
SYMPOSIUM:
MODELING HUMAN DECISIONMAKING
IN THE LAW

THE NEUROBIOLOGY OF OPINIONS:
CAN JUDGES AND JURIES BE
IMPARTIAL?

ISABELLE BROCAS^{*}

JUAN D. CARRILLO[†]

ABSTRACT

In this Article we build on neuroscience evidence to model belief formation and study decisionmaking by judges and juries. We show that physiological constraints generate posterior beliefs with properties that are qualitatively different from traditional Bayesian theory. In particular, decisionmakers will tend to reinforce their prior beliefs and to hold posteriors influenced by their preferences. We study the implications of the theory for decisions rendered by judges and juries. We show that early cases in judges' careers may affect their decisions later on, and that early evidence produced in a trial may matter more than late evidence. In the case of juries, we show that the well-known polarization effect is a direct

^{*} Associate Professor of Economics, University of Southern California; Research Fellow, Center for Economic Policy Research. We are grateful to participants in the Center for the Study of Law and Politics symposium on "Modeling Human Decisionmaking in the Law" for useful comments.

[†] Professor of Economics, University of Southern California; Research Fellow, Center for Economic Policy Research.

consequence of physiological constraints. It is more likely to be observed when information is mixed, as behavioral evidence suggests, and when prior beliefs and preferences are initially more divergent across jurors.

I. INTRODUCTION

Biases in beliefs and behavior have long been noted in many fields of psychology and economics. Although it is possible to categorize some specific biases and associate them with specific situations, the fundamental causes for biases are not well understood. Moreover, even though the standard Bayesian framework does not capture the way beliefs are revised after exposure to evidence, it is difficult to pinpoint an alternative model that does strictly better in a large range of situations.

In a 2012 study, we postulated that biases in beliefs and in decisionmaking may be mediated by physiological constraints that place a bound on the amount of information processed by neurons.¹ The idea is intuitive. The information about the world is encoded in the sensory system and processed and decoded by a variety of systems before being mapped into a decision. If the encoding-decoding procedure is constrained, then the information transmitted to the decisionmaker will differ from the information actually used to make a decision. In turn, the action of the decisionmaker will appear biased to an outside observer. Interestingly, the neurobiology literature shows that such constraints exist and act in a very specific way. In our prior work, *From Perception to Action: An Economic Model of Brain Processes* (“2012a study”), we linked the neurobiological evidence about information processing to systematic biases in decisionmaking and beliefs.² In *The Neuroeconomics of Strategic Decision-Making* (“2012b study”), we extended the analysis to situations in which several decisionmakers interact after privately making sense of the information they are exposed to.³ The objective of this Article is to build on those studies to analyze the behavior of judges and juries. We rely on results that have been demonstrated in those prior studies and illustrate how they apply to this particular case.

The task of judges and juries is to interpret the evidence produced and make a decision based on it. In a purely Bayesian world, judges may make different decisions due to different opinions or intensity of preferences.

1. Isabelle Brocas & Juan D. Carrillo, *From Perception to Action: An Economic Model of Brain Processes*, 78 GAMES & ECON. BEHAV. 81, 83 (2012).

2. *Id.* at 91–93.

3. Isabelle Brocas & Juan D. Carrillo, *The Neuroeconomics of Strategic Decision-Making* (2012) (unpublished manuscript) (on file with author).

However, posterior beliefs (“posteriors”) should be entirely driven by prior beliefs (“priors”) and information. In particular, two judges with the same prior belief and who are exposed to the same evidence should hold the same posterior belief. Similarly, in the Bayesian framework, we should see that members of a jury facing the same evidence will all update their beliefs in the same direction or have converging views; we should never, however, observe polarization. We show that these properties are not necessarily true when physiological constraints are taken into account.

More specifically, we show that physiological constraints generate biases in beliefs: decisionmakers will tend to form posteriors that confirm their priors and that are affected by the magnitude of their payoffs or preferences.⁴ Therefore, two individuals with the same prior who observe the same evidence will end up holding different posteriors if their preferences differ. Also, two individuals with different priors that have the same preferences may update their beliefs in opposite directions.⁵ These results offer a framework to address various implications on decisions rendered by judges and juries.

We show that physiological constraints make the order in which evidence is received critical.⁶ Therefore, cases analyzed in the early career of a judge may affect the decisions that this judge will take on later in a priori independent cases. Also, early evidence produced in a trial may matter more than late evidence. Hence, it is not the same to be exposed first to strong evidence a crime has been committed as it is to listen first to the childhood story of the criminal.

The case of juries is also interesting. In particular, we show that the distribution of preferences in a jury affects the way information is interpreted by individual jurors. If jurors are all willing to make the correct decision but have different priors or different preferences (and are therefore inclined to take different decisions), each will interpret the available evidence differently. This may result in polarization, defined as the fact that two subjects with either different preferences or different priors may move their beliefs further apart after being exposed to identical mixed evidence.⁷ Such an outcome is in line with a long series of studies,⁸ and cannot be reconciled with traditional Bayesian theories of decisionmaking.

4. See *infra* Part III.A.

5. See *infra* Part III.B.

6. See *infra* Part III.C.

7. See *infra* Part IV.A.

8. See, e.g., Robert M. Bray & Audrey M. Noble, *Authoritarianism and Decisions of Mock Juries: Evidence of Jury Bias and Group Polarization*, 36 J. PERSONALITY & SOC. PSYCHOL. 1424–30

Our study can also address a series of questions concerning the design of rules. We study the effect of the rule specified in cases in which a unanimous verdict is not reached by a jury, and we show that it also affects the way information is interpreted by individual jurors.⁹ Even though our model focuses on stylized cases, it suggests that the rule affects strongly the probability of polarization.

To illustrate our theory in the simplest possible terms, we consider a problem with two underlying states of the world: whether a crime has been committed or not. Evidence is produced, aggregated, and transmitted. It is positively but imperfectly correlated with the true state. Individuals who observe this evidence interpret it to make a decision (a verdict in the case of a judge or a recommendation in the case of a juror). Decisions are implemented yielding state-dependent payoffs for all the subjects involved.

The Article is organized as follows: In Part II, we present the model consistent with the neurobiology evidence. In Part III, we study the case in which a decision is delegated to a judge. In Part IV, we investigate recommendations by a jury. In Part V, we offer some concluding remarks.¹⁰

II. THE NEUROBIOLOGY OF DECISIONMAKING

To best fit the application to legal environments, we consider the following stylized situation. A person recently arrested for a crime is either guilty (state *A*) or innocent (state *B*). This person has to be evaluated by individual *i* (for example, a judge) who holds a prior belief p_i that the state is *A*.

When evidence is produced, *i* interprets the information to update the individual's belief, and takes the decision to convict (action *a*) or acquit (action *b*) the person. In this simple model, *a* is the correct action in state *A* (convict a guilty person) and *b* in state *B* (acquit an innocent person). We assume that *i* always wants to take the correct action. This is captured with the following payoffs:

(1978); John M. Darley & Paget H. Gross, *A Hypothesis-Confirming Bias in Labeling Effects*, 44 J. PERSONALITY & SOC. PSYCHOL. 20, 28–30 (1983); Charles G. Lord, Lee Ross & Mark R. Lepper, *Biased Assimilation and Attitude Polarization: The Effects of Prior Theories on Subsequently Considered Evidence*, 37 J. PERSONALITY & SOC. PSYCHOL. 2098, 2100–01 (1979); Eleanor C. Main & Thomas G. Walker, *Choice Shifts and Extreme Behavior: Judicial Review in the Federal Courts*, 91 J. SOC. PSYCHOL. 215, 220–21 (1973); Scott Plous, *Biases in the Assimilation of Technological Breakdowns: Do Accidents Make Us Safer?*, 21 J. APPLIED SOC. PSYCHOL. 1058, 1078 (1991).

9. See *infra* Part IV.B.

10. We will keep mathematical notations to a minimum in the main text. For details of the model and proofs of the results, see Brocas & Carrillo, *supra* note 1; Brocas & Carrillo, *supra* note 3; and *infra* Appendix.

$$U_i(a;A)=G_A^i>0, \quad U_i(b;B)=G_B^i>0, \quad U_i(b;A)=0, \quad U_i(a;B)=0$$

In words, wrong decisions are normalized to 0 and correct decisions yield positive payoffs. The parameters G_A^i and G_B^i capture the incremental utility of choosing correctly ($U_i(a;A)-U_i(b;A)$ and $U_i(b;B)-U_i(a;B)$). In this model, there is no scope for partisan recommendations: all individuals want to make the objectively correct decision and would do so under full information (a if $\Pr(A)=1$ and b if $\Pr(B)=1$).

We borrow the framework developed in our 2012a study to formalize the actual decisionmaking process.¹¹ Following the evidence from neurobiology, we model information processing and decisionmaking as a threshold mechanism in which the evidence produced is first encoded in the sensory system and then interpreted in reference to a threshold.¹² We now introduce it more formally.

Let us first take a step back and abstract from our application to describe the conceptual framework developed in neurobiology. It is based on experimental research in which subjects are asked to extract relevant information from noisy evidence before making a decision. A typical paradigm is subjects having to identify a color or the direction of a movement and report their perceptions. They are rewarded when their answer is correct. Both behavior and neural activity are recorded and correlated.

The evidence produced is a signal about the underlying state, which is first encoded in the sensory system. For instance, the decisionmaker must detect the color of an object placed in front of him and is given only a glimpse at the object under a certain light condition (the signal). Assume there may be two possible colors, black or white. Neurons detecting each color will react according to the strength of the signal. In particular, the light intensity and conditions will affect cell firing.¹³ We represent the encoded evidence by $c \in [0,1]$, which corresponds to the neuronal activity in the sensory system. The variable c can be interpreted as the ratio of neurons that detect the black color. The signal is imprecise but informative:

11. See generally Brocas & Carrillo, *supra* note 1 (establishing a formal economic framework capable of predicting decisions based on neurobiological premises of information processing).

12. *Id.* at 87.

13. Other factors affect cell variability. For instance, even when exposed to the same stimuli, neurons do not always fire in the same way. This can be understood as internal noise and may vary across individuals and across experiences. We will neglect this type of variability to concentrate on the effect of external variability on decisionmaking. For a detailed discussion, see Isabelle Brocas, *Information Processing and Decision-Making: Evidence from the Brain Sciences and Implications for Economics*, 83 J. ECON. BEHAV. & ORG. 292, 294–98 (2012).

a high-cell firing in favor of black is (stochastically) more likely to occur when black is the true color, while a low-cell firing in favor of black is (stochastically) more likely to occur when white is the true color. The same mechanism applies for the case we are interested in, with c representing the fraction of neurons supporting the hypothesis a crime has been committed given the evidence produced.

The next step is to understand how a decision is made based on the encoded information c . In a classical study, Doug Hanes and Jeffrey Schall use single-cell recording to analyze the neural processes responsible for the duration and variability of reaction times in monkeys.¹⁴ The authors find that movements are initiated when neural activity reaches a certain threshold activation level, in a winner-takes-all type of contest.¹⁵ This evidence suggests that the process can be schematically represented by a decision-threshold mechanism: it is as if there exists a threshold x such that action b is triggered when $c < x$ (for example, release the person when there is enough evidence the person is not guilty), and action a is triggered when $c \geq x$ (convict the person when there is enough evidence the person is guilty). At the same time, it filters information out. In other words, the mechanism provides an interpretation of the information. The sensory system collects c and the decision system interprets it as either $c < x$ or $c \geq x$, where the former is evidence of B and the latter is evidence of A . The decision system compares alternatives via this mechanism.¹⁶ This type of threshold mechanism is widely used in neurobiology to account for behavior and neural activity.

It is important to realize that the threshold x represents actual neuronal thresholds and synaptic connections. Those are the physical elements through which the threshold mechanism is implemented: they filter information just as the threshold does. A neuron fires if it receives enough input activity from neurons in previous layers, which requires a strong synaptic connectivity between neurons. Depending on the level of neuronal thresholds and the strength of the synaptic connections, neuronal activity will be stopped or propagated along a given path and will trigger one action

14. Doug P. Hanes & Jeffrey D. Schall, *Neural Control of Voluntary Movement Initiation*, 274 SCI. 427, 427 (1996).

15. *Id.* at 427–29.

16. See Jochen Ditterich, Mark E. Mazurek & Michael N. Shadlen, *Microstimulation of Visual Cortex Affects the Speed of Perceptual Decisions*, 6 NATURE NEUROSCIENCE 891, 894–96 (2003); Joshua I. Gold & Michael N. Shadlen, *Neural Computations that Underlie Decisions About Sensory Stimuli*, 5 TRENDS COGNITIVE SCI. 10, 13–15 (2001); Michael N. Shadlen et al., *A Computational Analysis of the Relationship Between Neuronal and Behavioral Responses to Visual Motion*, 16 J. NEUROSCIENCE 1486, 1491–1500 (1996).

or the other. This mechanism is economical in that it requires scarce knowledge to reach a decision.

We adopt this threshold mechanism and assume that the threshold x can be optimized to maximize expected utility. This presupposes that the information is interpreted taking into account information regarding prior beliefs and the intensity of preferences. This also presupposes that the way this information is taken into account is compatible with expected utility theory. Assuming neurons can perform this type of optimization is consistent with the work by Paul Platt and Michael Glimcher.¹⁷ The authors show that the brain represents the magnitude of the possible payoffs¹⁸ as well as their likelihood (prior beliefs and signals). They also show that neural activity correlates with decisions and fits an expected utility maximization model.¹⁹ In typical oculomotor tasks, rewards and information components modulate the activation of the lateral intraparietal area where the decision is computed given projections from the reward system (which evaluates payoffs) and the sensory system (which interprets the evidence).²⁰

Overall, decisionmaking in the brain can be represented by an “as if” model with the following timing of events: First, a threshold is set. Second, evidence is encoded in the sensory system. Third, evidence is compared to the threshold. And fourth, an action is implemented. An optimal threshold is one that takes into account all the relevant information and realizes that it is filtered out when making a decision. Said differently, it maximizes the expected utility of the individual given the likelihood of the events, that is, given the evidence about the state that is retained.

Recall that in the threshold mechanism, a decision is based on a coarse partition of information: either there is reasonable evidence in favor of A ($c \geq x$) or reasonable evidence in favor of B ($c < x$). Therefore, depending on these two possibilities, i can hold ex post either of two posterior beliefs: $\Pr(A | c \geq x) \equiv p_i(x)$ or $\Pr(A | c < x) \equiv p_i(x)$, depending on whether $c \geq x$ or $c < x$.

These beliefs are to be contrasted with those obtained in a pure Bayesian framework, which would be based on the actual evidence collected c . Formally, in the standard Bayesian model, the posterior belief

17. See Michael L. Platt & Paul W. Glimcher, *Neural Correlates of Decision Variables in Parietal Cortex*, 400 NATURE 233, 236–38 (1999).

18. G_A^i and G_B^i in our model.

19. *Id.* at 237–38.

20. *Id.* at 234.

based on the evidence c would be $\Pr(A|c) \equiv \pi_i(c)$. Our first result is as follows:

Theorem 1: *In the optimal threshold mechanism, decisionmaking is efficient but posterior beliefs are biased.*²¹

It is efficient to take action a if the probability that the true state is A is high enough. In the standard Bayesian model, it means that there exists a certain amount of evidence x_i^* such that it is efficient to take action a for all $c \geq x_i^*$ and to take action b for all $c < x_i^*$.

This means in particular that, for the purpose of decisionmaking, it is irrelevant whether the information obtained is slightly below x_i^* or far below it: the same decision is efficient either way (action b). Therefore, basing a decision on the coarse information partition $c < x_i^*$ or $c \geq x_i^*$ does not affect the ultimate decision.

However, the posterior emerging from the threshold mechanism is biased because the intensity of the signal is emphasized or deemphasized to keep just the information necessary to take the optimal action. Therefore, the individual will report an opinion that does not take into account the actual information the individual observed but reflects instead his or her (biased) interpretation.

We will exploit this result to analyze belief formation and decisionmaking in legal environments in which a decision must be made regarding a potential offense. We have in mind situations in which a judge must interpret the evidence available in order to make a decision. We will analyze such situations in Part III. We are also interested in situations in which a jury is delegated the decision and the verdict is subject to institutional rules (for instance, unanimous or majority voting). In such situations, jury members will interpret the evidence provided and recommend the decision they think is best. Those situations will be studied in Part IV.

III. DECISIONS RENDERED BY A JUDGE

The motivating example in this section is a judge who must review evidence and make a decision based on the judge's own judgment. The judge wants to take the correct action given the judge's limited information and is endowed with the preferences described earlier.

21. Brocas & Carrillo, *supra* note 1, at 86–89.

A. PRIORS, PREFERENCES, AND BIASES IN BELIEFS

We first address the question of whether posterior beliefs are biased in a systematic way. In particular, we are interested in determining if particular prior beliefs or particular attitudes towards outcomes may affect the interpretation of the evidence.

Proposition 1: *The interpretation of the evidence affects the judge's information processing in a way that (1) the judge's prior beliefs tend to be confirmed, and (2) the judge's posterior beliefs depend on the payoffs of the different actions.*²²

This result is a special case of our 2012a study.²³ Even though the judge is impartial, in the sense that the judge does not favor a decision and rather tries to take the one that best fits the state, the judge is subject to two biases. First, the judge tends to reinforce his or her prior: the higher the judge's confidence in state B , the higher the threshold x_i^* . This means that the judge is more prone to interpret information as evidence of B (technically, $c < x_i^*$ is more likely) and update his or her prior toward that state.

The intuition for the confirmatory bias effect is as follows: suppose the judge believes that the person is very likely to be innocent ($\Pr(B)$ is high) and is leaning towards acquitting the person given the current evidence. The judge will require a lot of evidence of the person being guilty in order to change the judge's mind. This is rational as the judge will need to be confident in the judge's change of mind to compensate for his or her current belief. This means that the threshold should be high. It is therefore very likely the evidence will fall below the threshold. As a consequence, the judge is likely to interpret the evidence in favor of his or her initial prior and acquit the person eventually.

Second, the judge tends to favor the state in which the highest payoffs can potentially be obtained: the higher the benefit of taking the correct action in state B , the higher the threshold. Again, the judge is likely to interpret the information as evidence of B and update his or her prior towards B . The logic of the argument is the same as before: suppose the judge is most interested in releasing innocent persons (and only to a lesser degree, convicting guilty persons). The judge will change his or her mind only if very strong evidence that the person is guilty is produced.

22. *Id.* at 93–95.

23. *Id.*

Therefore, it is optimal to set a high threshold which, again, makes a change of mind very unlikely.

As mentioned, decisions are not biased if the threshold is set optimally. That is, the judge will take the same action as in the standard Bayesian framework in which the signal is perfectly processed rather than interpreted. The main consequence, however, is that the judge will hold degrees of confidence in his or her decision that do not reflect the exact evidence produced. This would have an effect on future decisions if the judge were to reevaluate past choices or obtain new information.

B. OPINIONS AND IMPARTIAL JUDGMENTS

The previous section established that biases in beliefs can be linked to primitives: posterior beliefs are shaped not only by prior beliefs, but also by the magnitude of the state-contingent payoffs. Even though a judge tries to be impartial, the judge may end up holding a belief that reflects his or her own aversion to a particular crime and report an opinion that does not seem in line with the evidence from the perspective of an outside observer. The next result makes this statement more transparent:

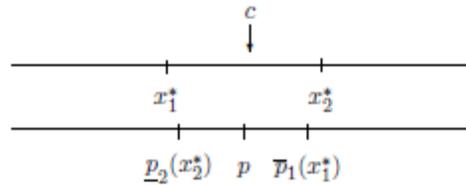
Proposition 2: *Two judges exposed to the same evidence may report different opinions. In particular, two judges sharing the same prior may end up having different posteriors, and two judges with different priors may interpret identical evidence in opposite ways.*

This proposition is a consequence of the results obtained in Part III.A.²⁴ Two judges holding the same prior belief and exposed to the same evidence would hold the same posterior belief in the standard Bayesian model. However, in the threshold model, they may end up with different posterior beliefs. This result can be due exclusively to a difference in preferences, and not to a difference in prior beliefs—a leading assumption to explain this type of phenomenon. Consider two judges, 1 and 2, who hold the same prior belief but different preferences, namely $G_1 > G_2$. Under Proposition 1, $x_1^* < x_2^*$. Assume also that the evidence is mixed in that it falls between the two thresholds, $c \in (x_1^*, x_2^*)$. In the standard model, the posterior beliefs are identical: $p_1(c) = p_2(c)$. In the threshold model, however, judge 1 revises his or her belief upwards (the threshold is surpassed) while judge 2 revises it downwards (the threshold is not reached). In other words, they disagree even though they share the same prior and are exposed to the same evidence.²⁵ Figure 1 depicts this.

24. See *supra* Part III.A.

25. Notice that they also take different actions (judge 1 chooses action a whereas judge 2

FIGURE 1. Judges Having the Same Prior p but Different Preferences Exposed to the Same Mixed Evidence c

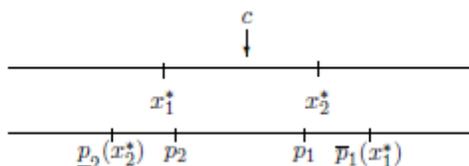


Also, two judges having the same preferences but whose prior opinions differ, $p_1 \neq p_2$, may interpret identical evidence in opposite ways. Consider for example two judges, 1 and 2, with priors $p_1 > p_2$. Again following Proposition 1, $x_1^* < x_2^*$. Assume now that evidence is mixed, $c \in (x_1^*, x_2^*)$. In the threshold model, judge 1 will update his or her belief upwards as his or her threshold is surpassed, while judge 2 will update his or her belief downwards as his or her threshold is not reached (see figure 2), thereby leading to a polarization of opinions. The proposition that people holding different priors may disagree after receiving the same evidence is, again, not novel. However, in the Bayesian model, subjects exposed to the same evidence would either update in the same direction (both upward or both downward) or else converge (the subject with the initially lowest belief increases it and the subject with the initially highest belief decreases it). Either way, they could never exhibit divergent beliefs.

The main lesson so far is that identical evidence will be interpreted differently by individuals with different preference intensities or different priors. Also, posterior beliefs are shaped by preferences. Individuals tend to confirm their priors and to develop strong confidence in their decisions to avoid the least beneficial outcomes.

chooses action b). However, this is not surprising: it is also true in the standard Bayesian setting that two individuals with different preferences facing the same information may choose different actions.

FIGURE 2. Judges Having Different Priors p_1 and p_2 but the Same Preferences Exposed to the Same Mixed Evidence c .



Interestingly, these results are reminiscent of observations made in a series of studies. Biases in beliefs and decisions have been reported in cases and experiments, suggesting in particular that the formation of prior beliefs is a key element in judicial decisions. For instance, Chris Guthrie, Jeffrey Rachlinski, and Andrew Wistrich report that federal judges may be subject to framing effects and anchoring.²⁶ Along the same lines, a large literature on the hindsight bias shows that decisionmakers assess the ex ante likelihood of a crime differently if they are revealed the ex post outcome.²⁷ This large body of work suggests that initial opinions—or priors—may be manipulated, and these manipulations will eventually play an important role in final decisionmaking. Even though our model does not address the formation of prior beliefs, it emphasizes the mechanism through which initial frames or information sets bias information acquisition and interpretation throughout trials, as well as at eventual decisions.²⁸ From the

26. Chris Guthrie, Jeffrey J. Rachlinski & Andrew J. Wistrich, *Inside the Judicial Mind*, 86 CORNELL L. REV. 777, 789–99 (2001). See also Chris Guthrie, *Prospect Theory, Risk Preference and the Law*, 97 NW. U. L. REV. 1115, 1123–27 (2003) (discussing the same for litigants involved in settlements).

27. See, e.g., Reid Hastie & W. Kip Viscusi, *What Juries Can't Do Well: The Jury's Performance as a Risk Manager*, 40 ARIZ. L. REV. 901, 914–17 (1988); Christine Jolls, Cass R. Sunstein & Richard Thaler, *A Behavioral Approach to Law and Economics*, 50 STAN. L. REV. 1471, 1523–27 (1988). For early studies of hindsight bias, see generally Hal R. Arkes et al., *Eliminating the Hindsight Bias*, 73 J. APPLIED PSYCHOL. 305 (1988), and Hal R. Arkes & Cindy Schipani, *Medical Malpractice v. the Business Judgment Rule: Differences in Hindsight Bias*, 73 OR. L. REV. 587 (1994). For a perspective on medical malpractice and hindsight bias, see generally Baruch Fischhoff, *Hindsight ≠ Foresight: The Effect of Outcome Knowledge on Judgment Under Uncertainty*, 1. J. EXPERIMENTAL PSYCHOL.: HUM. PERCEPTION & PERFORMANCE 288 (1975). For a review of legal literature and hindsight bias, see Russell B. Korobkin & Thomas S. Ulen, *Law and Behavioral Science: Removing the Rationality Assumption from Law and Economics*, 88 CALIF. L. REV. 1051, 1095–1100 (2000).

28. Several other biases have been documented in the legal literature. Among others, decisionmakers tend to blame action more than inaction and to overemphasize outcomes that spring from abnormal rather than normal circumstances, Robert A. Prentice & Jonathan J. Koehler, *A Normality Bias in Legal Decision Making*, 88 CORNELL L. REV. 583, 590–96 (2003), and are subject to category-bound thinking, Cass R. Sunstein et al., *Predictably Incoherent Judgments*, 54 STAN. L. REV. 1153, 1170–73 (2002). Judges may also be influenced by demanded sentences. Birte English & Thomas Mussweiler, *Sentencing Under Uncertainty: Anchoring Effects in the Courtroom*, 31 J. APPLIED SOC.

perspective of legal design, special attention should be put on the formation of initial beliefs.²⁹

C. DECISIONMAKING AND BELIEFS OVER TIME

Recall that decisionmaking is efficient conditional on prior beliefs and preferences. Said differently, if at time t the judge holds a belief μ_t , the threshold is set in such a way that the optimal decision is taken after the release of date t evidence. However, in many applications, a judge is receiving several pieces of evidence before making a decision. Given posteriors are built via a biasing filtering mechanism, the belief used when the last piece of evidence is interpreted reflects the previous filtering. Because such belief is biased at that point, the decision itself will also be biased from an earlier perspective.

Proposition 3: *Beliefs and decisions are path dependent. In particular, the order in which the evidence is received affects choices over time.*

To understand this proposition, consider a judge who initially believes the person is equally likely to be guilty or innocent and therefore sets an average threshold. The judge receives two pieces of evidence: one (weakly) in favor of A and the other (weakly) in favor of B .

If the judge is first exposed to the evidence in favor of B , the average threshold is not reached. The judge's confidence in B increases, and so does his or her threshold. The second piece of evidence (weakly in favor of A) does not reach this higher threshold either, so it does not reverse the judge's belief that b should be implemented. Conversely, if the judge is first exposed to the evidence in favor of A , the average threshold is surpassed. The threshold is lowered so the second piece of evidence (weakly in favor of B) also surpasses it. The judge ends up having two pieces of evidence above the thresholds and chooses action a .

Overall, the result shows that, in sharp contrast to a standard Bayesian framework, first impressions matter in the threshold model. Indeed, the thresholds endogenously generate an anchoring effect: the judge may "adopt" the interpretation of the first piece of evidence and reinforce this interpretation as new evidence comes in. Under some specifications, we

PSYCHOL. 1535, 1540–49 (2001).

29. This echoes many studies of the hindsight bias, Jolls, Sunstein & Thaler, *supra* note 27, at 1523–27; Kim A. Kamin & Jeffrey J. Rachlinski, *Ex Post ≠ Ex Ante: Determining Liability in Hindsight*, 19 L. & HUM. BEHAV. 89, 91, 99–101 (1995), which has been reported to be the cause of some verdicts. *See, e.g.*, *KSR Int'l Co. v. Teleflex, Inc.*, 550 U.S. 398, 421–23 (2007).

can show that thresholds also become more extreme over time.³⁰ Proposition 3 has many implications. We discuss two that are particularly relevant for our application.

1. Sequence of Information Revelation in Trials

A judge is never exposed to one single piece of evidence before making a recommendation or to a number of arguments presented simultaneously. The decision comes after a sequence of signals has been disclosed and interpreted. The previous result suggests that the order in which the evidence is produced may tilt the verdict. In particular, early disclosure of evidence that the person is innocent is more likely to be followed by a release. These anchoring effects are reinforced by the intensity of the preferences. A judge who values releasing innocent people relatively more than convicting guilty people and who is exposed first to favorable evidence will most likely release the person. Of course, this is true only if the decision is taken under a certain amount of doubt (a common scenario). If unambiguous evidence is produced, any judge in our model would take the “correct” decision.

Our results are also relevant to the presumption of innocence principle, which reduces roughly to working under the assumption that the prior $\Pr(A)$ is low enough so that in the absence of evidence, the person should be acquitted. Given our results, the principle should affect the way a judge forms an opinion based on evidence, because thresholds are such that priors tend to be reinforced. A verdict is more likely to be favorable if it is rendered under the presumption of innocence and where the very first evidence produced suggests innocence.

2. Wisdom or Stubbornness

Judges are often appointed for long periods of time. The evidence produced in a given trial sometimes contains information pertinent for other trials. The judge may therefore learn about the underlying state of the world throughout the judge’s mandate. The judge’s prior belief in the first case of his or her career will therefore be different from the judge’s prior belief in a later case. The biasing effects emphasized here suggest that judges may not only develop opinions that reinforce earlier pieces of evidence, but also take path-dependent biased decisions. A judge who convicts the first few persons is more likely to convict offenders in the

30. See Brocas & Carrillo, *supra* note 1, at 94–95.

future. This may have implications on the optimal length of judges' mandates from the perspective of a planner.

IV. DECISIONS RENDERED BY JURIES

We now turn to study the verdicts of juries. This problem has been analyzed theoretically elsewhere. Notably, Timothy Feddersen and Wolfgang Pesendorfer proposed a model of strategic voting by jurors to assess the merits of unanimous verdicts.³¹ Contrary to earlier literature, in which it is assumed that each juror behaves as if his or her vote alone determines the outcome, the authors assume that each juror possesses private information about the state (for example, technical knowledge) and makes a strategic vote.³² Our perspective is different. Instead of assuming jurors differ in their private information, we assume they differ in either their taste (preferences about outcomes) or initial opinion (prior beliefs). Also, rather than modeling the information produced during the trial as a private signal that is processed in a pure Bayesian manner, we model it as a common signal which is processed through a threshold mechanism. Our perspective is related, however, because we do presuppose a strategic interpretation of information. In other words, each juror takes into account that other jurors interpret the information to make a recommendation and that recommendations are aggregated given the rule.

The analysis in this section borrows results from our 2012b study.³³ With respect to the previous section, we extend the model in the following way: the jury is composed of n jurors indexed by i . When evidence is produced, each juror interprets the information to form an updated belief. Each juror makes a recommendation r_i , and the judge makes a decision based on the reports of all jurors. Jurors are heterogeneous. To keep the analysis simple and tractable, we assume that there are two possible types of jurors indexed by k . The proportion of type-1 and type-2 jurors is v and

31. Timothy J. Feddersen & Wolfgang Pesendorfer, *Convicting the Innocent: The Inferiority of Unanimous Jury Verdicts Under Strategic Voting*, 92 AM. POL. SCI. REV. 23, 23–30 (1998).

32. *Id.* at 24. Many other studies in political economy have analyzed the effect of strategic voting. See, e.g., David Austen-Smith & Jeffrey S. Banks, *Information Aggregation, Rationality, and the Condorcet Jury Theorem*, 90 AM. POL. SCI. REV. 34, 35–43 (1996) (analyzing the effect of strategic voting on juries and whether decisions made by jurors as a “collective” group are superior to decisions made by individual jurors); Timothy J. Feddersen & Wolfgang Pesendorfer, *The Swing Voter's Curse*, 86 AM. ECON. REV. 408, 412–18 (1996) (analyzing the effect of strategic abstention by uninformed voters in two-candidate elections); Roger B. Myerson, *Extended Poisson Games and the Condorcet Jury Theorem*, 25 GAMES & ECON. BEHAV. 111, 113–130 (1998) (analyzing strategic voting by jurors utilizing game theory).

33. Brocas & Carrillo, *supra* note 3.

1–v. Each type- k juror ($k = 1, 2$) has a prior belief p_k that the person is guilty. Preferences across groups are given by:

$$U_k(a; A) = G_A^k > 0, U_k(b; B) = G_B^k > 0, U_k(b; A) = 0, U_k(a; B) = 0$$

All jurors are exposed to the same evidence and the encoded information is c for all of them, so there are no differences in private signals. When making a recommendation, each juror anticipates how the verdict will be rendered given the recommendations made in his or her group and the recommendations made in the other group. Therefore, recommendations are strategic.

In the next subsections, we assume that given the primitives of the model $(G_A^1, G_B^1, p_1, G_A^2, G_B^2, p_2)$, were jurors delegated the decision, the thresholds x_1^* and x_2^* would be such that $x_1^* < x_2^*$.³⁴ That is, preferences are such that type-1 jurors are more willing to convict the person than type-2 jurors are.

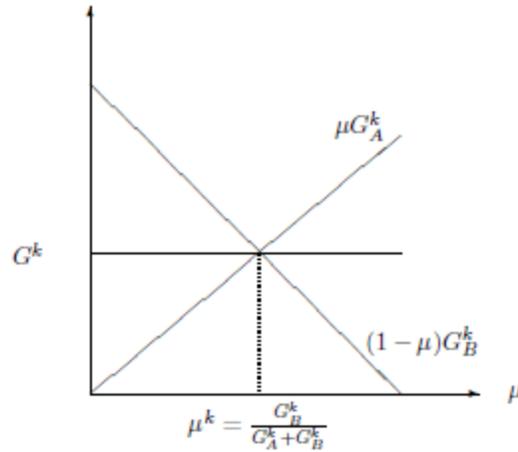
A. POLARIZATION

In this section, we restrict attention to the case where the judge requires a unanimous vote. Whenever jurors do not agree, a default option, o , is implemented. We assume that the default option is a compromise option. To capture this idea, we assume that $U_k(o|A) = U_k(o|B) = G^k$ and, to simplify the analysis, we consider the case illustrated in figure 3, $G^k = G_A^k G_B^k / (G_A^k + G_B^k)$.

In our construction, juror- k possessing a posterior belief below $G_B^k / (G_A^k + G_B^k)$ prefers b to o , and o to a . Similarly, juror- k possessing a posterior belief above $G_B^k / (G_A^k + G_B^k)$ prefers a to o and o to b .

34. For example, $G_A^1 > G_B^1, G_A^2 < G_B^2$ and $p_1 = p_2$.

FIGURE 3. Expected Payoffs as a Function of the Posterior Belief



Consider now the decision of a type- k juror facing evidence c . Evidence c below the threshold x_k triggers recommendation b , while evidence c above the threshold triggers recommendation a . Denote by x_2 the threshold used by type-2 jurors, and consider thresholds x_1 such that $x_1 < x_2$. If $c \leq x_1$, all jurors agree on recommendation b , and if $c \geq x_2$, all jurors agree on recommendation a . When $c \in (x_1, x_2)$, type-1 jurors recommend action a and type-2 jurors recommend action b , in which case o is implemented. From the perspective of a type-1 juror, the payoff obtained if $c \leq x_1$ is the expected payoff of taking action a given the posterior belief $\bar{p}_1(x_1)$. The payoff obtained if $c > x_1$, however, depends on whether the evidence was sufficient to trigger the same recommendation from type-2 jurors, that is, on whether $c \geq x_2$. We have the following result:

Proposition 4: *Type- k jurors set the optimal threshold at x_k^* . They recommend $r_k = b$ if $c < x_k^*$ and $r_k = a$ if $c \geq x_k^*$.*

The verdict is unanimous in favor of b when $c < x_1^$ and unanimous in favor of a when $c > x_2^*$. The vote is split when $c \in (x_1^*, x_2^*)$ and the default option o is implemented. In that case, jurors polarize.*

If type-1 jurors are a priori more in favor of convicting, their threshold will be lower than that of type-2 jurors. As a consequence, there exists a region of mixed evidence (x_1, x_2) where their best choice would be to take action a but action o will be taken instead. Given that the best action is always either action a or action b , the optimal thresholds need only to

discriminate between those two actions and they coincide with the thresholds obtained in the previous section.

The fact that thresholds differ implies that jurors will polarize in cases of disagreement. This result is consistent with evidence from various experimental studies showing that individuals who exhibit confirmatory biases may interpret the same information in opposite ways.³⁵ This “polarization effect” occurs when mixed evidence is given to subjects whose existing views lie on both sides of the evidence.³⁶ Their beliefs may then move further apart. In an early work, Charles Lord, Lee Ross, and Mark Lepper presented a set of typical arguments for and against the death penalty to a pool of subjects.³⁷ When asked about the merits of the death penalty, people who were initially in favor of (respectively against) capital punishment were more in favor of (respectively against) it after reading the studies.³⁸ The literature explains this effect in terms of cognitive biases and non-Bayesian information processing. In particular, it has been argued that individuals focus attention on the elements that support their original beliefs and (consciously or unconsciously) neglect the elements that contradict them.³⁹ Some researchers attribute polarization to heterogeneous prior beliefs⁴⁰ or to non-Bayesian updating.⁴¹ Our analysis suggests that attentional deficits, multiple priors, or biased information processing need not be at the origin of this result. Instead, our threshold mechanism can fully account for this behavior. Indeed, even if $p_1 = p_2$, jurors will polarize as long as they have different preferences (hence, different thresholds) and the evidence received is mixed. Moreover, note that in the threshold mechanism, the information that is retained is still processed in a Bayesian way, but is constrained because some information is filtered out. The result is illustrated in figure 4.

FIGURE 4. Optimal Thresholds in the Jury Model

35. See Darley & Gross, *supra* note 8, at 28–30; Main & Walker, *supra* note 8, at 220–21; Plous, *supra* note 8, at 1078; Cass R. Sunstein, *The Law of Group Polarization*, 10 J. POL. PHIL. 177–80 (2002).

36. See Sunstein, *supra* note 35, at 176.

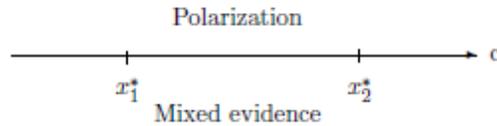
37. Lord, Ross & Lepper, *supra* note 8, at 2100–01.

38. *Id.* at 2103–04.

39. See *id.* at 2105 (describing the results of their study on biased assimilation of empirical evidence, “both proponents and opponents interpreted the additional information . . . as strongly supporting their own initial attitudes”).

40. E.g., Avinash K. Dixit & Jörgen W. Weibull, *Political Polarization*, 104 PROC. NAT’L ACAD. SCI. 7351, 7351–52 (2007).

41. E.g., Matthew Rabin & Joel L. Schrag, *First Impressions Matter: A Model of Confirmatory Bias*, 114 Q.J. ECON. 37, 38 (1999).



Interestingly, the polarization effect depends crucially on the intensity of preferences. Indeed, suppose that the preferences of both types of jurors become stronger in opposite directions. For example, both G_A^1 and G_B^2 increase. Type-1 and type-2 jurors now have higher incentives to recommend actions a and b , respectively. Therefore, type-1 jurors optimally decrease their threshold x_1^* whereas type-2 jurors optimally increase their threshold x_2^* , resulting in an increased gap between x_1^* and x_2^* . However, it is precisely when evidence falls in this region (what we called “mixed evidence”) that polarization occurs. Overall, stronger differences in preferences will result in higher levels of polarization following the release of evidence.

B. DECISION RULES AND BIASES IN BELIEFS

In this section, we study decisionmaking under different rules, concentrating on two examples. In the first example, we consider majority voting. To simplify the analysis, we assume that different types of jurors have identical preferences but different beliefs: $G_A^1 = G_B^1 = G$, $G_A^2 = G_B^2 = G$ and $p_1 \neq p_2$. To avoid cases of indifference, we assume that $v \neq 1/2$. The judge follows the recommendation of the majority. In this simple model, a juror prefers b to a if the juror’s posterior belief is below $1/2$ and a to b if the juror’s posterior belief is above $1/2$. Type- k jurors set a threshold x_k and recommend a if $c \geq x_k$ and b if $c < x_k$.

Assume $v < 1/2$. Then, type-2 jurors always obtain their preferred choice. The problem is equivalent to the case in which type-2 jurors are delegated the decision. Their threshold is therefore x_2^* . When type-1 jurors disagree with type-2 jurors, they are outnumbered. The degree of disagreement is also irrelevant to them as it does not affect the decision. If type-1 jurors are a priori relatively more in favor of action b , that is, $x_1^* < x_2^*$ then any threshold $\tilde{x}_1 \leq x_2^*$ is an equilibrium. Overall, type-1 jurors do not have any clear strategy, as their recommendation will never be followed, and we may observe either weak or strong polarization.

In the second example, we consider a situation in which the recommendation of one group is followed with some probability but the

outcome is attenuated to reflect that the decision is not unanimous. This corresponds to a situation in which jurors deliberate to reach a verdict and need to agree on a punishment. To keep the analysis simple, we consider the same preferences and beliefs as in the first example. However, we assume that, in case of disagreement, type-1 jurors convince the jury to recommend their attenuated preferred decision with probability η and the benefit of taking the correct, but attenuated, action is $g < G$ instead of G .

If η tends to zero, then type-1 jurors never obtain their preferred choice and we are back to the first example. When η tends to one and g tends to G , type-1 jurors always obtain their preferred decision and the optimal threshold is x_1^* . For intermediate cases, jurors need to take into account that they may not be able to voice their opinion, and even if they do, the final verdict will be attenuated: the person may be sentenced to a smaller punishment than the one recommended or be imposed a fine instead of being acquitted. These inefficient outcomes occur whenever jurors disagree. To decrease the likelihood of such outcomes, it becomes optimal to move the threshold in the direction of the threshold of type-2 jurors. The same argument applies to type-2 jurors and, at equilibrium, the thresholds are \tilde{x}_1 and \tilde{x}_2 such that $x_1^* < \tilde{x}_1 < \tilde{x}_2 < x_2^*$. The main message of this section is summarized in the next proposition:

Proposition 5: *The interpretation of the evidence depends on the rule.*

Even though the two examples above are very stylized, they illustrate the effect of the rule on the likelihood of polarization in juries. The rule affects the payoffs obtained in case of disagreement, and the optimal cutoff internalizes this externality. Therefore, jurors will process the information as a function of the rule, end up recommending different verdicts and hold different beliefs.

C. JURY SELECTION

The previous results demonstrate that the interpretation of the evidence produced in front of a jury varies as a function of the prior beliefs and the intensity of the preferences of the jurors. As such, jury selection affects the overall outcome of a trial. Some experts believe that a large share of cases litigated are won or lost in the jury-selection phase.⁴² During this process, attorneys will select jurors as they compete to secure opposite verdicts.

42. See, e.g., Herald P. Fahringer, "Mirror, Mirror on the Wall . . .": *Body Language, Intuition, and the Art of Jury Selection*, 17 AM. J. TRIAL ADVOC. 197, 197-98 (1993).

Informally, suppose there are two attorneys α and β . Attorney j 's preferences can be represented by a utility function over decisions $H_j(\cdot)$ that does not depend on the state realized. Attorney α 's utility is such that $H_\alpha(a) > H_\alpha(b) = 0$, while attorney β 's utility is such that $H_\beta(b) > H_\beta(a) = 0$ independently of the state. These utilities induce an indirect preference over juries. For instance, the defense attorney (here β) would like to secure a release (action b) and is willing to avoid jurors likely to interpret the evidence in favor of a conviction. The defense attorney's objective is to maximize the number of jurors with high thresholds and to minimize those with low thresholds. The prosecutor (α) has the opposite preferences and incentives.

Consider now the case with a default option, and let us assume that $H_\alpha(o) = H_\beta(o) = 0$ to reflect the fact that attorneys care exclusively about winning. From the perspective of α , for example, it is optimal to eliminate jurors with a strong bias against A and to keep jurors with a strong bias in favor of A . As attorneys have opposite incentives, they will eliminate the jurors with the strongest views on both sides of the spectrum. Overall, the adversarial jury selection process will keep jurors with moderate views. As a consequence and compared to the initial draw of jurors, the jurors selected in this way will (1) vote unanimously more often and (2) polarize their opinion less often after the release of evidence.

V. CONCLUSION

In this Article, we have built on neuroscience evidence to model belief formation and analyze the behavior of judges and juries. We have shown that physiological constraints generate posterior beliefs with qualitatively different properties from Bayesian posterior beliefs. In particular, decisionmakers will tend to reinforce their prior beliefs and to hold posterior beliefs shaped by their preferences over outcomes. Also, the well-known polarization effect is a direct consequence of the model and should be observed when evidence is mixed, as behavioral evidence suggests.

There are a series of implications for decisions rendered by judges and juries. We have shown that cases analyzed in the early career of a judge may affect future decisions on cases that are a priori independent. Also, early evidence produced in a trial may matter more than late evidence. For the case of juries, the distribution of preferences in a jury affects the way information is interpreted by individual jurors. Finally, we argue that the endogenous selection of jury members reduces the likelihood of polarization and split opinions.

This analysis could be extended in various ways. For instance, the results obtained in Part IV presuppose that all information is public knowledge. In the presence of uncertainty about the preferences and prior beliefs of other members in the jury, the thresholds should reflect the distribution of preferences and beliefs rather than the exact values. We have also restricted ourselves to situations in which decisionmakers want to take the correct action in each state. This provides the most likely scenario for avoiding biases and therefore constitutes the most natural benchmark. However, the results obtained may not fit all the available data on the polarization effect. If jurors have other types of preferences, or the rule in case of disagreement is different, then the group interaction may have different effects on posterior beliefs. This is one interesting alley for future research.

Our model abstracts from other important issues. For instance, we have assumed that decisionmakers (judges or jurors) can only determine whether a crime has been committed, but not the extent of the damage or the magnitude of the punishment. The ability to assess punitive damage has been studied in several articles, and it is not clear, for instance, whether judges or juries should be delegated these assessments.⁴³ More generally, we have not addressed the optimal design of legal procedures.⁴⁴ For example, our model could be extended to compare in greater detail the performance of judges and juries from the perspective of a planner with a specific objective.⁴⁵

43. See, e.g., A. Mitchell Polinsky & Steven Shavell, *Punitive Damages: An Economic Analysis*, 111 HARV. L. REV. 869, 891–93 (1998); Cass R. Sunstein, Daniel Kahneman & David Schkade, *Assessing Punitive Damages (with Notes on Cognition and Valuation in Law)*, 107 YALE L.J. 2071, 2094–99 (1998).

44. Recent articles have investigated how biases in judgment could be mitigated in different but related settings. See, e.g., Benjamin Lester, Nicola Persico & Ludo Visschers, *Information Acquisition and the Exclusion of Evidence in Trials*, 28 J.L. ECON. & ORG. 163 (2012) (analyzing how evidence should be optimally excluded when jurors suffer from a specific form of bounded rationality).

45. For an experimental analysis, see generally W. Kip Viscusi, *Jurors, Judges, and the Mistreatment of Risk by the Courts*, 30 J. LEGAL STUD. 107 (2001).

APPENDIX

In this appendix, we offer a brief description of the mathematical model and a sketch of some important results. We refer the reader to our 2012a¹ and 2012b² studies for additional details and long proofs.

In our 2012a study, information is modeled as follows:

When the state is S , the likelihood of c is $f(c|S)$ with $F(c|S) = \int_0^c f(y|S)dy$

representing the probability of a cell firing activity below c . A high-cell firing is more likely to occur when $S=A$, and a low-cell firing is more likely to occur when $S=B$, which is captured by the Monotone Likelihood Ratio Property $\frac{\partial}{\partial c} \left(\frac{f(c|B)}{f(c|A)} \right) < 0$ for all c .

A. EFFICIENT DECISIONMAKING

In the standard Bayesian framework, information is encoded and interpreted to its full extent. Namely, c represents the correct intensity of the signal and the decision is based on the realization of c .

Consider the decision of individual i who holds a posterior belief μ_i that the state is A after receiving evidence. Given the posterior belief, i 's expected payoffs of taking action $\gamma_i=a$ and $\gamma_i=b$ respectively are $V_i(a; \mu_i) = \mu_i G_A^i$ and $V_i(b; \mu_i) = (1 - \mu_i) G_B^i$.

The optimal decision is therefore

$$\gamma_i^* = \begin{cases} a & \text{if } \mu_i \geq \mu_i^* \equiv \frac{G_B^i}{G_A^i + G_B^i} \\ b & \text{if } \mu_i < \mu_i^* \end{cases} \quad (1)$$

1. See Isabelle Brocas & Juan D. Carrillo, *From Perception to Action: An Economic Model of Brain Processes*, 78 GAMES & ECON. BEHAV. 81, 83, 91–93 (2012).

2. See Isabelle Brocas & Juan D. Carrillo, *The Neuroeconomics of Strategic Decision-Making* (2012) (unpublished manuscript) (on file with author).

Suppose that i is exposed to a signal of intensity c . Conditional on interpreting the signal at its exact value, i 's posterior belief is therefore:

$$\gamma_i^* = \begin{cases} a & \text{if } \mu_i \geq \mu_i^* \equiv \frac{G_B^i}{G_A^i + G_A^i} \\ b & \text{if } \mu_i < \mu_i^* \end{cases} \quad (2)$$

$\pi_i(c)$ is increasing in c , so there exists x_i^* such that $\pi_i(c) < \mu_i^*$ when $c < x_i^*$ and $\pi_i(c) \geq \mu_i^*$ when $c \geq x_i^*$. Formally x_i^* satisfies

$$\frac{f(x_i^*|B)}{f(x_i^*|A)} = \frac{p_i}{1-p_i} \cdot \frac{G_A^i}{G_B^i} \quad (3)$$

B. DECISIONMAKING IN THE THRESHOLD MECHANISM

Given this mechanism, i can hold either of two posterior beliefs, $p_i(x) \equiv p_i$ or $p_i(x) \equiv \underline{p}_i(x)$, depending on whether $c \geq x$ or $c < x$. Formally,

$$\overline{p}_i(x) = \frac{[1-F(x|A)]p_i}{[1-F(x|A)]p_i + [1-F(x|B)](1-p_i)}; \quad \underline{p}_i(x) = \frac{F(x|A)p_i}{F(x|A)p_i + F(x|B)(1-p_i)}$$

The expected payoff associated with the threshold mechanism is therefore,

$$W(x) = Pr(c \geq x)V_i(a; \overline{p}_i(x)) + Pr(c < x)V_i(b; \underline{p}_i(x))$$

$$W(x) = p_i G_A^i (1-F(x|A)) + (1-p_i) G_B^i F(x|B)$$

The optimal threshold solves $W'(x) = 0$ and satisfies

$$\frac{f(x_i^*|B)}{f(x_i^*|A)} = \frac{p_i}{1-p_i} \cdot \frac{G_A^i}{G_B^i} \quad (4)$$

It is easy to see that actions a and b are taken under the same circumstances. However, the posterior beliefs are different. To establish Proposition 1, it is sufficient to differentiate the first order condition with respect to p_i and G_A^i / G_B^i to see that x_i^* is decreasing in the prior belief p_i and in the relative payoff G_A^i / G_B^i . We obtain Proposition 2 by direct inspection of $\overline{p}_i(x_i^*)$ and $\underline{p}_i(x_i^*)$ given equation (4). Now, given some information is filtered out, the prior belief used for a subsequent episode of information interpretation is a biased posterior. Given this, decisionmaking

in subsequent periods will not be efficient from a purely Bayesian perspective. Furthermore, the evidence used to compute the biased posterior will be reinforced in a subsequent period, and the order in which evidence is produced matters. The details of the argument and proof of Proposition 3 can be found in our 2012a study.³

C. JURIES IN THE THRESHOLD MODEL

Fix type-2 jurors' threshold x_2 . The expected payoff of type-1 choosing $x_1 < x_2$ is

$$\begin{aligned} W_1(x_1, x_2) &= \\ & Pr(c < x_1)(1 - p_1(x_1))G_B^1 + Pr(c \in (x_1, x_2))G + Pr(c > x_2) \overline{p}_1(x_2)G_A^1 \\ & = F(x_1|B)(1 - p_1)G_B^1 + p_1[F(x_2|A) - F(x_1|A)]G \\ & \quad + (1 - p_1)[F(x_2|B) - F(x_1|B)]G + [1 - F(x_2|A)]p_1G_1^A \end{aligned}$$

The optimal threshold for type-1 jurors satisfies the first order condition $\partial W_1 / \partial x_1 = 0$. The first order condition writes as

$$f(x_1|B)(1 - p_1)[G_B^1 - G] - p_1 f(x_1|A)G = 0$$

The solution to this equation is x_1^* . Similar calculations apply for type-2 optimal threshold: we find it is x_2^* . These arguments prove Proposition 4. A more detailed argument can be found in our 2012b study.⁴

3. See Brocas & Carrillo, *supra* note 1, at 93.

4. See Brocas & Carrillo, *supra* note 2.

D. MAJORITARIAN RULES

Type- k agents set a threshold x_k and recommend a if $c \geq x_k$ and b if $c < x_k$. If $v < 1/2$, type-2 agents are delegated the decision and set threshold x_2^* . Let us assume that $x_1 < x_2^*$ and let $p_1(x_1, x_2^*)$ be the posterior belief that the true state is A , conditional on the evidence being in (x_1, x_2^*) . The expected payoff of a type-1 juror is,

$$\begin{aligned} W_1(x_1, x_2^*) = & \\ & Pr(c > x_2^*)G \overline{p}_1(x_2^*) + Pr(c \in (x_1, x_2^*))G(1 - p_1(x_1, x_2^*)) + Pr(c < x_1)G(1 - p_1(x_1)) \\ & = (1 - p_1)F(x_2^*|B)G + p_1(1 - F(x_2^*|A))G \end{aligned}$$

Therefore the solution is any $\tilde{x}_1 \in (0, x_2^*)$. This proves Proposition 5.

E. DELIBERATION AND ATTENUATED VERDICTS

Fix type-2 jurors' threshold x_2 . The expected payoff of type-1 choosing $x_1 < x_2$ is

$$\begin{aligned} W_1(x_1, x_2) = & Pr(c < x_1)(1 - p_1(x_1))G + Pr(c > x_2) \overline{p}_1(x_2)G \\ & + Pr(c \in (x_1, x_2))g [\eta p_1(x_1, x_2) + (1 - \eta)(1 - p_1(x_1, x_2))] \end{aligned}$$

The optimal threshold solves $\partial W_1 / \partial x_1 = 0$ and satisfies

$$\frac{f(x_1|B)}{f(x_1|A)} = \frac{p_1}{1 - p_1} \frac{g\eta}{G - (1 - \eta)g} \quad (5)$$

It is sufficient to differentiate this expression with respect to η to see that the optimal threshold is decreasing in η . Similarly, if we differentiate this expression with respect to g , we find that the optimal threshold is decreasing in g . The optimal threshold is $\tilde{x}_1 > x_1^*$, and it exists if $\tilde{x}_1 < x_2$. Similarly, the expected payoff of type-2 choosing $x_2 > x_1$ is

$$\begin{aligned} W_2(x_1, x_2) = & Pr(c < x_1)(1 - p_2(x_1))G + Pr(c > x_2) \overline{p}_2(x_2)G \\ & + Pr(c \in (x_1, x_2))g [\eta p_2(x_1, x_2) + (1 - \eta)(1 - p_2(x_1, x_2))] \end{aligned}$$

The optimal threshold solves $\partial W_2 / \partial x_2 = 0$ and satisfies

$$\frac{f(x_2|B)}{f(x_2|A)} = \frac{p_2}{1-p_2} \frac{G-g\eta}{(1-\eta)g} \quad (6)$$

The optimal threshold is increasing in g and decreasing in η . The optimal threshold is $\tilde{x}_2 < x_2^*$, and it exists if $\tilde{x}_2 > x_1$. There exists an equilibrium if $\tilde{x}_1 < \tilde{x}_2$, that is if the parameters satisfy

$$\frac{p_1}{1-p_1} \frac{g\eta}{(G-(1-\eta)g)} > \frac{p_2}{1-p_2} \frac{G-g\eta}{(1-\eta)g}$$

For instance, this is satisfied if $p_2 = 0$ and $p_1 > 0$.

