - h < 0. \quad (9)$$

¹⁴The benefit from committing the crime can be considered explicitly as well without affecting any of the results.

¹⁵Guidelines 18 U.S.C § 3553 state that another goal of punishment is “to afford adequate deterrence to criminal conduct.”

¹⁶This can be done by imposing a longer or harsher imprisonment term. Even an execution can be made more severe by making it less humane. While an extreme sentence would maximize crime deterrence, it would also deter (or “chill”) desirable behavior (Kaplow (2011)) and excessively punish those individuals who were not deterred and committed the crime, perhaps because they were ignorant of the possible punishment or did not rationally assess the consequences of their crime before its commission. Formally, the optimal sentence will

3 An intermediate ‘guilty’ verdict

We consider adding a verdict that refines the ‘guilty’ verdict from the two-verdict system.¹⁷ Those defendants who would be convicted in the two-verdict system now receive one of two “guilty verdicts,” which we denote 1 and 2. Defendants who would be acquitted in the two-verdict system are still acquitted and are released.¹⁸ The distinction between the two ‘guilty’ verdicts may be based on the evidence available before and during the trial, so that among the collections of evidence that would lead to a conviction in the two-verdict system some lead to verdict 1 and the remaining to verdict 2.¹⁹ Denote by π_i^1 the probability that the defendant receives verdict 1 if he is innocent, and define π_i^2 , π_g^1 , and π_g^2 similarly.²⁰ Because the same set of defendants is acquitted as in the two-verdict case, we have

$$\pi_i = \pi_i^1 + \pi_i^2 \quad \text{and} \quad \pi_g = \pi_g^1 + \pi_g^2.$$

Without loss of generality²¹

$$\frac{\pi_g^2}{\pi_i^2} > \frac{\pi_g}{\pi_i} > \frac{\pi_g^1}{\pi_i^1},$$

so verdict 1 is an “intermediate verdict:” a guilty defendant is more likely to receive verdict 2, relative to an innocent defendant, than verdict 1.

Let s_j denote the sentence associated with verdict j . Given s_1 and s_2 , the interim social

be interior, even taking deterrence into account, if i) the maximal benefit from the crime exceeds the maximal disutility from the harshest possible sentence (e.g., benefits have an unbounded support and the defendant’s utility is bounded below), and ii) social welfare becomes sufficiently negative as defendants’ punishment becomes sufficiently harsh.

¹⁷Further intermediate verdicts may similarly be added, as discussed in Section 4.2.

¹⁸Section 8 discusses how to implement the additional verdict in a way that is likely not to affect jurors’ decision whether to acquit the defendant. It also discusses how the analysis might change if their decision is affected.

¹⁹Evidence leading to a homicide conviction in the two-verdict system may include, for example, the discovery, in the defendant’s house, of the gun from which the bullet was fired, a confession by the defendant, a death threat made by the defendant to the victim shortly before the murder, or any subset of these.

²⁰In keeping with most of the literature on trial design, we take a reduced-form approach to modeling these probabilities. In particular, we do not need to assume that the judge or jury are Bayesian. However, we provide a Bayesian micro-foundation for these probabilities in Appendix A.

²¹For any a, b, c, d of \mathbb{R}_{++} we have $\min\{a/b, c/d\} \leq (a+c)/(b+d) \leq \max\{a/b, c/d\}$, with strict inequalities if $a/b \neq c/d$, a generic condition which we will assume throughout (it is easy to impose conditions to guarantee it: for example, one can rank bodies of evidence in terms of the posterior that they generate, as in Appendix A).

welfare is given by

$$\begin{aligned} \mathcal{W}_3(s_1, s_2) = & \lambda [\pi_g^1 W(s_1, g) + \pi_g^2 W(s_2, g) + (1 - \pi_g)W(0, g)] + \\ & (1 - \lambda) [\pi_i^1 W(s_1, i) + \pi_i^2 W(s_2, i) + (1 - \pi_i)W(0, i)]. \end{aligned} \quad (10)$$

Our first result shows that if the sentence associated with a conviction in the two-verdict system is optimal, then there is an improvement that does not increase the sentence.

Proposition 1 1. *Suppose that s^* is the optimal interim sentence in the two-verdict system, i.e., the one that maximizes $\mathcal{W}_2(s)$. Then, there exists an $s_1 < s^*$ such that the interim welfare in the three-verdict system with sentences s_1 and s^* is higher than in the two-verdict system, i.e., $\mathcal{W}_3(s_1, s^*) > \mathcal{W}_2(s^*)$.* 2. *Suppose that s^{**} is the optimal ex-ante sentence in the two-verdict system. Then, there exists an $s_1 < s^{**}$ such that the ex-ante welfare in the three-verdict system with sentences s_1 and s^{**} is higher than in the two-verdict system.*

The proof of Proposition 1, below, shows that the results hold even when the original sentence is not optimal, as long as it is not too suboptimally lenient. Thus, it may be generally possible to improve upon the two-verdict system even under the strong restriction of not harming any defendant more than in the two-verdict system.

Proof. Consider the first part of the proposition. By construction s^* maximizes $\mathcal{W}_2(s)$ with respect to s . In particular, $s \leq \hat{s}$. Since all sentences are interior, s^* must satisfy the first-order condition

$$\lambda \pi_g W'(s^*, g) + (1 - \lambda) \pi_i W'(s^*, i) = 0. \quad (11)$$

Now consider the derivative of $\mathcal{W}_3(s_1, s^*)$ with respect to s_1 , evaluated at $s_1 = s^*$. From (15), we have

$$\left. \frac{\partial \mathcal{W}_3(s_1, s^*)}{\partial s_1} \right|_{s_1=s^*} = \lambda \pi_g^1 W'(s^*, g) + (1 - \lambda) \pi_i^1 W'(s^*, i). \quad (12)$$

Since $\frac{\pi_g^1}{\pi_i^1} < \frac{\pi_g}{\pi_i}$, $W'(s^*, g) > 0$, and $W'(s^*, i) < 0$, the first-order condition (11) implies that the right-hand side of (12) is strictly negative. This shows that decreasing s_1 below s^* strictly improves welfare, yielding the desired improvement.

For the second part of the proposition, differentiating (4) gives the first-order condition satisfied by s^{**} in the two-verdict system:

$$\begin{aligned} \frac{dH(s^{**})}{ds} [\eta_g (\pi_g W(s^{**}, g) + (1 - \pi_g)W(0, g)) + \eta_i (\pi_i W(s^{**}, i) + (1 - \pi_i)W(0, i)) - h] \\ + H(s^{**}) (\eta_g \pi_g W'(s^{**}, g) + \eta_i \pi_i W'(s^{**}, i)) = 0. \end{aligned} \quad (13)$$

The equivalent derivative for the three-verdict case with respect to s_1 at s^{**} is

$$\begin{aligned} \frac{dH_1(s^{**})}{ds_1} & [\eta_g (\pi_g W(s^{**}, g) + (1 - \pi_g)W(0, g)) + \eta_i (\pi_i W(s^{**}, i) + (1 - \pi_i)W(0, i)) - h] \\ & + H(s^{**}) (\eta_g \pi_g^1 W'(s^{**}, g) + \eta_i \pi_i^1 W'(s^{**}, i)), \end{aligned} \quad (14)$$

where $H_1(s_1)$ denotes the fraction of individuals who commit the crime in the three-verdict system as a function of the punishment s_1 when $s_2 = s^{**}$ (in particular, $H_1(s^{**}) = H(s^{**})$).

Dividing (13) by π_g we obtain

$$\begin{aligned} \frac{1}{\pi_g} \frac{dH(s^{**})}{ds} & [\eta_g (\pi_g W(s^{**}, g) + (1 - \pi_g)W(0, g)) + \eta_i (\pi_i W(s^{**}, i) + (1 - \pi_i)W(0, i)) - h] \\ & + H(s^{**}) \left(\eta_g W'(s^{**}, g) + \eta_i \frac{\pi_i}{\pi_g} W'(s^{**}, i) \right) = 0. \end{aligned}$$

Dividing (14) by π_g^1 we obtain

$$\begin{aligned} \frac{1}{\pi_g^1} \frac{dH_1(s^{**})}{ds_1} & [\eta_g (\pi_g W(s^{**}, g) + (1 - \pi_g)W(0, g)) + \eta_i (\pi_i W(s^{**}, i) + (1 - \pi_i)W(0, i)) - h] \\ & + H(s^{**}) \left(\eta_g W'(s^{**}, g) + \eta_i \frac{\pi_i^1}{\pi_g^1} W'(s^{**}, i) \right). \end{aligned}$$

From (3) and (7) we have

$$\frac{1}{\pi_g} \frac{dH(s^{**})}{ds} = \frac{1}{\pi_g^1} \frac{dH_1(s^{**})}{ds_1}.$$

Since $W'(s^{**}, i) < 0$, to prove that the derivative for the three-verdict case is negative, which would demonstrate the claim in the statement of the proposition, it suffices that $\pi_i^1/\pi_g^1 > \pi_i/\pi_g$, which holds by definition of verdict 1. ■

Proposition 1 shows the optimal two-verdict system can robustly be improved without increasing the punishment of any defendant, i.e., without increasing the coercive power of the state.

Our next result shows that welfare improvements exist even if the initial binary verdict was not optimal. Although the improvement we construct requires a higher sentence in one case, the probability of conviction is unchanged and the improvement is robust in a strong sense: welfare increases conditional on facing each type of the defendant, as stated below.

Proposition 2 *For any sentence $s > 0$ in the two-verdict system and any verdict technologies π_i, π_g, π_i^j , etc., there exists a three-verdict system with sentences s_1 and s_2 in which the interim*

welfare is higher than in the two-verdict system, i.e., $\mathcal{W}_3(s_1, s_2) > \mathcal{W}_2(s)$. Moreover, the welfare is higher conditional on the defendant being innocent and conditional on the defendant being guilty. If $s \leq \hat{s}$, then $s_1 < s < s_2$.

One key aspect of Proposition 2 is that it applies to all two-verdict systems, even those with a suboptimal sentence $s > 0$, and for any way of splitting of the conviction probabilities π_i and π_g . In particular, it applies whether s was chosen with an ex-ante or an interim welfare perspective in mind. Another key aspect of Proposition 2 is that the three-verdict system does not increase the probability of punishing the innocent relative to the two-verdict system. Instead, it modifies the sentence to reflect the richer information that verdicts 1 and 2 convey regarding the relative likelihood of the defendant being guilty or innocent.

Lando (2005) derives a similar result to Proposition 2 when the defendant is risk neutral with respect to the sentence and the prior about the defendant's type is uniform.²²

Proof. If $s > \hat{s}$, then setting $s_1 = s_2 = \hat{s}$ suffices, since the ex-post welfare $W(s, i)$ and $W(s, g)$ decreases in $s > \hat{s}$. Suppose that $s \leq \hat{s}$. First, observe that $\mathcal{W}_3(s, s) = \mathcal{W}_2(s)$: if we give verdicts 1 and 2 the sentence associated with the guilty verdict of the two-verdict case, then we clearly obtain the same welfare as in the two-verdict case. We are going to create a strict welfare improvement by slightly perturbing the sentences s_1 and s_2 . Consider any small $\varepsilon > 0$ and let $s_1 = s - \varepsilon$ and $s_2 = s + \varepsilon\gamma$. The welfare impact of this perturbation is

$$\mathcal{W}_3(s_1, s_2) = \mathcal{W}_2(s) + \lambda\varepsilon W'(s, g)(-\pi_g^1 + \gamma\pi_g^2) + (1 - \lambda)\varepsilon W'(s, i)(-\pi_i^1 + \gamma\pi_i^2) + o(\varepsilon), \quad (15)$$

where W' denotes the derivative of W with respect to its first argument. Since $W(\cdot, i)$ is decreasing, $W'(s, i)$ is negative. Similarly, because $s \leq \hat{s}$ and $W(\cdot, g)$ is increasing on that domain, $W'(s, g)$ is positive. Since $\pi_g^1/\pi_g^2 < \pi_i^1/\pi_i^2$, we can choose γ between these two ratios. Doing so guarantees that $-\pi_g^1 + \gamma\pi_g^2$ is positive and $-\pi_i^1 + \gamma\pi_i^2$ is negative, which shows the claim. ■

Proposition 2 considers interim social welfare, after the crime has taken place. The incentives to commit the crime may *a priori* be influenced by the introduction of a third verdict, as (6) indicates. The proof of Proposition 2 shows that the welfare-improving sentences in the three

²²The uniform prior assumption is our interpretation of Lando's formulas, which are non Bayesian, but which Lando microfounds in a special case in an appendix. Lando does not investigate the robustness properties of the improvement.

verdict system can in fact be chosen in a way that does not increase the set of individuals who commit the crime. To see this, recall that the range of welfare-improving ratios γ for $s < \hat{s}$ is $[\pi_g^1/\pi_g^2, \pi_i^1/\pi_i^2]$, which is independent of the function $W(\cdot, g)$. For any $s > 0$, choosing $\gamma = \pi_g^1/\pi_g^2$ would not change, to a first order, the welfare for a guilty defendant, and would increase the welfare for an innocent defendant. Replacing $W(\cdot, g)$ with the individual's utility function $u(\cdot)$ and setting $\gamma = \pi_g^1/\pi_g^2$ would make a guilty defendant indifferent between the two- and three-verdict systems, so the left-hand side of (6) would not change. Therefore, the three-verdict system would deter all the individuals deterred by the two-verdict system.²³

This observation immediately implies the following corollary of Proposition 2.

Corollary 1 *For any sentence $s > 0$ of the two-verdict system and any verdict technologies π_i, π_g, π_i^j , etc., there exists a three-verdict system with sentences s_1 and s_2 in which the set of individuals who commit the crime is no larger, and the interim and ex-ante welfare is strictly higher, than in the two-verdict system.*

While the improvement in Proposition 2 does not increase the probability of punishing an innocent defendant (or a guilty one), an erroneously convicted defendant may face a harsher sentence ex-post when $s_2 > s$.

4 Posterior-Based Verdicts

4.1 The Bayesian conviction model

The analysis of Section 3 did not impose any structure on how verdicts were determined. We now show how to specialize the setting to a class of verdicts based on the posterior probability that the defendant is guilty. We will use this in the remainder of the section. Starting with the prior probability $\lambda = \eta_g/(\eta_g + \eta_i)$, the trial generates evidence that is used to form the posterior. This is summarized by distributions $F(\cdot|g)$ and $F(\cdot|i)$, which describe the posterior based on

²³The utility of an innocent defendant would increase, so even more individuals would be deterred if the individual took into account the negligible probability he would be charged with the crime if he didn't commit it.

whether the defendant is actually guilty or innocent.²⁴ For expositional convenience, we assume that $F(\cdot|g)$ and $F(\cdot|i)$ have positive densities $f(\cdot|g)$ and $f(\cdot|i)$.

In a two-verdict system based on the defendant's posterior, it is natural to follow a cut-off rule. Appendix A shows that any "reasonable" verdict rule based on evidence in the two-verdict system can be formalized as a Bayesian model with posterior cut-off rule. If the posterior p is below a threshold p^* , then the defendant is acquitted, receiving a sentence of $s = 0$. If p exceeds p^* , then the defendant receives a sentence $s^* > 0$. The cutoff rule is a particular case of the previous section, with $\pi_g = Pr[p > p^*|g] = 1 - F(p^*|g)$ and $\pi_i = 1 - F(p^*|i)$.

The interim social welfare is given by

$$\begin{aligned} \mathcal{W}_2(p^*, s^*) = & \lambda [(1 - F(p^*|g))W(s^*, g) + F(p^*|g)W(0, g)] + \\ & (1 - \lambda) [(1 - F(p^*|i))W(s^*, i) + F(p^*|i)W(0, i)]. \end{aligned} \quad (16)$$

Similarly, the ex-ante social welfare is given by

$$H(s^*) (\eta_g ((1 - F(p^*|g))W(s^*, g) + F(p^*|g)W(0, g)) + \eta_i ((1 - F(p^*|i))W(s^*, i) + F(p^*|i)W(0, i)) - h). \quad (17)$$

In what follows, we will denote by (p^*, s^*) the cutoff and sentence used in the two-verdict system. These variables may be chosen to maximize (16) or (17). In that case, they correspond to the interim or ex-ante utilitarian optimum for the two-verdict case.

4.2 Multi-verdict systems

Our analysis can be extended to more than three verdicts. Granted an arbitrary number of verdicts, from an interim perspective one would wish to associate with each posterior belief p the sentence $s(p)$ that maximizes the welfare objective

$$pW(s, g) + (1 - p)W(s, i) \quad (18)$$

with respect to s . Rewriting the objective function as

$$\mathcal{W}(p, s) = p[W(s, g) - W(s, i)] + W(s, i),$$

²⁴In order to match the prior λ , the distributions must satisfy the conservation equation

$$\lambda = E[p] = \lambda \int_0^1 p dF(p|g) + (1 - \lambda) \int_0^1 p dF(p|i).$$

we notice that it is supermodular in (p, s) .²⁵ This implies that the selection of maximizers of (18) is isotone. In particular, there exists a nondecreasing selection $s(p)$ of optimal sentences. The same is true when choosing sentences to maximize the ex-ante welfare.

The arguments used for Propositions 2, Corollary 1, and 1 easily generalize to yield the following results. For $k \geq 2$, we define a k -verdict system by a vector $(p_0, s_0, p_1, s_1, \dots, p_{k-1}, s_{k-1})$ of strictly increasing cutoffs and sentences, with $p_0 = 0$, $p_{k-1} < 1$, $s_0 = 0$ and $s_{k-1} \leq \hat{s}$. In this system, a defendant receives sentence $s_{k'}$ whenever his posterior p lies in $(p_{k'}, p_{k'+1})$.

Proposition 3 *Suppose that the posterior distributions are continuous for both the guilty and innocent defendants. Then, for any k -verdict system there is a $k + 1$ verdict system that strictly increases ex-ante and interim welfare. Moreover, if a k -verdict system is optimal among all k -verdict systems and $k \geq 2$, then there is a $k + 1$ -verdict system that strictly increases ex-ante and interim welfare and in which any defendant receives a weakly lower sentence.*

5 Intermediate “not guilty” verdict

Suppose now that those defendants who would be acquitted in the current two-verdict system now receive one of two verdicts, which we denote 1 and 2. Both verdicts are associated with no jail time, i.e., with $s = 0$. Verdict 1, which we refer to as “not guilty,” obtains if the posterior is less than some cutoff $p^{iv} < p^*$, where p^* is the threshold for conviction, and verdict 2, which we refer to as “not proven,” obtains if the posterior is between p^{iv} and p^* . We denote by p_i the probability that a defendant is guilty conditional on verdict $i = 1, 2$. A conviction leads to the same sentence s^* as in the two-verdict system.

We assume that society observes the verdict at the end of the trial, but not the posterior regarding the defendant’s guilt. The stigmatization associated with being charged and tried is modeled by a cutoff p^s , such that the defendant is stigmatized if the probability he is guilty conditional on the verdict exceeds p^s . We take p^s as exogenous, and assume that convicting a defendant is more demanding than stigmatizing him, so $p^s < p^*$.²⁶ We also assume that if

²⁵ $W(s, g)$ increases in s over the relevant range $[0, \hat{s}]$ while $W(s, i)$ is decreasing in s . This implies that $\partial W / \partial p = W(s, g) - W(s, i)$ increases in s and, hence, supermodularity of $\mathcal{W}(p, s)$. See Milgrom and Shannon (1994).

²⁶This implies that the analysis of Section 3 does not change as a result of the stigma, since a defendant who receives verdicts 1 or 2 is stigmatized.

the defendant is completely cleared in the trial and the public were fully aware of this, then he would not be stigmatized. That is, $\underline{p} < p^s$, where \underline{p} is the lowest possible posterior. An innocent defendant who is stigmatized lowers welfare by $d^i > 0$, and a guilty defendant who is stigmatized increases welfare by $d^g > 0$.²⁷ We are interested in the optimal cutoff p^{iv} and the conditions under which introducing the additional verdict increases welfare. For expositional simplicity we will consider interim welfare; the same qualitative results hold for ex-ante welfare.

The relevant part of the welfare function in the two-verdict system is

$$\lambda [W(0, g) + 1_{p^{ng} > p^s} d^g] + (1 - \lambda) [W(0, i) - 1_{p^{ng} > p^s} d^i],$$

where p^{ng} is the probability that a defendant is guilty conditional on being acquitted, since whether an acquitted defendant is stigmatized depends on whether p^s is lower or higher than p^{ng} . We consider these two possibilities below.

Suppose first that $p^{ng} \geq p^s$, so an acquitted defendant in the two-verdict system is stigmatized. For any p^{iv} , it must be that $p_2 \geq p^{ng} \geq p^s$, so the defendant is stigmatized if he is found “not proven” in the three-verdict system. The split can have an effect on social welfare only if $p_1 \leq p^s$, in which case the defendant is not stigmatized if he is found “not guilty” in the three-verdict system. Therefore, consider p^{iv} such that $p_1 < p^s$. Eliminating the stigma when the defendant is found “not guilty” increases the relevant part of the welfare function by

$$-\lambda \sum_{p \leq p^{iv}} f(p|g) d^g + (1 - \lambda) \sum_{p \leq p^{iv}} f(p|i) d^i.$$

For a given posterior $p \leq p^{iv}$ the increase is

$$-\lambda f(p|g) d^g + (1 - \lambda) f(p|i) d^i > 0 \iff \frac{f(p|g)}{f(p|i)} < \frac{(1 - \lambda) d^i}{\lambda d^g}. \quad (19)$$

Since $f(p|g)/f(p|i)$ increases in the posterior p , a fact we show in Appendix A.1, we obtain the following result.

Proposition 4 *Suppose that being acquitted in the two-verdict system carries a stigma. Then, optimally splitting the acquittal into “not guilty” and “not proven” increases interim welfare if and only if $\frac{f(\underline{q}|g)}{f(\underline{q}|i)} < \frac{(1-\lambda)d^i}{\lambda d^g}$.*

If the condition in Proposition 4 holds, then the optimal cutoff p^{iv} is the minimum between the highest posterior for which (19) holds and the highest posterior such that $p_1 \leq p^s$. Notice

²⁷A similar analysis can be conducted for $d^i \leq 0$ and/or $d^g \leq 0$.

that the condition in Proposition 4 is satisfied more easily if the defendant is more likely to be innocent (λ decreases), the stigma for the innocent increases, or the stigma for the guilty decreases.

Now suppose that $p^{ng} < p^s$, so an acquitted defendant in the two-verdict system is not stigmatized. The split can have an effect on social welfare only if $p_2 > p^s$, in which case the defendant is stigmatized if he is found “not proven” in the three-verdict system. Therefore, consider p^{iv} such that $p_2 > p^s$. Stigmatizing the defendant when he is found “not proven” increases the relevant part of the welfare function by

$$\lambda \sum_{p > p^{iv}} f(p|g) d^g - (1 - \lambda) \sum_{p > p^{iv}} f(p|i) d^i.$$

For a given posterior $p > p^{iv}$ the increase is

$$\lambda f(p|g) d^g - (1 - \lambda) f(p|i) d^i > 0 \iff \frac{f(p|g)}{f(p|i)} > \frac{(1 - \lambda) d^i}{\lambda d^g}. \quad (20)$$

Since $f(p|g)/f(p|i)$ increases in the posterior p , we obtain the following result.

Proposition 5 *Suppose that being acquitted in the two-verdict system does not carry a stigma. Then, optimally splitting the acquittal into “not guilty” and “not proven” increases interim welfare if and only if $\frac{f(p^*|g)}{f(p^*|i)} > \frac{(1-\lambda)d^i}{\lambda d^g}$.*

If the condition in Proposition 5 holds, then the optimal p^{iv} is the maximum between the lowest posterior for which (20) holds and the lowest posterior such that $p_2 \geq p^s$. Notice that the condition in Proposition 5 is satisfied more easily if the defendant is more likely to be guilty (λ increases), the stigma for the innocent decreases, or the stigma for the guilty increases.

6 Value of evidence with a third verdict

The previous sections have taken as given the technology that generates evidence regarding the defendant’s guilt. Gathering evidence is costly, however, and the amount of evidence generated in a case depends on the incentives of the agents involved in the evidence-gathering process: law enforcement officers, prosecutors, experts, etc.

Setting aside the possible biases in these agents’ behavior, the socially optimal amount of evidence to be gathered in a case clearly depends on the verdict structure. For example, a trial

system in which a single verdict is given regardless of the evidence produced clearly eliminates any value of gathering evidence.

This section investigates the impact on evidence gathering of introducing a third verdict. For simplicity, we focus on the setting of Section 3 with the Bayesian conviction model of Section 4.1. We consider interim welfare, since for any particular crime investigated it is plausible that the agents involved in the discovery stage of the trial are primarily concerned with the facts pertaining to the specific case and not with deterrence.

A (possibly multi-) verdict system leads to welfare

$$w(p) = pW(s(p), g) + (1 - p)W(s(p), i), \quad (21)$$

where $p \mapsto s(p)$ is a step function that starts at zero, has two levels in a two-verdict system, and three levels in a three-verdict system. The welfare function $w(p)$ is piecewise linear. It starts at 0, and decreases until a kink at which the sentence jumps from 0 to a positive level. Figure 1 represents the welfare function for the optimal two-verdict system when $W(\cdot, g)$ and $W(\cdot, i)$ are quadratic, for parameters given in the appendix.

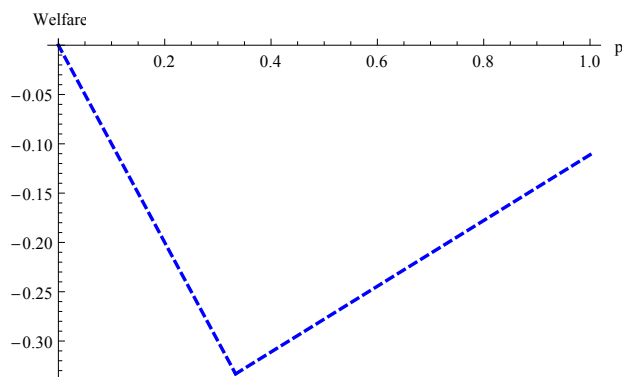


Figure 1: Welfare function, 2 verdicts.

The kink occurs at the cutoff $p^* = 1/3$, at which the sentence jumps from 0 to $2/3$. Figure 2 represents the welfare function for the optimal three-verdict system obtained by adding an intermediate verdict and keeping the highest sentence at $2/3$. The first cut-off is $p_1 = p^* = 1/3$, and the second cut-off is $p_2 = 1/2$. The welfare function is discontinuous at p_1 : this reflects the fact that p_1 is not chosen optimally, but is rather “inherited” from the two-verdict system. In contrast, because p_2 is chosen optimally, the welfare function is kinked but continuous at p_2 .

Actual evidence formation processes are complex, involving various actors of different types – forensic experts, lawyers, witnesses—and different forms of evidence. To model evidence

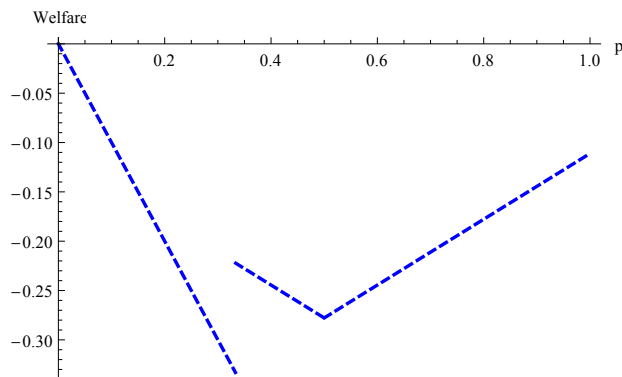


Figure 2: Welfare function, 3 verdicts.

formation, we must abstract from much of this complexity. Instead, we take the viewpoint of a social planner who may gather information until a verdict is reached.

The tradeoff at the heart of this task is clear: effort spent gathering evidence is costly, but provides information about the defendant’s guilt. We discuss two ways to model this tradeoff (there are, of course, many others). The first is a one-shot binary evidence-gathering decision, which already captures the rough intuition for why two-verdict and three-verdict systems differ in their effects on evidence gathering. The second is a continuous evidence-gathering process, which provides a more visually appealing representation of the impact of a third verdict on evidence gathering.

6.0.1 One-shot evidence gathering

Suppose the planner decides whether to gather evidence, at a cost of $c > 0$. Starting with a prior p_0 , the evidence returns a higher probability of guilt, say $p_0 + \Delta$ with probability $1/2$, and a lower probability $p_0 - \Delta$ also with probability $1/2$. The belief process is a martingale: the mean of the posterior p' is equal to the prior: $\frac{1}{2}(p + \Delta) + \frac{1}{2}(p - \Delta) = p$.

When is evidence gathering socially desirable? Suppose first that the prior is close to 0, so that the posterior p' surely lies below the cutoff p_1 . Then, the additional evidence has no value as the defendant will be acquitted in all cases. Similarly, if p_0 is high enough for p' to lie above the cutoff p_1 no matter what, the additional evidence has no value as the defendant will be convicted regardless of p' . For p_0 slightly below p_1 and Δ such that $p_0 + \Delta$ lies above p_1 , the value of evidence is positive, since it can lead to a conviction and increase welfare (relative to an acquittal) when it does. Similarly, evidence is valuable for p_0 slightly above p_1 . Thus, evidence

is valuable around the kink, where the welfare function is convex.

Consider now the case of three verdicts. For p_0 slightly below p_1 , the value of evidence is higher than in the two-verdict case because a positive belief update triggers a large improvement in welfare (see Figure 2). For p in a neighborhood of p_2 , the value of evidence is also positive due to the convex kink there, whereas it is 0 (for Δ small enough) in the two-verdict case.

For p_0 slightly above p_1 however, additional evidence may be more valuable in the two-verdict case, which creates a “doughnut hole:” additional evidence is more valuable in the three-verdict case than in the two-verdict case for more extreme beliefs, and less valuable in some intermediate region. This result is easier to visualize in the next model, where evidence gathering is more gradual.

6.0.2 Continuous evidence gathering

Now suppose that evidence is gathered continuously. As long as evidence is gathered, a flow cost of c is incurred. During this time the belief p_t that the defendant is guilty evolves continuously in a way that is consistent with Bayesian updating, so that the closer the belief is to 0 or 1 the more slowly it evolves. Let the value function $v(p)$ denote the social welfare that arises from stopping optimally when the current belief is p . If it is optimal to stop immediately, then $v(p)$ coincides with $w(p)$ given by (21). Otherwise, it is optimal to continue collecting evidence, which leads to p changing continuously. In this case, $v(p)$ is the expectation of the value of stopping optimally in the future.

As in the one-shot setting, the value of gathering evidence in the two-verdict setting (formally analyzed in Appendix B) is high when the belief is very close to the convex kink p_1 in Figure 1. For beliefs slightly farther from the kink the value is still high because collecting evidence there leads with high probability to beliefs that are closer to the kink, for which the value of gathering evidence is high. The continuous evidence collection structure smooths the value of evidence collection as a function of the belief. This value decreases continuously as the belief moves away from the kink. For beliefs that are sufficiently far from the kink, it is optimal to stop immediately. Less informative evidence and higher costs of evidence collection lead to lower value of evidence collection, which correspond to functions v with to lower values. This is depicted in Figure 3 for parameters given in the appendix.

Consider now the case of three verdicts in Figure 2. Around the kink p_2 , the value function v

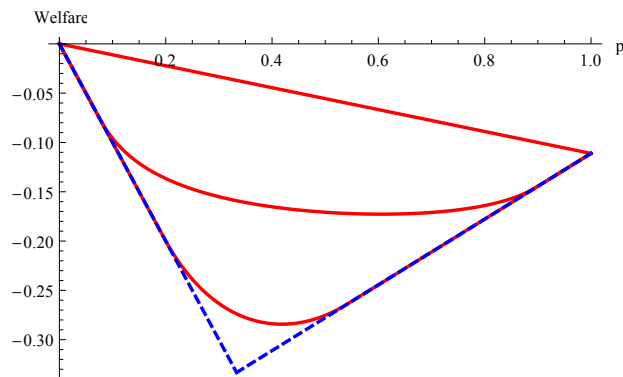


Figure 3: Value function, 2 verdicts, for varying cost levels.

behaves similarly to the two-verdict case. Around the discontinuity p_1 the behavior is different, just like in the one-shot setting. Immediately to the left of the discontinuity the value of collecting evidence is high, so v substantially exceeds w , and this value decreases continuously for lower beliefs. Immediately to the right of the discontinuity there is no value in collecting evidence, so v coincides with w . Higher beliefs have a positive value of evidence collection, because of the kink. This value continuously decreases to 0 as the beliefs increase above the kink. This is depicted in Figure 4.

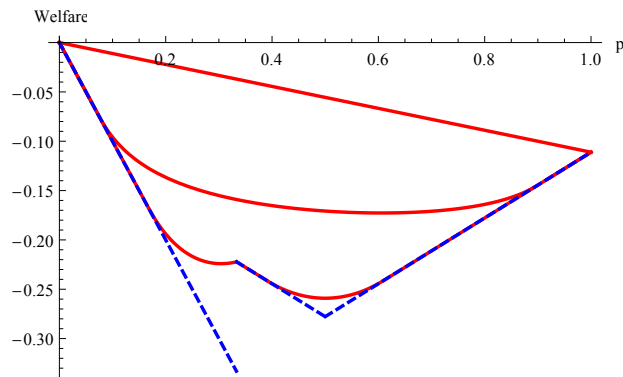


Figure 4: Value function, 3 verdicts, for varying cost levels.

In conclusion, the impact of switching to a three-verdict system by splitting the guilty verdict depends on how informative the evidence is and how costly it is to collect. When evidence is very informative, the posterior is unlikely to end up in the middle region, so the intermediate verdict has little impact. When finding new evidence is very costly, however, the posterior may end up in the middle region. The third-verdict system then increases the value of gathering evidence in two regions, below p_1 and around p_2 , and decreases the value immediately above p_1 . Overall, because $\tilde{p}_0 < \hat{p}_1$ and $\tilde{p}_2 > \hat{p}_2$, the three-verdict system results in evidence gathering at

more extreme beliefs, where in the two-verdict evidence gathering has already stopped.

7 Reflecting residual doubt in the current justice system

Numerous scholars including Fisher (2012), Lando (2005) have already noted the existence of ad hoc probabilistic sentencing in existing criminal justice systems. The most explicit inclusion of residual doubt in the U.S. criminal justice system concerns the determination of death sentences. In capital cases, juries must decide, after returning a guilty verdict, whether the defendant should get the death penalty. In this penalty phase, residual or “lingering” doubt may be used as a mitigating circumstance to reject the death penalty.²⁸ The Capital Jury Project—an academic survey of past jurors in capital cases—has found that lingering doubt was the most important mitigating factor identified by jurors.

There is, however, wide variation in how residual doubt is applied. First, the U.S. penal code (Title 18, §3592) does not explicitly mention residual doubt in its list of mitigating factors, although it does state that mitigating circumstances are not limited to this list. In some cases, jurors are not informed that lingering doubt is a valid mitigating circumstance.²⁹ In *Franklin v. Lynaugh* (1988), the U.S. Supreme Court rejected a defendant’s right to invoke residual doubt at the penalty stage, while in *People v. McDonald* (Supreme Court of Illinois, 1995) a trial judge refused to answer jurors’ question on the issue, a decision which was later affirmed by the Supreme Court of Illinois.

Compounding this inconsistency, there is empirical evidence that many jurors get confused with the voting rules used to establish aggravating and mitigating circumstances at the penalty stage. While the unanimity rule is required to find a circumstance aggravating, no such standard exists for mitigating circumstances. The Capital Jury Project found, however, that 45% of jurors failed to understand that they were allowed to consider any mitigating evidence during the sentencing phase of the trial, not just the factors listed in the instructions.³⁰

²⁸The more demanding requirement of proving guilty beyond “all doubt” has been discussed in some states, such as the bill proposed in 2003 by then Illinois House Republican leader Tom Cross. Some death penalty advocates have countered that it was impossible to prove anything beyond *all* doubt, and that the bill would in effect rule out the death penalty. Various degrees of lingering doubt have been discussed (e.g., Sand and Rose, 2003) without any mathematical formalism.

²⁹See, e.g., *People v. Gonzales and Soliz*, California Supreme Court, 2009.

³⁰The CJP’s findings concerning jurors’ understanding of instructions are summarized at

When sentencing is performed by a trial judge, the invocation of residual doubt can be highly controversial. In *State v. Krone* (Arizona Supreme Court, 1995) a trial judge sentenced to life in prison a defendant found guilty of murder, citing doubt about whether he was the true killer. In their legal textbook, Dressler and Thomas (2010, pp. 57–61) comment that this decision “borders on the unbelievable.” They do not, however, suggest an alternative solution.

In non-capital cases, only five states permit juries to make the sentencing decision. Outside of these states, residual doubt can thus only be expressed by the sentencing judge, whose opinion does not necessarily reflect the views of the jury. Again, residual doubt is not listed as a mitigating factor in sentencing guidelines.

The fact that residual doubt should only be considered in capital cases seems largely arbitrary.³¹ Even comparably less serious cases can carry long sentences, resulting in extreme punishments for defendants who are found guilty but for whom residual doubt remains. For example, in *State v. May* (Arizona Superior Court, 2007) a thirty-five-year-old defendant was sentenced to 75 years in jail after being found guilty of touching, in a residential swimming pool, the clothing of four children in the vicinity of their genitals (Nelson, 2013). Jurors had doubts about the guilt of the defendant: they were twice unable to reach a verdict within the first three days of deliberation. The explicit inclusion of residual doubt in sentencing would have likely avoided such an extreme outcome.

Some felonies provide an indirect way of expressing doubt by using the lesser-included-offense rule: juries can return a manslaughter verdict, rather than a first- or second-degree murder verdict, or a larceny verdict instead of a robbery verdict. However, each of these verdicts corresponds to a precise charge (e.g., whether premeditation and malice aforethought were involved) and doubt about a particular charge can only be imperfectly expressed by returning a guilty verdict on a lower count. These instruments only offer a limited and, in fact, improper, way of reflecting residual doubt. Furthermore, the less-included-offense rule is not a constitutional right of the defendant; its application is therefore to some extent arbitrary and depends on the

<http://www.capitalpunishmentincontext.org/issues/juryinstruct>.

³¹Capital sentences are unique in their irreversibility, which creates an additional reason for avoid this sentence in case of lingering doubt: exonerating evidence may appear after the execution of the defendant, preventing any release and compensation. In practice, however, this fundamental difference is attenuated by the fact that death-row defendants spend many years in jail before their execution until all recourses have been exhausted, while non-capital defendants serving long sentences may die in jail, which also prevents any release or compensation.

inclination of the jury (see Mascolo (1986)).

Even when the lesser-included-offense rule does not apply, residual doubt may be reflected by returning a guilty verdict only on a subset of the charges brought against the defendant. There is anecdotal evidence that such compromise is sometimes used by the jurors to reflect doubt. In the aforementioned *State v. May*, for instance, Nelson (2013) notes that “it seems likely that the defendant molested either all of the children or none of them. So why did the jury ultimately reach a verdict of guilty on five counts and not guilty on two? The answer is that the jurors compromised.” Dropping some charges is, however, a very coarse instrument to incorporate residual doubt: for example, this approach cannot be used to reduce the sentence of a defendant facing a single but severe count, while it may be used for a defendant facing several counts, the sum of which adds to the same aggregate maximal sentence as in the single-count case.³² Even when it is feasible, the approach exposes the defendant to another idiosyncratic component of the jury—whether it is sophisticated or willing enough to use this compromise strategy—introducing a source of jury heterogeneity in trial outcomes even for otherwise identical cases.³³

The U.S. justice system incorporates residual doubt about a defendant’s guilt in two other ways. First, a defendant found not guilty in a criminal trial may still be found guilty in a civil suit, which uses the less demanding preponderance-of-evidence standard of proof. However, civil suit sentences carry no jail time and thus may be more limited in preventing recidivism. Furthermore, the connection between criminal and civil trials is generally limited, preventing any coordination and coherent decision across these trials. Second, residual doubt variations also imply different likelihoods of post-trial events such as successful appeals and exonerations, which affect the defendant’s ultimate punishment. These events are largely beyond the control of the first court and are not a close substitute for the additional verdicts introduced here.

In summary, the current criminal justice system includes various ways of reflecting residual doubt in outcomes and it appears that these ways are used purposefully by some actors of the system. However, these ways are largely arbitrary, inconvenient, and uncoordinated. This paper

³²The set of charges leveled at the defendant may also be affected by the strategic decisions of the prosecutor, which increases the prosecutor’s power and adds to the complexity of this problem.

³³It should also be noted that under the current law, such compromise is actually illegal if it results from a bargaining between pro-acquittal and pro-conviction jurors. Such an arrangement currently violates the rights of the defendant if the pro-acquittal jurors still believe that the defendant should be found not guilty (Mascolo (1986)).

proposes a structured, systematic approach for the consideration of residual doubt in criminal justice decisions and explicit designs which are shown to improve welfare in many settings.

8 Implementation and jurors' reactions to additional verdicts

Implementation: verdicts vs. sentences

Formalizing the intermediate sentence introduced in this paper as an intermediate *verdict* is consistent with the not-proven verdict, discussed in Section 5, used by some criminal justice systems. In this formulation, the jury must decide, according to some collective rule, among the three verdicts.

An alternative “two-step” implementation maintains the current separation between the fact-finding and sentencing stages. The verdict outcome is still binary (“guilty” or “not guilty”), and residual doubt is expressed in the form of intermediate sentences decided in the sentencing stage.

The second implementation presents a significant advantage: in principle, the jury can be given exactly the same instructions as in the current system, which allows to cleanly split the set of cases which would receive a “guilty” verdict under the current system into multiple sentence levels reflecting the strength of evidence, and thus leaves unchanged the probability of acquitting the defendant.

Intermediate sentences can be decided in a variety of ways, which may involve a sentencing judge, sentencing guidelines (e.g., automatically rule out the death penalty if the evidence is solely based on a confession), or a jury.

Regardless of the implementation, a potential concern is how the jury may react to additional verdicts. The remainder of our discussion focuses on this issue.

Jurors' reaction to additional verdicts

Jury decisions involve collective and psychological considerations: jurors may have limited and uneven ability to understand jury instructions or interpret the evidence, have varied tolerance for erroneous convictions and acquittals, and are subject to individual biases and to persuasion and group-think dynamics, to cite only a few issues. Even abstracting from these

issues, jury decisions are difficult to analyze.³⁴

The literature on criminal trial design varies from fully rational to completely reduced-form models of jury behavior. At the most “rational” extreme, Lee (2015) considers jurors who perfectly take into account how prosecutors select the pool of defendants who go to trial. Prosecutors can influence this pool by choosing the plea sentence that they propose to defendants before the trial.³⁵ Other papers on trial design (Kaplow (2011), Daughety and Reinganum (2015a,b), Da Silveira (2015), Silva (2018)) abstract from any jury decision, focusing on reduced-form thresholds or on a mechanism design approach without jurors.

A key observation is that our Propositions 1 and 2 continue to hold under the two-step implementation mentioned above, provided that jurors are given the same instructions as in the current system to decide between the guilty and not-guilty verdicts, and react to these instructions in the same way, no matter how imperfect, as they currently do. No matter how “tough on crime” or otherwise biased each juror is, and what voting, persuasion or other collective processes are at play, all these components would play out in exactly the same way at the fact-finding stage, under a standard binary verdict, as in the first step of the two-step approach, guaranteeing that no more defendants are found guilty in the three-verdict system than in the current one.

The main question, therefore, is to what extent jurors would know and incorporate in the fact-finding stage the fact that residual doubt may play a significant role in the sentencing stage.

In practice, there is little evidence that jurors incorporate sentencing considerations into their verdict decisions. On the contrary, in recent history judicial practice has been to keep the jury uninformed about the punishment faced by the defendant (Sauer (1995)). In *United States v. Patrick* (D.C. Circuit, 1974), the court affirmed that the jury’s role is limited to a determination of guilt or innocence. Instructions entirely focus on describing the procedure for finding facts. In many cases—such as *People v. May* above—jurors are unaware of the minimum-punishment guidelines relevant for the case.

³⁴Austen-Smith and Banks (1996), Feddersen and Pesendorfer (1996, 1997), and Gerardi and Yariv (2007) identify important informational effects, which may arise even when all jurors have identical preferences. A central mechanism in this literature is that, conditional on being pivotal in a vote, a rational juror may put so much weight on other jurors’ signals that he significantly discounts, and potentially discards, his own information.

³⁵The approach presumes that jurors are aware of the plea sentence offered to the defendant. In practice, the jury is often instructed to consider only the evidence produced at trial.

There is also empirical evidence that harsher sentences do not result in lower conviction rates. In a study of non-homicide violent case-level data of North Carolina Superior Courts, Da Silveira (2015) finds that the probability of conviction of defendants going to trial in fact increases with the sentence that they face.³⁶ Such a correlation cannot be easily explained away by prosecutor behavior: if, in particular, prosecutors attached more importance to obtaining a conviction when the case is more severe, they would send to trial defendants who are more likely to be found guilty and obtain a guilty plea from the other ones, and one would expect the probability of plea settlements to increase with the severity of the trial sentence. This relation seems contradicted by the data.³⁷

More generally, there is strong evidence that jurors have a limited understanding of the sentences faced by defendants. For example, the aforementioned Capital Jury Project found that most jurors “grossly underestimated” the amount of time spent in jail entailed by a guilty verdict. It is reasonable to believe that jurors would be as unaware of, say, maximum-sentencing guidelines as they currently are of minimum-sentencing guidelines.

Finally, if contrary to expectations jurors incorporated the intermediate verdict into their decision, they might adopt a different standard of proof to convict defendants, knowing that the corresponding cases would result in a different sentence than in the current system. To the extent that jurors did this with the social welfare objective in mind, such a change would likely be beneficial. Jurors may, however, have their own objective in mind. For example, they may ignore, from the interim perspective in which they are placed, the deterrence value, *ex ante*, of higher expected punishments—this issue arises even in a two-verdict system, and may explain the fact, mentioned earlier, that jurors are specifically asked to focus on finding facts and left relatively uninformed about the strength of the punishment implied by a guilty verdict. Jurors may also worry about the length of deliberation, and be willing to continue deliberation only if the social value of doing so is high. The analysis of Section 6 suggests that this value is not lowered by the introduction of an intermediate verdict, and may in fact be higher, for a wide range of beliefs.

³⁶Da Silveira’s analysis excludes the most and least severe cases to focus on a relatively homogeneous pool of cases.

³⁷Elder (1989) finds evidence that circumstances that may aggravate punishment *reduce* the probability of settlement. Similarly, Boylan (2012) finds that a 10-month increase in prison sentences raises trial rates by 1 percent.

A Foundation of the Bayesian Conviction Model

We now study whether actual court proceedings can be translated into a Bayesian updating process and a threshold. We address this by considering an evidence-based trial technology. There is a set X of evidence elements, and “evidence collection” refers to a subset of X . The court technology is a mapping $D : 2^X \rightarrow \{G, N\}$, which for every evidence collection decides whether the defendant is guilty or not guilty.³⁸ Distributions P_θ on 2^X , for $\theta \in \{g, i\}$, describe the probability that different evidence collections arise conditional on the defendant being actually guilty or innocent. We assume that both distributions have full support. Letting π_θ^k denote the probability that a defendant of type θ receive verdict k , we have $\pi_\theta^k = P_\theta(D^{-1}(k))$ for each type θ and verdict k in $\{G, N\}$. Recall that $\pi_i^G < \pi_g^G$, i.e., $P_i(D^{-1}(G)) < P_g(D^{-1}(G))$, and that λ is the prior that the defendant is guilty. We ask several questions.

1. Given D , P_i , P_g , and λ , can D be rationalized as the result of Bayesian updating with a threshold on the posterior for determining guilt? At a minimum, this would require D to respect “incriminating” and “exculpatory” evidence sets, which are determined by whether they indicate that the defendant is more likely to be guilty than innocent.
2. Given D and λ , can P_i and P_g be chosen to rationalize D as the result of Bayesian updating with a threshold on the posterior for determining guilt?
3. Given λ , can D , P_i , and P_g be chosen to rationalize D as the result of Bayesian updating with a threshold on the posterior for determining guilt?

To answer these questions, we formally order defendant types i and g so that $i < g$, and we order verdicts as $N < G$. Then, we say that D **can be rationalized** as the result of Bayesian updating with a threshold on the posterior if for every $E, E' \subseteq X$ we have $D(E) < D(E')$ if and only if the posterior that the defendant is guilty is higher under E' than under E , i.e.,

$$\frac{\lambda P_g(E)}{\lambda P_g(E) + (1 - \lambda) P_i(E)} < \frac{\lambda P_g(E')}{\lambda P_g(E') + (1 - \lambda) P_i(E')}.$$

This condition is equivalent to $\lambda P_g(E) (\lambda P_g(E') + (1 - \lambda) P_i(E')) < \lambda P_g(E') (\lambda P_g(E) + (1 - \lambda) P_i(E))$ and, after rearranging, to

$$\frac{P_g(E)}{P_i(E)} < \frac{P_g(E')}{P_i(E')}.$$

The likelihood ratios are thus ordered independently of λ . For every evidence set $E \subseteq X$, denote by $r(E) = P_g(E) / P_i(E)$ its likelihood ratio. This shows the following proposition.

Proposition 6 *D can be rationalized if and only if for every $E, E' \subseteq X$ the following holds:*

$$r(E) \leq r(E') \Rightarrow D(E) \leq D(E').$$

While we started with a Bayesian definition of rationalizability, this concept is in fact non-Bayesian: it is purely based on the likelihood ratio of guilty given the observed evidence and, in particular, is independent of any prior.

Equipped with this result, we can answer the questions above. For 1, the answer is “yes” if and only if

$$\max \{r(E) : D(E) = N\} < \max \{r(E) : D(E) = G\}. \tag{22}$$

³⁸The analysis can be generalized to stochastic decisions.

For 2, the answer is “yes:” choose P_g and P_i so that (22) holds. Since 2 implies 3, that answer to 3 is also “yes.”

Incriminating and exculpatory evidence: definitions and properties

When D can be rationalized, we say that evidence $e \in X$ is **D -incriminating** if for every $E \subseteq X$ with $e \notin E$, $D(E) = g$ implies that $D(E \cup \{e\}) = g$. We say that evidence $e \in X$ is *P -incriminating* if for every $E \subseteq X$ with $e \notin E$ we have that $r(E) \leq r(E \cup \{e\})$. Decision- and belief-based notions of exculpatory evidence are defined similarly. The following result establishes the logical connection between these concepts.

Proposition 7 *If D is rationalized by P , any P -incriminating evidence is also D -incriminating.*

The reverse need not hold: one can easily construct examples in which some evidence collection E suffices to convict the defendant (i.e., $D(E) = g$) and the additional piece of evidence e reduces the ‘guilt’ ratio ($r(E \cup \{e\}) < r(E)$), but not enough to change the decision ($D(E \cup \{e\}) = g$).

Our definition and characterization of rationalization extend without change to probabilistic functions D , in which the image of D is the probability that the defendant is found guilty.

A.1 Ordering posterior distributions with the MLRP

In the Bayesian conviction model, the posterior belief is formed by combining a prior with the signals observed about the defendant. One may view each evidence collection E as a signal, and signals may be ordered according to the likelihood ratio $r(E)$. The distributions P_i and P_g over evidence collections can then be mapped into distributions over likelihood ratios r . In a Bayesian conviction model, only the likelihood ratio matters for the decision, and one can thus without loss identify any signal with r . Thus, without loss, signals may be ranked according to this likelihood ratio. Let R_g and R_i denote the distributions of r , conditional on being guilty and innocent, respectively. When the signal distributions, conditional on being guilty or innocent, are continuous, let ρ_g and ρ_i denote their densities. By construction, we have $\rho_g(r)/\rho_i(r) = r$. In statistical terms, this means that R_g and R_i are ranked according the MLRP: the ratio of their density is increasing in the signal. Moreover, because the posterior $p(r)$, given a signal r , is equal to the conditional probability of $\theta = g$ given r , it inherits the MLRP.³⁹ Let F_g and F_i denote the distributions of p , conditional on being guilty and innocent, respectively, and let f_g and f_i denote the densities of F_g and F_i (which exist as long as R_g and R_i are continuous), we have $f_g(p)/f_i(p)$ is increasing in p .

Proposition 8 *Suppose that both signal distributions, conditional on being guilty and innocent, are continuous. Then both distributions of the posterior p are continuous, and their density functions satisfy the MLRP.*

This property, which holds without loss (except for the continuity assumption, of a technical nature), plays a key role for subsequent results.

B Modeling continuous evidence gathering

As long as evidence is gathered, the belief p_t that the defendant is guilty evolves as a martingale as in Bolton and Harris (1999):

$$dp_t = Dp_t(1 - p_t)dB_t,$$

³⁹This fact is well-known and straightforward to establish.: if θ is the state of the world, r is a signal, and the conditional distributions $\rho(r|\theta)$ are ranked according to MLRP, then the posterior distributions $\rho(\theta|r)$ are also ranked according to the MLRP.

where B is the standard Brownian motion and D is a measure of the quality of the signal: the higher D is, the faster p evolves toward the true probability that the defendant is guilty (0 or 1). At some time T , the evidence formation process is stopped and the verdict is chosen based on the posterior p_T , which results in social welfare $w(p_T)$.

Adapting the arguments of Bolton and Harris (1999) to our environment, the value function $v(\cdot)$ must satisfy the Bellman equation

$$0 = \max\{w(p) - v(p); -rv(p) - c + \frac{1}{2}D^2p^2(1-p)^2v''(p)\}, \quad (23)$$

where r is a discount rate that captures the idea that longer judicial processes are penalizing for all parties. The first part of the equation implies that $v(p) \geq w(p)$, which means that the value function always exceeds the welfare obtained by stopping immediately. This is natural, since the option of stopping is available at any time. The second part of the equation describes the evolution of the value function while evidence is accumulated:

$$0 = -rv(p) - c + \frac{1}{2}D^2p^2(1-p)^2v''(p).$$

All solutions to this equation are in closed form when $D^2/r = 3/2$:

$$v(p) = -\frac{c}{r} + \left(A_1 + A_2 \left(p - \frac{1}{2} \right) (1-p)^{-2} \right) p^{-\frac{1}{2}}(1-p)^{\frac{3}{2}}, \quad (24)$$

where A_1 and A_2 are free integration constants. For simplicity, in what follows we set $r = 1$ and $D^2 = 3/2$ and vary the cost c .

The region in which evidence is gathered and value functions are determined by the conditions that v is continuous, weakly above w , and when it hits w , it satisfies the smooth pasting property whenever w is continuously differentiable at the hitting point.

Starting with the two-verdict case, one should expect v to coincide with w when p is either close to 0 or close to 1: in this case, there is a high degree of confidence in the defendant's guilt and the value of further evidence gathering is low. Near w 's kink (i.e., the threshold p^* at which the sentence switches), however, the value of additional evidence is high, so v should be strictly above w . Thus, it suffices to connect v and w on both sides of p^* . At the connection points, \hat{p}_1 and \hat{p}_2 such that $\hat{p}_1 < p^* < \hat{p}_2$, v must be equal to w (this is the "value matching" condition) and the derivatives must also coincide (this is the smooth pasting condition).

This imposes four conditions (two value matching and two smooth pasting), and there also four free parameters: the cutoffs \hat{p}_1 and \hat{p}_2 , and the constants A_1 and A_2 arising in equation (24). The result is depicted in Figure 3.

Now consider the three-verdict case. Around the kink p_2 , we still have a two-way smooth connection between w and v , as in the two-verdict case. Around $p_1 = p^*$, however, w is discontinuous, jumping upward from $\underline{w} = -1/3$ to $\bar{w} = -2/9$ as p passes p_1 . In this case, if $v(p_1) > \bar{w}$ (the cost is low), then the situation is exactly as in the two-verdict case. Intuitively, the cost is low enough that the intermediate verdict doesn't matter: evidence is gathered until either the not guilty or the guilty verdict is reached. This a situation in which the trial technology is quite accurate, so a two-verdict system suffices.

For larger costs, however, v hits w exactly at $p_1 = p^*$, due to the upward jump. The smooth pasting condition is violated, because the left derivative of v is higher than its right derivative at p_1 , and v is equal to w on a right neighborhood of p_1 . Intuitively, this kink in the value function reflects the fact that $p_1 = p^*$ was not chosen optimally for the three-verdict system, but rather inherited from the two-verdict system.

The evidence-gathering region now has two parts. When p is below p_1 , there is a large incentive to gather evidence, because such evidence can change the sentence from 0 to s_1 , and s_1 was tailored to provide a fairer

sentence around p_1 than both 0 and s_2 . This also implies that not gathering evidence in a right-neighborhood of p_1 is optimal. The second evidence-gathering region is around p_2 , as before.⁴⁰

Because the first region violates the smooth pasting condition at p_1 , its determination is slightly different. We must determine the threshold \tilde{p}_0 at which the region begins, and we know that the region ends at the cutoff p_1 . At \tilde{p}_0 , we have two conditions: the value matching and the smooth pasting conditions. At p_1 , however, we only have the value matching condition $v(p_1) = \bar{w}$, since the smooth pasting condition is violated. This gives three conditions. There are also three free parameters: the cutoff \tilde{p}_0 and the constants \hat{A}_1 and \hat{A}_2 in (24) for that region. The result is depicted in Figure 4.

Because the welfare w_3 is always higher than the welfare w_2 , it is straightforward to establish that the value function v_3 in the three-verdict case is (weakly) higher than the two-verdict value function v_2 . This matters for high enough cost, i.e., when $v(p_1) = \bar{w}$. In that case, v_3 is strictly above v_2 around p_1 , and it is also strictly above v_2 in the second evidence-gathering region, closer to p_2 . This implies that the cutoff \tilde{p}_0 is lower than the cutoff \hat{p}_1 of the two-verdict case, and the right cutoff \tilde{p}_2 of the second evidence-gathering region in the three-verdict case is greater than \hat{p}_2 .

⁴⁰As the search cost decreases, the two search regions become connected when $v(p_1) \geq \bar{w}$.

C Parameters for the welfare functions of Section 6

We set to 1 the ideal sentence \bar{s} for the guilty and use quadratic loss functions: $W(s, g) = -(1-s)^2$, $W(s, i) = -s^2$. We also assume that the prior is equal to 1/2: the defendant is equally likely to be guilty or innocent ex ante. To obtain simple expressions for the optimal cutoffs and sentences, we reverse-engineer the signal structure. The optimal cutoff p^* is given by the indifference condition

$$p^*W(s^*, g) + (1 - p^*)W(s^*, i) = p^*W(0, g),$$

or $p^*(1 - (s^*)^2) + (1 - p^*)(-s^*)^2 = p^*$. The optimal sentence s^* solves

$$s^* \in \arg \max_s \frac{1}{2}Pr(p \geq p^*|g)W(s|g) + \frac{1}{2}Pr(p \geq p^*|i)W(s|i),$$

which yields the first-order condition, $(1 - F(p^*|g))(1 - s^*) = (1 - F(p^*|i))s^*$, where $F(\cdot|\theta)$ denotes the probability distribution of the posterior p conditional on the defendant's type θ . By choosing $F(\cdot, g)$ and $F(\cdot, i)$ so that the ratio $q = \frac{1 - F(p^*|i)}{1 - F(p^*|g)}$ is equal to 1/2 when evaluated at $p = 1/3$, we verify that $p^* = 1/3$ and $s^* = 2/3$ solve the maximization problem. Note that q must be less than 1, from MLRP.

With three verdicts, we impose the restrictions $p_1 = 1/3$ and $s_2 = 2/3$, so that we are indeed splitting the guilty verdict, without increasing the guilty sentence, and optimize over the remaining two parameters, p_2 and s_1 . These parameters are again characterized by the indifference equation for p_2 , given the sentences s_1 and s_2 that are given above and below p_2 ,

$$p_2W(s_1, g) + (1 - p_2)W(s_1, i) = p_2W(s_2, g) + (1 - p_2)W(s_2, i),$$

and by the optimality condition for s_1 , which is

$$s_1 \in \arg \max_s \frac{1}{2}Pr(p \in [p_1, p_2]|g)W(s|g) + \frac{1}{2}Pr(p \in [p_1, p_2]|i)W(s|i),$$

which yields the first-order condition

$$F([p_1, p_2]|g)(1 - s_1) = F([p_1, p_2]|i)s_1.$$

Using reverse engineering again, we choose $F(\cdot|g)$ and $F(\cdot, i)$ so that the ratio $q' = \frac{F([p_1, p_2]|i)}{F([p_1, p_2]|g)}$ evaluated at $p_1 = 1/3$ and $p_2 = 1/2$ be equal to 2. This inequality guarantees that $s_1 = 1/3$ and $p_2 = 1/2$ satisfy the optimality conditions. Note that the ratio q' must be greater than q , by MLRP.

This yields the welfare functions $w_2(p) = w_3(p) = -p$ for $p < 1/3$, $w_2(p) = -p/9 - (1 - p) \times 4/9$ for $p \geq 1/3$, and $w_3(p) = -p/9 - (1 - p) \times 4/9$ for $p \geq 1/2$, and $w_3(p) = -p\frac{4}{9} - (1 - p)\frac{1}{9}$ for $p \in [1/3, 1/2)$.

References

- ATHEY, S. (2002) "Monotone Comparative Statics under Uncertainty," *Quarterly Journal of Economics*, Vol. 117, pp. 187–223.
- AUSTEN-SMITH, D., BANKS, J. (1996) "Information Aggregation, Rationality, and the Condorcet Jury Theorem," *American Political Science Review*, Vol. 90, pp. 34–45.
- BOLTON, P., HARRIS, C. (1999) "Strategic Experimentation," *Econometrica*, Vol. 67, pp. 349–374.

- BOYLAN, R. (2012) “The Effect of Punishment Severity on Plea Bargaining,” *Journal of Law and Economics*, Vol. 55, pp. 565–591.
- BRAY, S. (2005) “Not Proven: Introducing a Third Verdict,” *University of Chicago Law Review*, Vol. 72, pp. 1299–1329.
- BURNS, R. (2009) *The Death of the American Trial*, University of Chicago Press.
- COONS, J.E. (1963) “Approaches to Court Imposed Compromise—The Uses of Doubt and Reason,” *Northwestern University Law Review*, Vol. 58, pp.750–794.
- DA SILVEIRA, B. (2015) Bargaining with Asymmetric Information: An Empirical Study of Plea Negotiations,” *Working Paper*, Washington University.
- DAUGHETY, A., REINGANUM, J. (2015a) “Informal Sanctions on Prosecutors and Defendants and the Disposition of Criminal Cases,” forthcoming in the *Journal of Law, Economics, and Organization*.
- DAUGHETY, A., REINGANUM, J. (2015b) “Selecting Among Acquitted Defendants: Procedural Choice vs. Selective Compensation,” *Working Paper*, Vanderbilt University.
- DAVIS, J. (2002) “Taking Uncertainty Seriously: Revising Injunction Doctrine,” *Rutgers Law Journal*, Vol. 34, pp.363–462.
- DRESSLER, J., THOMAS, G. (2010) “Does the Process Protect the Innocent,” in *Criminal Procedure: Prosecuting Crime*, Fourth Edition, West Academic Publishing.
- ELDER, H. (1989) “Trials and Settlement in the Criminal Courts: an Empirical Analysis of Dispositions and Sentencing,” *Journal of Legal Studies*, Vol. 18, pp. 191–208.
- FEDDERSEN, T., PESENDORFER, W. (1996) “The Swing Voter’s Curse,” *American Economic Review*, Vol. 86, pp. 408–424.
- FEDDERSEN, T., PESENDORFER, W. (1997) “Voting Behavior and Information Aggregation in Elections with Private Information,” *Econometrica*, Vol. 65, pp. 1029–1058.
- FISCHER, D. (1993) “Proportional Liability: Statistical Evidence and the Probability Paradox,” *Vanderbilt Law Review*, Vol. 46, pp.1201–1226.
- FISCHER, D. (2001) “Tort Recovery for Loss of a Chance,” *Wake Forest Law Review*, Vol. 36, pp. 605–655.
- FISHER, T. (2011) “Constitutionalism and the Criminal Law: Rethinking Criminal Trial Bifurcation,” *University of Toronto Law Journal*, Vol. 61, pp. 811–843.
- FISHER, T. (2012) “Conviction Without Conviction,” *Minnesota Law Review*, Vol. 96, pp. 833–885.
- GERARDI, D., YARIV, L. (2007) “Deliberative Voting,” *Journal of Economic Theory*, Vol. 134, pp. 317–338.
- GROGGER, J. (1992) “Arrests, Persistent Youth Joblessness, and Black-White Employment Differentials,” *Review of Economics and Statistics*, Vol. 74, pp. 100–106.
- GROGGER, J. (1995) “The Effect of Arrest on the Employment and Earnings of Young Men,” *Quarterly Journal of Economics*, Vol. 90, pp. 51–72.

- GROSS, S., O'BRIEN, B., HU, C., AND E. KENNEDY (2014) "Rate of False Conviction of Criminal Defendants who are Sentenced to Death," *Proceedings of the National Academy of Sciences*, Vol. 111, pp. 7230–7235.
- GROSSMAN, G., AND KATZ, M. (1983) "Plea Bargaining and Social Welfare," *American Economic Review*, Vol. 73, pp. 749–757.
- KAMENICA, E., AND GENTZKOW, M. (2011) "Bayesian Persuasion," *American Economic Review*, Vol. 101, pp. 2590–2615.
- KAPLOW, L. (2011) "On the Optimal Burden of Proof," *Journal of Political Economy*, Vol. 119, pp. 1104–1140.
- KAYE, D. (1982) "The Limits of the Preponderance of the Evidence Standard: Justifiably Naked Statistical Evidence and Multiple Causation," *American Bar Foundation Research Journal*, Vol. 7, pp. 487–516.
- LANDO, H. (2005) "The Size of the Sanction should Depend on the Weight of the Evidence," *Review of Law and Economics*, Vol. 1, pp. 277–292.
- LEE, S. (2014) "Plea Bargaining: On the Selection of Jury Trials," *Economic Theory*, Vol. 57, pp. 59–88.
- LOTT, J. (1990) "The Effect of Conviction on the Legitimate Income of Criminals," *Economics Letters*, Vol. 34, pp. 381–385.
- MASCOLO, E. (1985) "Procedural Due Process and the Lesser-Included Offense Doctrine," *Albany Law Review*, Vol. 50, pp. 263–304.
- MILGROM, P., SEGAL, I. (2002) "Envelope Theorems for Arbitrary Choice Sets," *Econometrica*, Vol. 70, pp. 583–601.
- MILGROM, P., SHANNON, C. (1994) "Monotone Comparative Statics," *Econometrica*, Vol. 62, pp. 157–180.
- NELSON, W. (2013) "Political Decision Making by Informed Juries." *William and Mary Law Review*, Vol. 55, pp. 1149–1166.
- QUAH, J., AND STRULOVICI, B. (2012) "Discounting, Values, and Decisions," *Journal of Political Economy*, Vol. 121, pp. 898–939.
- REIMAN, J. AND VAN DEN HAAG, E. (1990) "On the common saying that it is better that ten guilty persons escape than that one innocent suffer: pro and con," *Social Philosophy and Policy*, Vol. 7, pp. 226–248.
- ROSTRON, A. (2004) "Beyond Market Share Liability: A Theory of Proportional Share Liability for Nonfungible Products," *UCLA Law Review*, Vol. 52, pp.151–215.
- SAND, L., ROSE, D. (2003) "Proof Beyond All Possible Doubt: Is there a Need for Higher Burden of Proof When the Sentence May Be Death," *Chicago-Kent Law Review*, Vol. 78, pp. 1359–1376.
- SAUER, K. (1995) "Informed Conviction: Instructing the Jury About Mandatory Sentencing Consequences," *Columbia Law Review*, Vol. 95, pp. 1232–1272.
- SHAVELL, S. (1985) "Uncertainty over Causation and the Determination of Civil Liability," *Journal of Law and Economics*, Vol. 28, pp. 587–609.
- SIEGEL, R., AND STRULOVICI, B. (2019) "Judicial Mechanism Design," *Working Paper*.

SILVA, F. (2018) "If We Confess Our Sins," *International Economic Review*, forthcoming.

SPOTTSWOOD, M. (2019) "Towards a Continuous Burden of Proof," *Working Paper*.

TEICHMAN, D. (2017) "Convicting with Reasonable Doubt: An Evidentiary Theory of Criminal Law," *Notre Dame Law Review*, Vol. 93, pp. 757–810.