

## Informal Sanctions on Prosecutors and Defendants and the Disposition of Criminal Cases

by Andrew F. Daughety and Jennifer F. Reinganum

### Technical Appendix

The basic notation and the building blocks of the payoff functions are given in the text; some of this material is repeated here for easy reference. The D of type  $t$ 's payoff function from trial is:

$$\pi_T^D(t) = S_c(1 - F_t) + k^D + r^D\mu(G | c)(1 - F_t) + r^D\mu(G | a)F_t, t \in \{I, G\}. \quad (\text{TA.1})$$

D's payoff from accepting a plea bargain of  $S_b$  is:

$$\pi_b^D = S_b + r^D\mu(G | b). \quad (\text{TA.2})$$

D's expected payoff following rejection (given his type) is:

$$\pi_R^D(t) = \rho^P\pi_T^D(t) + (1 - \rho^P)\pi_d^D, \quad (\text{TA.3})$$

where  $\pi_d^D = r^D\mu(G | d)$ .

The prosecutor's payoff from going to trial (given her beliefs following the defendant's rejection of her plea offer) can be written as:

$$\begin{aligned} \pi_T^P = & v(G | R)\{S_c(1 - F_G) - k^P - r_I^P\mu(I | c)(1 - F_G) - r_G^P\mu(G | a)F_G\} \\ & + v(I | R)\{S_c(1 - F_I) - k^P - r_I^P\mu(I | c)(1 - F_I) - r_G^P\mu(G | a)F_I\}. \end{aligned} \quad (\text{TA.4})$$

P's payoff from dropping the case is simply:

$$\pi_d^P = -r_G^P\mu(G | d). \quad (\text{TA.5})$$

P's expected payoff following rejection is given by:

$$\pi_R^P = \rho^P\pi_T^P + (1 - \rho^P)\pi_d^P. \quad (\text{TA.6})$$

#### *A Preliminary Result*

**Remark 1.**  $\pi_T^D(I) < \pi_T^D(G)$ . That is, an innocent defendant expects a smaller loss at trial than a guilty defendant (for given beliefs on the part of the observers and P).

Proof. First, note that for arbitrary positive values of  $\rho_G^D$  and  $\rho_I^D$ , the observers' posterior probability of guilt is higher following a conviction than following an acquittal if and only if  $F_I > F_G$ . More formally,  $\mu(G | c) (>, =, <) \mu(G | a)$  as  $F_I (>, =, <) F_G$ . To see this, notice that, by Bayes' Rule,

$$\mu(G | c) = \rho_G^D(1 - \lambda)(1 - F_G) / [\rho_G^D(1 - \lambda)(1 - F_G) + \rho_I^D\lambda(1 - F_I)], \quad (\text{TA.7})$$

whereas

$$\mu(G | a) = \rho_G^D(1 - \lambda)F_G / [\rho_G^D(1 - \lambda)F_G + \rho_I^D\lambda F_I]. \quad (\text{TA.8})$$

Simple though tedious algebra indicates that  $\mu(G | c) > \mu(G | a)$  if and only if  $F_I > F_G$ , which is a maintained assumption. As  $\rho_I^D$  goes to zero, both  $\mu(G | c)$  and  $\mu(G | a)$  go to 1, whereas as  $\rho_G^D$  goes to zero, both  $\mu(G | c)$  and  $\mu(G | a)$  go to 0.

Second, the expression for  $\pi_I^D(t)$  can be differentiated with respect to  $F_I$  to obtain:

$$\partial \pi_I^D(t) / \partial F_I = -S_c - r^D \{ \mu(G | c) - \mu(G | a) \}.$$

Since the term in curly brackets has been proved to be non-negative, the entire expression is negative. Since  $F_I > F_G$ , the result follows. QED

Note that an innocent defendant and a guilty defendant expect the same loss if they accept the plea offer of  $S_b$  (for given beliefs on the part of the observers and P). This is because accepting the plea bargain results in case disposition b, which yields a payoff of  $\pi_A^D = S_b + r^D \mu(G | b)$ , independent of D's true type.  $S_b$  is the formal sanction, whereas  $r^D \mu(G | b)$  is the informal sanction imposed by the observer, who believes that D's type is G with probability  $\mu(G | b)$  if he accepts the plea bargain.

### *Maintained Restrictions*

MR1.  $(S_c - r_I^P)(1 - F_I) - k^P < 0$ . This means that, if it were common knowledge (or commonly-believed) that D is innocent, then P would prefer to drop the case rather than proceed to trial.

To see this, substitute  $v(G | R) = \mu(G | a) = \mu(G | d) = 0$  and  $v(I | R) = \mu(I | c) = 1$  into the formulas for  $\pi_T^P$  and  $\pi_d^P$ . This yields  $\pi_T^P = S_c(1 - F_I) - k^P - r_I^P(1 - F_I)$ , whereas  $\pi_d^P = 0$ . Then MR1 implies that  $\pi_T^P < \pi_d^P$ .

MR2.  $(1 - \lambda)[(S_c + r_G^P)(1 - F_G) - k^P] + \lambda[(S_c - r_I^P)(1 - F_I) - k^P] > 0$ . This means that, if it were common knowledge (or commonly-believed) that the fraction of guilty defendants among those that rejected the plea offer is  $1 - \lambda$ , then P would prefer to take the case to trial rather than drop it.

To see this, substitute  $v(G | R) = 1 - \lambda$  and  $v(I | R) = \lambda$  into the formula for  $\pi_T^P$ . Moreover, substitute  $\mu(G | a) = (1 - \lambda)F_G / [(1 - \lambda)F_G + \lambda F_I]$  and  $\mu(I | c) = \lambda(1 - F_I) / [(1 - \lambda)(1 - F_G) + \lambda(1 - F_I)]$  into the formula for  $\pi_T^P$ . This yields:  $\pi_T^P = [(1 - \lambda)(1 - F_G) + \lambda(1 - F_I)]S^T - r_I^P\lambda(1 - F_I) - r_G^P(1 - \lambda)F_G - k^P$ . Finally, substitute  $\mu(G | d) = 1 - \lambda$  into the formula for  $\pi_d^P$  to obtain  $\pi_d^P = -r_G^P(1 - \lambda)$ . Then MR2 implies that  $\pi_T^P > \pi_d^P$ .

MR1 and MR2 imply the following, but we state it here for easy reference.

MR3.  $(S_c + r_G^P)(1 - F_G) - k^P > 0$ . This means that, if it were common knowledge (or commonly-believed) that D is guilty, then P would prefer to take the case to trial rather than drop it.

To see this, substitute  $v(G | R) = \mu(G | a) = \mu(G | d) = 1$  and  $v(I | R) = \mu(I | c) = 0$  into the formulas for  $\pi_T^P$  and  $\pi_d^P$ . This yields  $\pi_T^P = \{S_c(1 - F_G) - k^P - r_G^P F_G\}$ , whereas  $\pi_d^P = -r_G^P$ . Then MR1 and MR2 imply that  $\pi_T^P < \pi_d^P$ .

The following remark considers arbitrary mixed strategies and the value to P of going to trial rather than dropping the case.

Remark 2. For arbitrary mixing probabilities  $(\rho_G^D, \rho_I^D)$ , where  $\rho_I^D$  is the probability that type t rejects the plea offer, the expression  $(\pi_T^P - \pi_d^P)$  (i.e., the difference between P's payoff from taking the case to trial and dropping it) is decreasing in  $\rho_I^D$  and increasing in  $\rho_G^D$ . Also,  $\text{num}(\pi_T^P - \pi_d^P)$  is decreasing

in  $\rho_1^D$  and increasing in  $\rho_G^D$ .

Proof. For arbitrary mixing probabilities  $(\rho_G^D, \rho_1^D)$ , P's payoff from trial is given by:

$$\begin{aligned} \pi_T^P = & v(G | R) \{S_c(1 - F_G) - k^P - r_1^P \mu(I | c)(1 - F_G) - r_G^P \mu(G | a)F_G\} \\ & + v(I | R) \{S_c(1 - F_I) - k^P - r_1^P \mu(I | c)(1 - F_I) - r_G^P \mu(G | a)F_I\}, \quad (\text{TA.9}) \end{aligned}$$

where  $v(G | R) = \rho_G^D(1 - \lambda)/[\rho_G^D(1 - \lambda) + \rho_1^D\lambda]$ ;  $\mu(I | c) = \rho_1^D\lambda(1 - F_I)/[\rho_G^D(1 - \lambda)(1 - F_G) + \rho_1^D\lambda(1 - F_I)]$ ; and  $\mu(G | a) = \rho_G^D(1 - \lambda)F_G/[\rho_G^D(1 - \lambda)F_G + \rho_1^D\lambda F_I]$ . P's payoff from dropping the case is given by  $\pi_d^P = -r_G^P \mu(G | d)$ , where  $\mu(G | d) = \rho_G^D(1 - \lambda)/[\rho_G^D(1 - \lambda) + \rho_1^D\lambda]$ . After much substitution, it can be shown that  $\text{num}(\pi_T^P - \pi_d^P) = S_c[\rho_G^D(1 - \lambda)(1 - F_G) + \rho_1^D\lambda(1 - F_I)] - r_1^P\rho_1^D\lambda(1 - F_I) - k^P[\rho_G^D(1 - \lambda) + \rho_1^D\lambda] + r_G^P\rho_G^D(1 - \lambda)(1 - F_G)$  and  $\text{denom}(\pi_T^P - \pi_d^P) = \rho_G^D(1 - \lambda) + \rho_1^D\lambda$ . Differentiation, algebra, and MR1-MR3 yield the result that  $(\pi_T^P - \pi_d^P)$  goes down as  $\rho_1^D$  goes up and  $(\pi_T^P - \pi_d^P)$  goes up as  $\rho_G^D$  goes up. Considering only  $\text{num}(\pi_T^P - \pi_d^P)$ , this clearly goes down as  $\rho_1^D$  goes up (by MR1) and  $\text{num}(\pi_T^P - \pi_d^P)$  goes up as  $\rho_G^D$  goes up (by MR3). QED

### *Candidates for Equilibria*

There are 9 distinct candidate forms for equilibria. Candidates 1-4 are pure-strategy equilibria; that is, each type of D plays a particular strategy with probability 1. Candidates 5-9 involve at least one type of D mixing between accepting and rejecting the plea offer.

1. Types I and G accept the plea offer.
2. Type I rejects the plea offer, whereas type G accepts it.
3. Type I accepts the plea offer, whereas type G rejects it.
4. Types I and G reject the plea offer.
5. Type I rejects the plea offer, whereas type G mixes.
6. Type G accepts the plea offer, whereas type I mixes.

7. Type G rejects the plea offer, whereas type I mixes.
8. Type I accepts the plea offer, whereas type G mixes.
9. Both types mix.

We argue that the only candidate forms that can actually be equilibria are Candidates 4 and 5. We postpone characterization of these equilibria until after we dispose of those candidate forms that cannot be equilibria.

1. Types I and G accept the plea offer. In this putative equilibrium, the dispositions {a, c, d} are all out-of-equilibrium events. What should P believe if D unexpectedly rejects the plea offer? And what should observers believe if they unexpectedly observe a disposition of a, c or d? Since type I expects a lower loss from trial than type G, whereas both types expect the same loss from a dropped case, the equilibrium refinement D1 (Cho and Kreps, 1987) implies that unexpected rejection of the plea offer (and unexpected dispositions a, c or d) should be assigned to type I. This is because type I would be willing to risk a larger probability of trial  $\rho^P$  to defect from accepting to rejecting the plea offer than would type G. Formally, this means that in this putative equilibrium,  $v(G | R) = \mu(G | a) = \mu(G | c) = \mu(G | d) = 0$ . Consequently,  $\pi_T^P = S_c(1 - F_I) - k^P - r_I^P(1 - F_I)$  and  $\pi_d^P = 0$ ; by MR1, this means that P will prefer to drop the case. Basically, if both types are expected to accept the plea bargain, then rejecting it is taken as a clear signal of innocence and P will therefore drop the case (P does not have to worry about informal sanctions from dropping the case, as observers take this disposition as a clear signal of innocence). But if P will drop the case following rejection of the plea offer, then both types will defect from this putative equilibrium to rejecting the plea offer. Thus, there cannot be an equilibrium of this form.

2. Type I rejects the plea offer, whereas type G accepts it. In this putative equilibrium, the

observer's beliefs following the dispositions {a, c, d} (which could only occur following a rejection of the plea offer) are that D is surely of type I. Moreover, P also believes that a rejection implies type I. Thus,  $v(G | R) = \mu(G | a) = \mu(G | c) = \mu(G | d) = 0$ . Consequently,  $\pi_T^P = S_c(1 - F_I) - k^P - r_1^P(1 - F_I)$  and  $\pi_d^P = 0$ ; by MR1, this means that P will prefer to drop the case. Basically, if type G is expected to accept the plea bargain, then rejecting it is taken as a clear signal of innocence and P will therefore drop the case (P does not have to worry about informal sanctions from dropping the case, as observers take this disposition as a clear signal of innocence). Again, if P is expected to drop the case following a rejection of the plea offer, then the type G defendant will defect from this putative equilibrium to rejecting the plea offer. Thus, there cannot be an equilibrium of this form.

3. Type I accepts the plea offer, whereas type G rejects it. In this putative equilibrium, the observer's beliefs following the dispositions {a, c, d} (which could only occur following a rejection of the plea offer) are that D is surely of type G. Moreover, P also believes that a rejection implies type G. By MR3, P will prefer to take the case to trial rather than to drop it. Recall that – holding beliefs constant – type I faces a lower expected cost of trial than does type G, whereas they expect the same loss by accepting the plea bargain. This implies that if type G prefers trial to the plea bargain (or is indifferent), then type I must strictly prefer trial to the plea bargain. Therefore, the type I defendant will defect from this putative equilibrium to rejecting the plea offer. Thus, there cannot be an equilibrium of this form.

6. Type G accepts the plea offer, whereas type I mixes. The argument is exactly the same as for candidate 2, and will be omitted.

7. Type G rejects the plea offer, whereas type I mixes. Recall that – holding beliefs constant – type I faces a lower expected cost of trial than does type G, whereas they expect the same loss by

accepting the plea bargain (or by having the case dropped). If type I is indifferent between accepting the plea offer and rejecting it, then type G must strictly prefer the plea bargain, as long as P takes the case to trial with positive probability following rejection. Recall that MR2 ensures that if all type G's and all type I's rejected the plea offer, P would take the case to trial (rather than dropping it). Since the mixture of defendant types in this putative equilibrium puts more weight on type G (and less on type I) relative to the prior, Remark 2 implies that trial is even more attractive to P (relative to dropping the case), so P would take the case to trial. This implies that if type G prefers to reject the plea bargain (or is indifferent), then type I must strictly prefer to reject it. Therefore, the type I defendant will defect from this putative equilibrium to rejecting the plea offer. Thus, there cannot be an equilibrium of this form.

8. Type I accepts the plea offer, whereas type G mixes. In this putative equilibrium, any rejection of the plea offer, and any outcome  $\{a, c, d\}$  are attributed to type G. By MR3, P would take the case to trial following a rejection. If type G is mixing, then he must be indifferent between the plea offer and trial. But since type I expects a smaller loss than type G from trial (but the same loss as G from accepting the plea bargain), it must be that type I strictly prefers trial to accepting the plea bargain. Therefore, the type I defendant will defect from this putative equilibrium to rejecting the plea offer. Thus, there cannot be an equilibrium of this form.

9. Both types mix. Again – holding beliefs constant – type I faces a lower expected cost of trial than does type G, whereas they expect the same loss by accepting the plea bargain (or by having the case dropped). If type I is indifferent between accepting the plea offer and rejecting it, then type G must strictly prefer the plea bargain, as long as P takes the case to trial with positive probability following rejection, which would cause type G to defect from this putative equilibrium to accepting

the plea.

Only if P drops the case with probability one following rejection can both defendant types be made indifferent. Such an equilibrium can be ruled out if we assume that, when both types are indifferent, they reject the plea offer (or are believed, by both P and the observers, to reject the plea offer) with the same probability. In this case, the mixture of rejecting types is (believed to be) the same as the prior, and P prefers to take the case to trial rather than dropping it. But then it cannot be that both D types are indifferent, as type I expects a smaller loss from trial than type G. Alternatively, a stronger version of MR2 could guarantee that P prefers to make a demand that is unacceptable to both types, rather than dropping the case against both types. While provoking rejection by both types does not change the observers' beliefs (because rejection is on the equilibrium path), it does change P's posterior beliefs. Her payoff from trial is now:

$$\begin{aligned} \bar{\pi}_T^P = & (1 - \lambda)\{S_c(1 - F_G) - k^P - r_I^P\mu(I | c)(1 - F_G) - r_G^P\mu(G | a)F_G\} \\ & + \lambda\{S_c(1 - F_I) - k^P - r_I^P\mu(I | c)(1 - F_I) - r_G^P\mu(G | a)F_I\}, \end{aligned}$$

where  $\mu(I | c) = \rho_I^D\lambda(1 - F_I)/[\rho_G^D(1 - \lambda)(1 - F_G) + \rho_I^D\lambda(1 - F_I)]$  and  $\mu(G | a) = \rho_G^D(1 - \lambda)F_G/[\rho_G^D(1 - \lambda)F_G + \rho_I^D\lambda F_I]$ . P's payoff from dropping the case is still given by  $\pi_d^P = -r_G^P\mu(G | d)$ , where  $\mu(G | d) = \rho_G^D(1 - \lambda)/[\rho_G^D(1 - \lambda) + \rho_I^D\lambda]$ . If  $[(1 - \lambda)(1 - F_G) + \lambda(1 - F_I)]S_c$  is assumed to be sufficiently large, then P will prefer to provoke trial against both types rather than dropping the case against both types. This will also undermine a putative equilibrium wherein both types of defendant mix between accepting and rejecting the plea offer.

### *Characterizing Equilibria*

The only remaining candidate forms for an equilibrium are Candidates 4 (types I and G reject the plea offer) and 5 (type I rejects the plea offer, whereas type G mixes). Thus, an innocent

defendant rejects the plea offer with probability one, but a guilty defendant may accept the plea offer with positive probability. Because Candidate 4 is a limiting case of Candidate 5, we can focus on Candidate 5.

The timing of the game is such that each type of D chooses to accept or reject the plea offer, taking as given the likelihood that P takes the case to trial following rejection; and P chooses to take the case to trial or drop it, given her beliefs about the posterior probability that D is of type G, given rejection. Both of these decisions are taken following P's choice of plea offer,  $S_b$ , so both parties must take this offer as given at subsequent decision nodes.

We first characterize the equilibrium in the continuation game, given  $S_b$ , allowing for mixed strategies for both P ( $\rho^p$ ) and the D of type G, ( $\rho_G^D$ ; type I will always reject the plea offer in this putative equilibrium). Since the observers' beliefs will depend on their conjectured value for  $\rho_G^D$ , we will augment the notation for the observers' beliefs to reflect these conjectures. Other functions that also depend on these conjectures through the observers' beliefs will be similarly augmented.

Suppose that observers conjecture that the D of type G rejects the plea offer with probability  $\rho_G^{D\theta}$ . Then  $\mu(G | c; \rho_G^{D\theta}) = \rho_G^{D\theta}(1 - \lambda)(1 - F_G)/[\rho_G^{D\theta}(1 - \lambda)(1 - F_G) + \lambda(1 - F_I)]$ ;  $\mu(G | a; \rho_G^{D\theta}) = \rho_G^{D\theta}(1 - \lambda)F_G/[\rho_G^{D\theta}(1 - \lambda)F_G + \lambda F_I]$ ;  $\mu(G | d; \rho_G^{D\theta}) = \rho_G^{D\theta}(1 - \lambda)/[\rho_G^{D\theta}(1 - \lambda) + \lambda]$ ; and  $\mu(G | b; \rho_G^{D\theta}) = 1$ . Moreover, suppose that type G anticipates these beliefs, and also expects that P will take the case to trial following rejection with probability  $\rho^p$ . Then type G will be indifferent, and hence willing to mix, between accepting and rejecting the offer  $S_b$  if  $\pi_R^D(G; \rho_G^{D\theta}) = \rho^p \pi_T^D(G; \rho_G^{D\theta}) + (1 - \rho^p) \pi_d^D(\rho_G^{D\theta}) = \pi_b^D(\rho_G^{D\theta})$ . That is, if:  $\rho^p[\pi_T^D(G; \rho_G^{D\theta}) - \pi_d^D(\rho_G^{D\theta})] = \pi_b^D(\rho_G^{D\theta}) - \pi_d^D(\rho_G^{D\theta})$ . Substitution and simplification yields:

$$\rho^p \{S_c(1 - F_G) + k^D + r^D \mu(G | c; \rho_G^{D\theta})(1 - F_G) + r^D \mu(G | a; \rho_G^{D\theta})F_G\} + (1 - \rho^p)r^D \mu(G | d; \rho_G^{D\theta}) = S_b + r^D.$$

Upon collecting terms, the value of  $\rho^P$  that results in this equality is:

$$\rho^P(S_b; \rho_G^{D\Theta}) = \frac{\{S_b + r^D(1 - \mu(G | d; \rho_G^{D\Theta}))\}}{\{S_c(1 - F_G) + k^D + r^D[\mu(G | c; \rho_G^{D\Theta})(1 - F_G) + \mu(G | a; \rho_G^{D\Theta})F_G - \mu(G | d; \rho_G^{D\Theta})]\}} \quad (\text{TA.10})$$

The numerator of the expression  $\rho^P(S_b; \rho_G^{D\Theta})$ , which is the difference between type G's payoff from accepting the plea offer versus having his case dropped, is clearly positive, meaning that D would prefer to have his case dropped. The denominator of the expression  $\rho^P(S_b; \rho_G^{D\Theta})$  is the difference between type G's payoff from trial versus having his case dropped. This denominator is also positive (see Remark 3 below), which implies that type G would prefer that P drop the case against him rather than take it to trial.

Remark 3. The denominator of the expression  $\rho^P(S_b; \rho_G^{D\Theta})$  is positive.

Proof. A sufficient condition for the denominator to be positive is that

$$\begin{aligned} & \mu(G | c; \rho_G^{D\Theta})(1 - F_G) + \mu(G | a; \rho_G^{D\Theta})F_G - \mu(G | d; \rho_G^{D\Theta}) \\ & = [\mu(G | c; \rho_G^{D\Theta}) - \mu(G | d; \rho_G^{D\Theta})](1 - F_G) + [\mu(G | a; \rho_G^{D\Theta}) - \mu(G | d; \rho_G^{D\Theta})]F_G > 0. \end{aligned}$$

Recall that  $\mu(G | c; \rho_G^{D\Theta}) = \rho_G^{D\Theta}(1 - \lambda)(1 - F_G)/[\rho_G^{D\Theta}(1 - \lambda)(1 - F_G) + \lambda(1 - F_I)]$ ;  $\mu(G | a; \rho_G^{D\Theta}) = \rho_G^{D\Theta}(1 - \lambda)F_G/[\rho_G^{D\Theta}(1 - \lambda)F_G + \lambda F_I]$ ; and  $\mu(G | d; \rho_G^{D\Theta}) = \rho_G^{D\Theta}(1 - \lambda)/[\rho_G^{D\Theta}(1 - \lambda) + \lambda]$ . Let  $X = [\rho_G^{D\Theta}(1 - \lambda)(1 - F_G) + \lambda(1 - F_I)]$  and let  $Y = [\rho_G^{D\Theta}(1 - \lambda)F_G + \lambda F_I]$ ; so  $X + Y = [\rho_G^{D\Theta}(1 - \lambda) + \lambda]$ . Then:

$$\begin{aligned} & [\mu(G | c; \rho_G^{D\Theta}) - \mu(G | d; \rho_G^{D\Theta})](1 - F_G) + [\mu(G | a; \rho_G^{D\Theta}) - \mu(G | d; \rho_G^{D\Theta})]F_G \\ & = \{[\rho_G^{D\Theta}(1 - \lambda)(1 - F_G)/X] - [\rho_G^{D\Theta}(1 - \lambda)/(X + Y)]\}(1 - F_G) \\ & \quad + \{[\rho_G^{D\Theta}(1 - \lambda)F_G/Y] - [\rho_G^{D\Theta}(1 - \lambda)/(X + Y)]\}F_G > 0 \end{aligned}$$

if and only if (after some algebra)  $[(1 - F_G)(X + Y)](1 - F_G)Y + [F_G(X + Y) - Y]F_GX > 0$ , which can be verified by substituting for X, Y and X + Y, collecting terms, and recalling that  $F_I > F_G$ . QED

Since the observers' beliefs are based on their conjectures  $\rho_G^{D\Theta}$  and the case disposition, and NOT on  $S_b$ , which they do not observe, the expression  $\rho^P(S_b; \rho_G^{D\Theta})$  is an increasing function of  $S_b$ .

That is, when  $S_b$  is higher, P must take the case to trial following rejection with a higher probability in order to make the D of type G indifferent about accepting or rejecting  $S_b$ . Notice that even a plea offer of  $S_b = 0$  requires a positive probability of trial following a rejection in order to induce the D of type G to be willing to accept it; this is because acceptance of a plea offer comes with a sure informal sanction of  $r^D$  (as only a truly guilty D is expected to accept the plea).

Now consider P's decision about trying versus dropping the case. Again suppose that observers – and P – both conjecture that type G rejects the plea offer with probability  $\rho_G^{D\ominus}$  in this candidate for equilibrium; thus  $v(G | R; \rho_G^{D\ominus}) = \rho_G^{D\ominus}(1 - \lambda)/[\rho_G^{D\ominus}(1 - \lambda) + \lambda]$ . Since these conjectures must be the same (and correct) in equilibrium, it is valid to equate them at this point in order to identify what common beliefs will make P indifferent, and hence willing to mix, between trying and dropping the case following a rejection. P will be indifferent between these two options if  $\pi_T^P(\rho_G^{D\ominus}) = \pi_d^P(\rho_G^{D\ominus})$ ; that is, if:

$$\begin{aligned} & v(G | R; \rho_G^{D\ominus}) \{S_c(1 - F_G) - k^P - r_1^P \mu(I | c; \rho_G^{D\ominus})(1 - F_G) - r_G^P \mu(G | a; \rho_G^{D\ominus}) F_G\} \\ & + v(I | R; \rho_G^{D\ominus}) \{S_c(1 - F_I) - k^P - r_1^P \mu(I | c; \rho_G^{D\ominus})(1 - F_I) - r_G^P \mu(G | a; \rho_G^{D\ominus}) F_I\} \\ & = - r_G^P \mu(G | d; \rho_G^{D\ominus}). \end{aligned} \quad (\text{TA.11})$$

Substituting for the beliefs and simplifying yields (see also the proof of Remark 2):

$$\begin{aligned} \text{num}(\pi_T^P(\rho_G^{D\ominus}) - \pi_d^P(\rho_G^{D\ominus})) &= S_c[\rho_G^{D\ominus}(1 - \lambda)(1 - F_G) + \lambda(1 - F_I)] - r_1^P \lambda(1 - F_I) \\ & - k^P[\rho_G^{D\ominus}(1 - \lambda) + \lambda] + r_G^P \rho_G^{D\ominus}(1 - \lambda)(1 - F_G), \end{aligned} \quad (\text{TA.12})$$

and  $\text{denom}(\pi_T^P(\rho_G^{D\ominus}) - \pi_d^P(\rho_G^{D\ominus})) = [\rho_G^{D\ominus}(1 - \lambda) + \lambda]$ . The expression  $\text{num}(\pi_T^P(\rho_G^{D\ominus}) - \pi_d^P(\rho_G^{D\ominus}))$  is increasing in  $\rho_G^{D\ominus}$  by MR3. Moreover, we know from MR1 that this expression is negative for  $\rho_G^{D\ominus} = 0$  (that is, when a D of type G is never expected to reject), and we know from MR2 that it is positive for  $\rho_G^{D\ominus} = 1$  (that is, when a D of type G is always expected to reject). Therefore, the value of  $\rho_G^{D\ominus}$  that will

make P indifferent between trying and dropping the case is given by:

$$\rho_G^{D0} = -\lambda[(S_c - r_1^P)(1 - F_1) - k^P]/(1 - \lambda)[(S_c + r_G^P)(1 - F_G) - k^P], \quad (\text{TA.13})$$

where the numerator is positive by MR1; the denominator is positive by MR3; and the ratio is a fraction by MR2. For any  $\rho_G^{D0} > \rho_G^{D0}$ , P will strictly prefer to take the case to trial following a rejection, and for any  $\rho_G^{D0} < \rho_G^{D0}$ , P will strictly prefer to drop the case following a rejection.

To summarize, type G is willing to mix between accepting and rejecting the plea offer  $S_b$  if he anticipates that the observers' beliefs are  $\rho_G^{D0} = \rho_G^{D0}$  and he expects that P will take the case to trial following rejection of offer  $S_b$  with probability  $\rho^P(S_b; \rho_G^{D0})$ . P is indifferent between trying and dropping the case if she (and the observers) believes that type G rejects the plea offer with probability  $\rho_G^{D0}$ . Thus, the mixed-strategy equilibrium, given  $S_b$ , is  $(\rho_G^{D0}, \rho^P(S_b; \rho_G^{D0}))$ .

We can now move back to the decision node at which P chooses the plea offer  $S_b$ , anticipating that it will be following by the mixed-strategy equilibrium  $(\rho_G^{D0}, \rho^P(S_b; \rho_G^{D0}))$  in the continuation game. P's payoff from making the plea offer  $S_b$  is:

$$(1 - \rho_G^{D0})(1 - \lambda)S_b + (\rho_G^{D0}(1 - \lambda) + \lambda)[\rho^P(S_b; \rho_G^{D0})\pi_T^P(\rho_G^{D0}) + (1 - \rho^P(S_b; \rho_G^{D0}))\pi_d^P(\rho_G^{D0})]. \quad (\text{TA.14})$$

The set of feasible  $S_b$  values is bounded below by 0 and above by  $S_b = \pi_T^D(G; \rho_G^{D0}) - r^D$ , where  $\pi_T^D(G; \rho_G^{D0})$  is the expression for  $\pi_T^D(G)$ , evaluated at the beliefs  $\mu(G | c; \rho_G^{D0}) = \rho_G^{D0}(1 - \lambda)(1 - F_G) / [\rho_G^{D0}(1 - \lambda)(1 - F_G) + \lambda(1 - F_1)]$ ; and  $\mu(G | a; \rho_G^{D0}) = \rho_G^{D0}(1 - \lambda)F_G / [\rho_G^{D0}(1 - \lambda)F_G + \lambda F_1]$ . This is because accepting the plea offer results in a combined sanction of  $S_b + r^D$  (since only guilty D's accept the plea offer) and thus any plea offer higher than  $\pi_T^D(G; \rho_G^{D0}) - r^D$  will be rejected for sure (rather than with probability  $\rho_G^{D0}$ ). At this upper bound, the function  $\rho^P(S_b; \rho_G^{D0})$  just reaches 1. In order to have a non-empty feasible range, we need  $\pi_T^D(G; \rho_G^{D0}) - r^D \geq 0$ ; or, equivalently,  $r^D[1 - \mu(G | c; \rho_G^{D0})(1 - F_G) - \mu(G | a; \rho_G^{D0})F_G] \leq S_c(1 - F_G) + k^D$ . Since the term in brackets on the left-hand-side can be re-written

as  $(1 - F_G)(1 - \mu(G | c; \rho_G^{D0})) + F_G(1 - \mu(G | a; \rho_G^{D0}))$ , it is clearly positive.

Condition 1. In order for P to be able to induce a D of type G to accept a plea offer, it must be that

$$r^D \leq [S_c(1 - F_G) + k^D] / [1 - \mu(G | c; \rho_G^{D0})(1 - F_G) - \mu(G | a; \rho_G^{D0})F_G].$$

The expression  $r^D[1 - \mu(G | c; \rho_G^{D0})(1 - F_G) - \mu(G | a; \rho_G^{D0})F_G]$  is the increment in informal sanctions that the D of type G suffers by accepting a plea (which only a true G is expected to do) rather than going to trial (where there is a chance of conviction and a chance of acquittal, with corresponding informal sanctions). If there were no informal sanctions for D, then Condition 1 would be satisfied automatically. Thus, informal sanctions on D constrain P's ability to settle cases via plea bargain.

Returning to P's payoff as a function of  $S_b$  (i.e., equation (TA.14)), notice two things. First, since  $\rho_G^{D0}$ , which is independent of  $S_b$ , renders P indifferent between trying and dropping the case following rejection, the term in square brackets simply equals  $\pi_d^P(\rho_G^{D0}) = -r_G^P \mu(G | d; \rho_G^{D0})$ , where  $\mu(G | d; \rho_G^{D0}) = \rho_G^{D0}(1 - \lambda) / [\rho_G^{D0}(1 - \lambda) + \lambda]$ . Thus, the optimal  $S_b$  that supports some plea bargaining is the upper limit of the feasible range,  $S_b(\rho_G^{D0}) = \pi_T^D(G; \rho_G^{D0}) - r^D$ ; this is rejected by type G with probability  $\rho_G^{D0}$ , and P goes to trial with certainty following a rejection. Note that a D of type I would always reject this plea offer, consistent with the hypothesized form of the equilibrium.

Every plea offer in the feasible set  $[0, \pi_T^D(G; \rho_G^{D0}) - r^D]$  is consistent with a mixed-strategy equilibrium in which some G-types accept, and others reject, the offer. But P could make a higher demand that would provoke certain rejection. We need to verify that P prefers the hypothesized equilibrium described above to the "defection payoff" she would obtain if all cases went to trial.

In the hypothesized equilibrium, P settles with  $(1 - \rho_G^{D0})(1 - \lambda)$  guilty defendants and goes to trial against the rest of the guilty defendants and all of the innocent defendants; if P defects and

provokes rejection by all, then she will simply replace the settlement  $S_b(\rho_G^{D0}) = \pi_T^D(G; \rho_G^{D0}) - r^D$  with the expected payoff from taking a guilty defendant to trial (holding the observers' beliefs fixed at the levels implied by  $\rho_G^{D0}$ , because trial is already on the equilibrium path). Thus, P prefers (at least weakly) the hypothesized equilibrium to defection as long as:

$$\begin{aligned} \pi_T^D(G; \rho_G^{D0}) - r^D &= S_c(1 - F_G) + k^D + r^D\mu(G | c; \rho_G^{D0})(1 - F_G) + r^D\mu(G | a; \rho_G^{D0})F_G - r^D \\ &\geq S_c(1 - F_G) - k^P - r_1^P\mu(I | c; \rho_G^{D0})(1 - F_G) - r_G^P\mu(G | a; \rho_G^{D0})F_G. \end{aligned} \quad (\text{TA.15})$$

Rearranging, we can write this as:

$$r^D[1 - \mu(G | c; \rho_G^{D0})(1 - F_G) - \mu(G | a; \rho_G^{D0})F_G] \leq k^P + k^D + r_1^P\mu(I | c; \rho_G^{D0})(1 - F_G) + r_G^P\mu(G | a; \rho_G^{D0})F_G.$$

Condition 2. For P to find it preferable to settle with a D of type G rather than provoking a trial, it must be that:

$$r^D \leq [k^P + k^D + r_1^P\mu(I | c; \rho_G^{D0})(1 - F_G) + r_G^P\mu(G | a; \rho_G^{D0})F_G] / [1 - \mu(G | c; \rho_G^{D0})(1 - F_G) - \mu(G | a; \rho_G^{D0})F_G].$$

Again, if D faced no informal sanctions, then Condition 2 would be satisfied. High informal sanctions on D can undermine P's desire to settle using plea bargaining.

Finally, P could also defect by dropping all cases; again, this would not change the observers' beliefs, but we need to verify that P prefers the hypothesized equilibrium outcome to what she would get by defecting to dropping all cases. However, Condition 1 is sufficient to imply this preference. To see why, notice that in the hypothesized equilibrium, P's payoff is:

$$(1 - \rho_G^{D0})(1 - \lambda)[\pi_T^D(G; \rho_G^{D0}) - r^D] + (\rho_G^{D0}(1 - \lambda) + \lambda)\pi_T^P(\rho_G^{D0}). \quad (\text{TA.16})$$

We already know that  $\pi_T^P(\rho_G^{D0}) = \pi_d^P(\rho_G^{D0})$  by construction (and  $\pi_T^P(\rho_G^D) > \pi_d^P(\rho_G^D)$  for  $\rho_G^D > \rho_G^{D0}$ ). Then Condition 1 implies that the settlement offer  $S_c(\rho_G^{D0}) = [\pi_T^D(G; \rho_G^{D0}) - r^D]$  is non-negative, whereas P's payoff from dropping a case is  $-\pi_T^P(\rho_G^{D0})$ , which is strictly negative.

Both Conditions 1 and 2 are restrictions on  $r^D$ ; however, we have been unable to determine

which right-hand-side provides the tighter constraint.

*Comparative Statics for the Selected Equilibrium*

Here we summarize comparative static effects of parameter changes in  $r_1^p$ ,  $r_G^p$ ,  $r^D$ ,  $k^p$ ,  $k^D$ ,  $F_G$ ,  $F_I$ ,  $\lambda$ , and  $S_c$  on equilibrium strategies such as the plea offer and the likelihood of plea bargaining success. Recall that the likelihood of plea bargaining failure is  $\rho_G^{D0}$ , where

$$\rho_G^{D0} = -\lambda[(S_c - r_1^p)(1 - F_I) - k^p]/(1 - \lambda)[(S_c + r_G^p)(1 - F_G) - k^p], \quad (\text{TA.17})$$

and the equilibrium plea offer is

$$S_b(\rho_G^{D0}) = \pi_1^D(G; \rho_G^{D0}) - r^D = S_c(1 - F_G) + k^D + r^D\mu(G | c; \rho_G^{D0})(1 - F_G) + r^D\mu(G | a; \rho_G^{D0})F_G - r^D. \quad (\text{TA.18})$$

First, we consider the impact of changes in the parameters on  $\rho_G^{D0}$ . Recall that a higher value of  $\rho_G^D$  makes trial more attractive to P (relative to dropping the case) following a rejection, and  $\rho_G^{D0}$  makes P indifferent between these two decisions (even though she goes to trial with probability one in equilibrium). Therefore, any parameter change that would tip P toward one decision or the other must be counter-balanced by a change in  $\rho_G^{D0}$  that restore's P's indifference.

$$\partial\rho_G^{D0}/\partial\lambda = -[(S_c - r_1^p)(1 - F_I) - k^p]/(1 - \lambda)^2[(S_c + r_G^p)(1 - F_G) - k^p] > 0;$$

$$\partial\rho_G^{D0}/\partial r_1^p = \lambda(1 - F_I)/(1 - \lambda)[(S_c + r_G^p)(1 - F_G) - k^p] > 0;$$

$$\partial\rho_G^{D0}/\partial k^p = -\lambda\{[(S_c - r_1^p)(1 - F_I) - k^p] - [(S_c + r_G^p)(1 - F_G) - k^p]\}/(1 - \lambda)[(S_c + r_G^p)(1 - F_G) - k^p]^2 > 0;$$

$$\partial\rho_G^{D0}/\partial F_G = -\lambda[(S_c - r_1^p)(1 - F_I) - k^p](S_c + r_G^p)/(1 - \lambda)[(S_c + r_G^p)(1 - F_G) - k^p]^2 > 0.$$

An increase in  $\lambda$  (the fraction of innocent among those arrested),  $r_1^p$  (the sanction rate for punishing an innocent defendant),  $k^p$  (P's cost of trial), or  $F_G$  (the probability that a guilty defendant is acquitted) has the direct effect of making trial less attractive, so the fraction of guilty types in the pool of those rejecting must increase to restore P's willingness to go to trial.

$$\partial\rho_G^{D0}/\partial r_G^p = \lambda[(S_c - r_1^p)(1 - F_I) - k^p](1 - F_G)/(1 - \lambda)[(S_c + r_G^p)(1 - F_G) - k^p]^2 < 0.$$

$$\partial \rho_G^{D0} / \partial S_c = \frac{\{(1 - \lambda)[(S_c + r_G^P)(1 - F_G) - k^P](-\lambda(1 - F_I)) + \lambda[(S_c - r_I^P)(1 - F_I) - k^P](1 - \lambda)(1 - F_G)\}}{\{(1 - \lambda)[(S_c + r_G^P)(1 - F_G) - k^P]\}^2} < 0.$$

An increase in  $r_G^P$  (the sanction rate for failing to punish a guilty defendant) or  $S_c$  (the formal sanction) has the direct effect of making trial more attractive (relative to dropping the case), so the fraction of guilty types in the pool of those rejecting can decrease and yet maintain P's willingness to go to trial. Some comparative statics are ambiguous (i.e., they could go either way, depending on the relative magnitude of parameters) or we are unable to determine the direction of the impact.

$$\partial \rho_G^{D0} / \partial F_I = \lambda(S_c - r_I^P) / (1 - \lambda)[(S_c + r_G^P)(1 - F_G) - k^P] > 0 \text{ (resp., } < 0) \text{ if } (S_c - r_I^P) > 0 \text{ (resp., } < 0);$$

Finally, as  $r^D$  and  $k^D$  do not appear in P's payoffs, they do not have an impact on  $\rho_G^{D0}$ ; that is,  $\partial \rho_G^{D0} / \partial k^D = 0$  and  $\partial \rho_G^{D0} / \partial r^D = 0$ .

Now consider the impact of parameter changes on  $S_b(\rho_G^{D0})$ , which is an increasing function.

There are several parameters that affect the equilibrium plea offer only indirectly through  $\rho_G^{D0}$ .

$$\partial S_b(\rho_G^{D0}) / \partial \lambda = S_b'(\rho_G^{D0})(\partial \rho_G^{D0} / \partial \lambda) > 0;$$

$$\partial S_b(\rho_G^{D0}) / \partial r_I^P = S_b'(\rho_G^{D0})(\partial \rho_G^{D0} / \partial r_I^P) > 0;$$

$$\partial S_b(\rho_G^{D0}) / \partial k^P = S_b'(\rho_G^{D0})(\partial \rho_G^{D0} / \partial k^P) > 0;$$

$$\partial S_b(\rho_G^{D0}) / \partial r_G^P = S_b'(\rho_G^{D0})(\partial \rho_G^{D0} / \partial r_G^P) < 0;$$

$$\partial S_b(\rho_G^{D0}) / \partial F_I = S_b'(\rho_G^{D0})(\partial \rho_G^{D0} / \partial F_I) > 0 \text{ (resp., } < 0) \text{ if } (S_c - r_I^P) > 0 \text{ (resp., } < 0).$$

The parameters  $r^D$  and  $k^D$  affect the equilibrium offer only directly, as  $\rho_G^{D0}$  does not depend on them.

$$\partial S_b(\rho_G^{D0}) / \partial r^D = \mu(G | c; \rho_G^{D0})(1 - F_G) + \mu(G | a; \rho_G^{D0})F_G - 1 < 0;$$

$$\partial S_b(\rho_G^{D0}) / \partial k^D = 1 > 0;$$

The parameters  $F_G$  and  $S_c$  affect the plea offer both directly and indirectly.

$$\partial S_b(\rho_G^{D0}) / \partial S_c = (1 - F_G) + S_b'(\rho_G^{D0})(\partial \rho_G^{D0} / \partial S_c) = ???, \text{ as the direct effect is positive and the}$$

indirect effect is negative.

$\partial S_b(\rho_G^{D0})/\partial F_G = \{-S_c - r^D[\mu(G | c; \rho_G^{D0}) - \mu(G | a; \rho_G^{D0})]\} + S_b'(\rho_G^{D0})(\partial \rho_G^{D0}/\partial F_G) = ???$ , as the direct effect (in curly brackets) is negative and the indirect effect is positive.

Recall that the right-hand-side of Condition 1 is:

$$[S_c(1 - F_G) + k^D]/[1 - \mu(G | c; \rho_G^{D0})(1 - F_G) - \mu(G | a; \rho_G^{D0})F_G].$$

Since the denominator is decreasing in  $\rho_G^{D0}$ , and  $\rho_G^{D0}$  is increasing (resp., decreasing) in  $r_1^P$  (resp.,  $r_G^P$ ), it follows that the r.h.s. of Condition 1 is increasing in  $r_1^P$  and decreasing in  $r_G^P$ .

The right-hand-side of Condition 2 is:

$$[k^P + k^D + r_1^P \mu(I | c; \rho_G^{D0})(1 - F_G) + r_G^P \mu(G | a; \rho_G^{D0})F_G]/[1 - \mu(G | c; \rho_G^{D0})(1 - F_G) - \mu(G | a; \rho_G^{D0})F_G].$$

Again, the denominator is decreasing in  $r_1^P$ , and decreasing in  $r_G^P$ . Since the numerator is increasing in  $r_1^P$  (the only questionable term is  $r_1^P \mu(I | c; \rho_G^{D0})$  because  $\mu(I | c; \rho_G^{D0})$  is decreasing, but overall this product is increasing in  $r_1^P$ ), it follows that the r.h.s. of Condition 2 is increasing in  $r_1^P$ . We are unable to determine the impact of an increase in  $r_G^P$  on the r.h.s. of Condition 2, because both the numerator and the denominator are increasing in  $r_G^P$ .

Finally, we show that an increase in  $S_c$  decreases the expected informal sanctions facing a D of type I, and increases the expected informal sanctions facing a D of type G. To see this, we first write the ex ante expected informal sanctions facing a D of unknown type. This is given by:

$$r^D \lambda \{(1 - F_I) \mu(G | c; \rho_G^{D0}) + F_I \mu(G | a; \rho_G^{D0})\} \\ + r^D (1 - \lambda) \{(1 - \rho_G^{D0}) + \rho_G^{D0} (1 - F_G) \mu(G | c; \rho_G^{D0}) + \rho_G^{D0} F_G \mu(G | a; \rho_G^{D0})\}.$$

The terms in the first line represent the contribution to ex ante expected informal sanctions generated by the D of type I (who never accepts the plea), and the terms in the second line represent the contribution generated by the D of type G (who sometimes accepts the plea). Collecting coefficients on the expressions  $\mu(G | c; \rho_G^{D0})$  and  $\mu(G | a; \rho_G^{D0})$  yields the following version:

$$r^D \{\lambda(1 - F_I) + (1 - \lambda)\rho_G^{D0}(1 - F_G)\} \mu(G | c; \rho_G^{D0}) + r^D \{\lambda F_I + (1 - \lambda)\rho_G^{D0} F_G\} \mu(G | a; \rho_G^{D0}) + r^D(1 - \lambda)(1 - \rho_G^{D0}).$$

Upon recalling the form of the beliefs,  $\mu(G | c; \rho_G^{D0})$  and  $\mu(G | a; \rho_G^{D0})$ , we see that the expressions in curly brackets above are the same as the denominators in the beliefs they multiply; thus, the expression above reduces to:

$$r^D \{(1 - \lambda)\rho_G^{D0}(1 - F_G)\} + r^D \{(1 - \lambda)\rho_G^{D0} F_G\} + r^D(1 - \lambda)(1 - \rho_G^{D0}) = r^D(1 - \lambda).$$

That is, the ex ante expected informal sanctions facing a D of unknown type are independent of  $S_c$ .

Since it is clear that the expected informal sanctions facing a D of type I, which are given by

$r^D \{(1 - F_I)\mu(G | c; \rho_G^{D0}) + F_I \mu(G | a; \rho_G^{D0})\}$ , are decreasing in  $S_c$  (because  $\mu(G | c; \rho_G^{D0})$  and  $\mu(G | a; \rho_G^{D0})$  are both increasing in  $\rho_G^{D0}$  and  $\rho_G^{D0}$  is decreasing in  $S_c$ ), then the fact that the ex ante expected informal sanctions facing a D of unknown type are independent of  $S_c$  implies that the expected informal sanctions facing a D of type G, given by  $r^D \{(1 - \rho_G^{D0}) + \rho_G^{D0}(1 - F_G)\mu(G | c; \rho_G^{D0}) + \rho_G^{D0} F_G \mu(G | a; \rho_G^{D0})\}$ , must be increasing in  $S_c$ .

### *The Scottish Verdict: Three-Outcome Regime*

The three outcomes are “guilty” (denoted g), “not guilty” (denoted ng) and “not proven” (denoted np). The relevant threshold for a finding of g is denoted  $\gamma_g$  (this is assumed to be the same as the threshold  $\gamma_c$  for conviction in the two-outcome regime), so  $1 - F_t(\gamma_g)$ ,  $t \in \{I, G\}$ , is the probability that the D of type t is found guilty. Similarly, the threshold for a finding of ng is denoted  $\gamma_{ng}$ , so  $F_t(\gamma_{ng})$ ,  $t \in \{I, G\}$ , is the probability that the D of type t is found not guilty. Finally, the probability the D of type t receives a verdict of not proven is given by  $\Delta_t \equiv F_t(\gamma_g) - F_t(\gamma_{ng})$ ,  $t \in \{I, G\}$ .

For the three-outcome regime, we assume the Strict Monotone Likelihood Ratio Property (SMLRP). That is,  $f(e | G)/f(e | I)$  is strictly increasing in  $e$  for  $e \in (0, 1)$ . The assumption of

SMLRP implies Strict First-Order Stochastic Dominance; that, in the two-outcome regime,  $F_G(e) < F_I(e)$  for all  $e \in (0, 1)$ . Evaluating at  $e = \gamma_g$  yields  $F_G(\gamma_g) < F_I(\gamma_g)$ ; that is, an innocent D is more likely to be acquitted at trial than a guilty D. SMLRP further implies the following relationships that will be used in the three-outcome regime.

Strict Reverse Hazard Rate Dominance (SRHRD):  $f_G(e)/F_G(e) > f_I(e)/F_I(e)$  for all  $e \in (0, 1)$ .

Strict Hazard Rate Dominance (SHRD):  $f_G(e)/[1 - F_G(e)] < f_I(e)/[1 - F_I(e)]$  for all  $e \in (0, 1)$ .

The effect of dividing the former “acquittal” evidence interval into two sub-intervals corresponding to “not proven” and “not guilty” is to change the D of type t’s payoff function from trial to the following form:

$$\pi_t^D(t) = S_c(1 - F_t(\gamma_g)) + k^D + r^D\mu(G | g)(1 - F_t(\gamma_g)) + r^D\mu(G | np)\Delta_t + r^D\mu(G | ng)F_t(\gamma_{ng}). \text{(TA.19)}$$

Since  $\gamma_g = \gamma_c$ , the effect is basically to replace the expression  $r^D\mu(G | a)F_t(\gamma_c)$  with  $r^D\mu(G | np)\Delta_t + r^D\mu(G | ng)F_t(\gamma_{ng})$ .

For arbitrary mixing probabilities  $(\rho_I^D, \rho_G^D)$ , the beliefs are now:

$$\mu(G | g) = \rho_G^D(1 - \lambda)(1 - F_G(\gamma_g))/[\rho_G^D(1 - \lambda)(1 - F_G(\gamma_g)) + \rho_I^D\lambda(1 - F_I(\gamma_g))];$$

$$\mu(G | np) = \rho_G^D(1 - \lambda)\Delta_G/[\rho_G^D(1 - \lambda)\Delta_G + \rho_I^D\lambda\Delta_I];$$

$$\mu(G | ng) = \rho_G^D(1 - \lambda)F_G(\gamma_{ng})/[\rho_G^D(1 - \lambda)F_G(\gamma_{ng}) + \rho_I^D\lambda F_I(\gamma_{ng})].$$

A sufficient condition for  $\pi_t^D(G) > \pi_t^D(I)$  is that  $\mu(G | ng) \leq \mu(G | np) \leq \mu(G | g)$ . First, notice that  $\mu(G | np) \geq \mu(G | ng)$  if and only if  $F_I(\gamma_{ng})/F_G(\gamma_{ng}) \geq \Delta_I/\Delta_G = [F_I(\gamma_g) - F_I(\gamma_{ng})]/[F_G(\gamma_g) - F_G(\gamma_{ng})]$  or, equivalently, if and only if  $F_I(\gamma_{ng})/F_G(\gamma_{ng}) \geq F_I(\gamma_g)/F_G(\gamma_g)$ . These expressions are equal at  $\gamma_{ng} = \gamma_g$ , and SMLRP (SRHRD) implies that the ratio  $F_I(e)/F_G(e)$  is strictly decreasing in  $e$ . Thus,  $F_I(\gamma_{ng})/F_G(\gamma_{ng}) > F_I(\gamma_g)/F_G(\gamma_g)$  for all  $\gamma_{ng} < \gamma_g$ . Next, notice that  $\mu(G | g) \geq \mu(G | np)$  if and only if  $\Delta_I/\Delta_G = [F_I(\gamma_g) - F_I(\gamma_{ng})]/[F_G(\gamma_g) - F_G(\gamma_{ng})] \geq [1 - F_I(\gamma_g)]/[1 - F_G(\gamma_g)]$  or, equivalently, if and only if  $[1 - F_I(\gamma_{ng})]/[1 -$

$F_G(\gamma_{ng})] \geq [1 - F_I(\gamma_g)]/[1 - F_G(\gamma_g)]$ . These expressions are equal at  $\gamma_{ng} = \gamma_g$ , and SMLRP (SHRD) implies that the ratio  $[1 - F_I(e)]/[1 - F_G(e)]$  is strictly decreasing in  $e$ . Thus,  $[1 - F_I(\gamma_{ng})]/[1 - F_G(\gamma_{ng})] > [1 - F_I(\gamma_g)]/[1 - F_G(\gamma_g)]$  for all  $\gamma_{ng} < \gamma_g$ . We therefore conclude that  $\mu(G | ng) < \mu(G | np) < \mu(G | g)$  and thus  $\pi_T^D(G) > \pi_T^D(I)$ .

Since type G expects a worse outcome at trial than does type I, the equilibrium will be of the same form as before; that is, all Ds of type I will go to trial, along with a fraction of Ds of type G, denoted  $\rho_G^D$ . The plea offer will make a D of type G indifferent about accepting the plea deal and going to trial. We will again select the lowest value of  $\rho_G^D$  consistent with incentivizing P to go to trial rather than dropping the case following a rejected plea offer. Incorporating  $\Theta$ 's beliefs and P's beliefs (which are given by  $v(G | R; \rho_G^{D\Theta}) = \rho_G^{D\Theta}(1 - \lambda)/[\rho_G^{D\Theta}(1 - \lambda) + \lambda]$ ), we can write P's expected payoff from trial as follows:

$$\begin{aligned} & v(G | R; \rho_G^{D\Theta}) \{ S_c(1 - F_G(\gamma_g)) - k^P - r_1^P \mu(I | g; \rho_G^{D\Theta})(1 - F_G(\gamma_g)) - r_G^P \mu(G | np; \rho_G^{D\Theta}) \Delta_G \\ & \quad - r_G^P \mu(G | ng; \rho_G^{D\Theta}) F_G(\gamma_{ng}) \} \\ & + v(I | R; \rho_G^{D\Theta}) \{ S_c(1 - F_I(\gamma_g)) - k^P - r_1^P \mu(I | g; \rho_G^{D\Theta})(1 - F_I(\gamma_g)) - r_G^P \mu(G | np; \rho_G^{D\Theta}) \Delta_I \\ & \quad - r_G^P \mu(G | ng; \rho_G^{D\Theta}) F_I(\gamma_{ng}) \}. \end{aligned} \quad (TA.19)$$

Substituting for the beliefs and collecting terms yields:

$$\begin{aligned} \text{num}(\pi_T^P(\rho_G^{D\Theta})) &= S_c[\rho_G^{D\Theta}(1 - \lambda)(1 - F_G(\gamma_g)) + \lambda(1 - F_I(\gamma_g))] - r_1^P \lambda(1 - F_I(\gamma_g)) \\ & \quad - k^P[\rho_G^{D\Theta}(1 - \lambda) + \lambda] - r_G^P \rho_G^{D\Theta}(1 - \lambda) F_G(\gamma_g), \end{aligned} \quad (TA.20)$$

and  $\text{denom}(\pi_T^P(\rho_G^{D\Theta})) = [\rho_G^{D\Theta}(1 - \lambda) + \lambda]$ . Notice that P's expected payoff from trial is independent of the fact that the acquittal interval has been subdivided into intervals pertaining to outcomes of not proven and not guilty. Moreover, since P's payoff from dropping the case is still  $= -r_G^P \mu(G | d; \rho_G^{D\Theta})$ , it follows that the mixing probability for the D of type G that just makes P indifferent between trial

and dropping the case is exactly the same as in the two-verdict case:

$$\rho_G^{D0} = -\lambda[(S_c - r_1^p)(1 - F_1(\gamma_g)) - k^p]/(1 - \lambda)[(S_c + r_G^p)(1 - F_G(\gamma_g)) - k^p]. \quad (\text{TA.22})$$

Although the form of the equilibrium plea offer is still the same,  $S_b(\rho_G^{D0}) = \pi_T^D(G; \rho_G^{D0}) - r^D$ , recall that the function  $\pi_T^D(G; \rho_G^{D0})$  in the three-outcome regime replaces the expression  $r^D\mu(G | a)F_G(\gamma_c)$  with  $r^D\mu(G | np)\Delta_G + r^D\mu(G | ng)F_G(\gamma_{ng})$ . Because the beliefs are, in both regimes, evaluated at the same value of  $\rho_G^{D0}$  (and because  $\gamma_g = \gamma_c$ ), we only need to compare  $\mu(G | np)\Delta_G + \mu(G | ng)F_G(\gamma_{ng})$  with  $\mu(G | a)F_G(\gamma_g)$  in order to determine whether the equilibrium plea offer is higher or lower under the three-outcome regime. It will be useful to write:

$$\mu(G | a) = \rho_G^{D0}(1 - \lambda)F_G(\gamma_g)/A, \text{ where } A \equiv [\rho_G^{D0}(1 - \lambda)F_G(\gamma_g) + \lambda F_1(\gamma_g)];$$

$$\mu(G | ng) = \rho_G^{D0}(1 - \lambda)F_G(\gamma_{ng})/B, \text{ where } B \equiv [\rho_G^{D0}(1 - \lambda)F_G(\gamma_{ng}) + \lambda F_1(\gamma_{ng})]; \text{ and}$$

$$\mu(G | np) = \rho_G^{D0}(1 - \lambda)\Delta_G/C, \text{ where } C \equiv [\rho_G^{D0}(1 - \lambda)\Delta_G + \lambda\Delta_1].$$

Then  $\mu(G | np)\Delta_G + \mu(G | ng)F_G(\gamma_{ng}) > \mu(G | a)F_G(\gamma_g)$  if and only if  $[(\Delta_G)^2/C] + [(F_G(\gamma_{ng}))^2/B] > [(F_G(\gamma_g))^2/A]$ , which holds if and only if  $[F_G(\gamma_g)B - F_G(\gamma_{ng})A]^2 > 0$ . The term in brackets is nonzero because the ratio  $F_G(e)/[\rho_G^{D0}(1 - \lambda)F_G(e) + \lambda F_1(e)]$  is increasing by RHRD. Thus, the D of type G faces a higher expected punishment at trial under the three-outcome regime than under the two-outcome regime, and this also results in P making a higher plea offer in equilibrium.

On the other hand, a D of type I will face a lower expected punishment at trial under the three-outcome regime than under the two-outcome regime if  $\mu(G | np)\Delta_1 + \mu(G | ng)F_1(\gamma_{ng}) < \mu(G | a)F_1(\gamma_g)$  or, equivalently, if and only if  $(1 - \mu(I | np))\Delta_1 + (1 - \mu(I | ng))F_1(\gamma_{ng}) < (1 - \mu(I | a))F_1(\gamma_g)$ . This inequality holds if and only if  $[(\Delta_1)^2/C] + [(F_1(\gamma_{ng}))^2/B] > [(F_1(\gamma_g))^2/A]$ , which holds if and only if  $[F_1(\gamma_g)B - F_1(\gamma_{ng})A]^2 > 0$ . The term in brackets is nonzero because the ratio  $F_1(e)/[\rho_G^{D0}(1 - \lambda)F_G(e) + \lambda F_1(e)]$  is strictly decreasing by SRHRD.

Both Conditions 1 and 2 are easier to fulfill in the three-outcome regime. This is because the right-hand-side of Condition 1 becomes:

$$S_c(1 - F_G) + k^D/[1 - \mu(G | g; \rho_G^{D0})(1 - F_G(\gamma_g)) - \mu(G | np; \rho_G^{D0})\Delta_G - \mu(G | ng; \rho_G^{D0})F_G(\gamma_{ng})],$$

and we have just shown that this denominator is smaller than the corresponding expression under the two-verdict regime. The denominator in the right-hand-side of Condition 2 is the same as in the right-hand-side of Condition 1, and (as argued above) this has become smaller with the addition of the third outcome. The numerator in the right-hand-side of Condition 2 is now:

$$k^P + k^D + r_1^P \mu(I | g; \rho_G^{D0})(1 - F_G(\gamma_g)) + r_G^P \mu(G | np; \rho_G^{D0})\Delta_G + r_G^P \mu(G | ng; \rho_G^{D0})F_G(\gamma_{ng}),$$

which is larger than in the two-outcome regime. Thus, both Conditions 1 and 2 hold for larger ranges of the parameter  $r^D$ .

Finally, the outside observers' expected loss from misclassification is lower under the three-outcome regime.

$$\begin{aligned} M(\rho_G^D) &= \lambda(1 - F_I(\gamma_g))r^D \mu(G | g) + \lambda r^D [\mu(G | np)\Delta_I + \mu(G | ng)F_I(\gamma_{ng})] \text{ (this term is lower)} \\ &+ \rho_G^D(1 - \lambda)(1 - F_G(\gamma_g))[r^D - r^D \mu(G | g)] \text{ (this term is the same)} \\ &+ \rho_G^D(1 - \lambda) \{ \Delta_G[r^D - r^D \mu(G | np)] + F_G(\gamma_{ng})[r^D - r^D \mu(G | ng)] \} \text{ (this term is lower)} \\ &+ \lambda(1 - F_I(\gamma_g))[r_1^P - r_1^P \mu(I | g)] + \lambda r_G^P [\mu(G | np)\Delta_I + \mu(G | ng)F_I(\gamma_{ng})] \text{ (this term is lower)} \\ &+ \rho_G^D(1 - \lambda)(1 - F_G(\gamma_g))[r_1^P \mu(I | g)] \text{ (this term is the same)} \\ &+ \rho_G^D(1 - \lambda) \{ \Delta_G[r_G^P - r_G^P \mu(G | np)] + F_G(\gamma_{ng})[r_G^P - r_G^P \mu(G | ng)] \}. \text{ (this term is lower)} \end{aligned}$$

To summarize, we find that the form of the equilibrium is substantially the same under both regimes. Plea bargaining is successful for a greater range of parameters under the three-outcome regime. The G-type prefers the two-outcome regime, whereas the I-type prefers the three-outcome regime. P prefers the three-outcome regime, as she obtains the same expected payoff from trial,

whereas the plea offer is higher in the three-outcome regime and is accepted with the same probability. Finally, the expected loss due to misclassification experienced by outside observers is lower under the three-outcome regime.

*Equilibrium Plea Acceptance by Innocent Defendants*

In the base model, D has two possible types: G and I. When we add the idea of a strong (S) and a weak (W) version of D, with  $\omega$  being the probability that D is weak, we end up with four possible types: GS, GW, IS, and IW. For the base model, equations (A.1a)-(A.1d) in the Appendix describe  $\Theta$ 's posterior belief that D is G, given the case disposition a, b, c, or d. These depend on  $\Theta$ 's conjectures about the probability that a D of type G (resp., I) would reject the plea offer, which is denoted by  $\rho_G^{D\Theta}$  (resp.,  $\rho_I^{D\Theta}$ ). When we expand our type set as above, we will need four type-indexed probabilities of rejection:  $\rho_{GS}^{D\Theta}$  denotes  $\Theta$ 's conjectures about the probability that a D of type GS would reject the plea offer. The expressions  $\rho_{GW}^{D\Theta}$ ,  $\rho_{IS}^{D\Theta}$ , and  $\rho_{IW}^{D\Theta}$  are similarly defined. The relevant equations, for arbitrary conjectures, are modified as follows:

$$\begin{aligned}\mu(G | a) &= [\omega\rho_{GW}^{D\Theta} + (1-\omega)\rho_{GS}^{D\Theta}](1-\lambda)F_G / \{[\omega\rho_{GW}^{D\Theta} + (1-\omega)\rho_{GS}^{D\Theta}](1-\lambda)F_G + [\omega\rho_{IW}^{D\Theta} + (1-\omega)\rho_{IS}^{D\Theta}]\lambda F_I\}; \\ \mu(G | b) &= [\omega(1-\rho_{GW}^{D\Theta}) + (1-\omega)(1-\rho_{GS}^{D\Theta})](1-\lambda) / \{[\omega(1-\rho_{GW}^{D\Theta}) + (1-\omega)(1-\rho_{GS}^{D\Theta})](1-\lambda) + [\omega(1-\rho_{IW}^{D\Theta}) + (1-\omega)(1-\rho_{IS}^{D\Theta})]\lambda\}; \\ \mu(G | c) &= [\omega\rho_{GW}^{D\Theta} + (1-\omega)\rho_{GS}^{D\Theta}](1-\lambda)(1-F_G) / \{[\omega\rho_{GW}^{D\Theta} + (1-\omega)\rho_{GS}^{D\Theta}](1-\lambda)(1-F_G) + [\omega\rho_{IW}^{D\Theta} + (1-\omega)\rho_{IS}^{D\Theta}]\lambda(1-F_I)\}; \\ \text{and } \mu(G | d) &= [\omega\rho_{GW}^{D\Theta} + (1-\omega)\rho_{GS}^{D\Theta}](1-\lambda) / \{[\omega\rho_{GW}^{D\Theta} + (1-\omega)\rho_{GS}^{D\Theta}](1-\lambda) + [\omega\rho_{IW}^{D\Theta} + (1-\omega)\rho_{IS}^{D\Theta}]\lambda\}.\end{aligned}$$

As long as the share of weak types ( $\omega$ ) is sufficiently small, the equilibrium will still involve P making a plea offer that renders the D of type GS indifferent between acceptance and going to trial; the D of type I rejects this demand for sure ( $\rho_I^{D\Theta} = 1$ ) and both weak types accept it for sure ( $\rho_{GW}^{D\Theta} = \rho_{IW}^{D\Theta} = 0$ ). Substituting these into the equations above gives the following:

$$\mu(G | a) = \rho_{GS}^{D\Theta}(1-\lambda)F_G / [\rho_{GS}^{D\Theta}(1-\lambda)F_G + \lambda F_I]; \quad (\text{TA.22a})$$

$$\mu(G | b) = [\omega + (1 - \omega)(1 - \rho_{GS}^{D\Theta})](1 - \lambda) / \{[\omega + (1 - \omega)(1 - \rho_{GS}^{D\Theta})](1 - \lambda) + \omega\lambda\}; \quad (\text{TA.22b})$$

$$\mu(G | c) = \rho_{GS}^{D\Theta}(1 - \lambda)(1 - F_G) / [\rho_{GS}^{D\Theta}(1 - \lambda)(1 - F_G) + \lambda(1 - F_I)]; \quad (\text{TA.22c})$$

and  $\mu(G | d) = \rho_{GS}^{D\Theta}(1 - \lambda) / [\rho_{GS}^{D\Theta}(1 - \lambda) + \lambda]. \quad (\text{TA.22d})$

It is straightforward to compare these equations with equations (A.1a)-(A.1d), after substituting therein the equilibrium value  $\rho_I^{D\Theta} = 1$ . In particular, the system of equations above is the same as the system (A.1a)-(A.1d), with one exception. In the base model, equation (A.1b) becomes  $\mu(G | b) = 1$ ; the acceptance of a plea offer is a clear signal of guilt. Whereas equation (TA.22b) provides a value for  $\mu(G | b)$  that is less than 1 for all  $\omega > 0$ ; that is, acceptance of a plea offer is no longer a sure sign of guilt, as a fraction  $\omega$  of innocent defendants also accept the plea offer. Moreover,  $\mu(G | b)$  is a decreasing function of  $\omega$ ; as the fraction of weak defendants increases, accepting the plea offer is an increasingly weak signal of guilt.

Now consider P's beliefs upon observing a rejection of her plea offer. In the base model, for arbitrary conjectures these beliefs are given by:  $v(G | R; \rho_G^{D\Theta}, \rho_I^{D\Theta}) = \rho_G^{D\Theta}(1 - \lambda) / [\rho_G^{D\Theta}(1 - \lambda) + \rho_I^{D\Theta}\lambda]$  (recall, P and  $\Theta$  have common conjectures about D). In equilibrium,  $\rho_I^{D\Theta} = 1$ , so  $v(G | R; \rho_G^{D\Theta}, \rho_I^{D\Theta} = 1) = \rho_G^{D\Theta}(1 - \lambda) / [\rho_G^{D\Theta}(1 - \lambda) + \lambda]$ . With four D types, this expression becomes:  $v(G | R; \rho_{GS}^{D\Theta}, \rho_{GW}^{D\Theta}, \rho_{IS}^{D\Theta}, \rho_{IW}^{D\Theta})$

$$= [\omega\rho_{GW}^{D\Theta} + (1 - \omega)\rho_{GS}^{D\Theta}](1 - \lambda) / \{[\omega\rho_{GW}^{D\Theta} + (1 - \omega)\rho_{GS}^{D\Theta}](1 - \lambda) + [\omega\rho_{IW}^{D\Theta} + (1 - \omega)\rho_{IS}^{D\Theta}]\lambda\}.$$

In equilibrium,  $\rho_{IS}^{D\Theta} = 1$  and  $\rho_{GW}^{D\Theta} = \rho_{IW}^{D\Theta} = 0$ ; making these substitutions results in:

$v(G | R; \rho_{GS}^{D\Theta}, \rho_{GW}^{D\Theta} = 0, \rho_{IS}^{D\Theta} = 1, \rho_{IW}^{D\Theta} = 0) = \rho_{GS}^{D\Theta}(1 - \lambda) / [\rho_{GS}^{D\Theta}(1 - \lambda) + \lambda]$ . This is exactly the same form as in the base model. Given that the fraction  $\omega$  of both G-types and I-types accept the plea offer for sure, and all of the remaining (i.e., strong) I-types reject the plea offer, the mixture of innocent and guilty defendants among those that rejected the plea offer has exactly the same form.

This allows us to write P's indifference condition between taking a case to trial versus dropping it following rejection as follows (this is simply equation (TA.11) with  $\rho_{GS}^{D\theta}$  in place of  $\rho_G^{D\theta}$ ):

$$\begin{aligned} & v(G | R; \rho_{GS}^{D\theta}) \{S_c(1 - F_G) - k^P - r_1^P \mu(I | c; \rho_{GS}^{D\theta})(1 - F_G) - r_G^P \mu(G | a; \rho_{GS}^{D\theta}) F_G\} \\ & + v(I | R; \rho_{GS}^{D\theta}) \{S_c(1 - F_I) - k^P - r_1^P \mu(I | c; \rho_{GS}^{D\theta})(1 - F_I) - r_G^P \mu(G | a; \rho_{GS}^{D\theta}) F_I\} \\ & = - r_G^P \mu(G | d; \rho_{GS}^{D\theta}). \end{aligned} \quad (TA.23)$$

Since we have already verified that all of the expressions above are the same as in the base model, we can conclude that, in equilibrium, the D of type GS will mix between accepting and rejecting the plea offer, rejecting it with exactly the same probability as before. That is, the D of type GS rejects the plea offer with probability  $\rho_{GS}^{D\theta} = \rho_G^{D0}$ ; the computed value of  $\rho_G^{D0}$  is given in equation (TA.18).

The equilibrium rate of plea acceptance is now  $\omega\lambda + [\omega + (1 - \omega)(1 - \rho_G^{D0})](1 - \lambda)$ , which is higher than in the base model wherein this rate is  $(1 - \rho_G^{D0})(1 - \lambda)$ . P is able to obtain a plea agreement with more guilty defendants, but also unavoidably sweeps up some innocent defendants as well.

The equilibrium plea offer is also affected because a defendant accepting a plea offer is no longer inferred to be guilty for sure. The plea offer in the base model is  $S_b(\rho_G^{D0}) = \pi_T^D(G; \rho_G^{D0}) - r^D$ , whereas the new plea offer is:  $S_b(\rho_G^{D0}) = \pi_T^D(GS; \rho_G^{D0}) - r^D \mu(G | b; \rho_G^{D0})$ . Note that  $\pi_T^D(G; \rho_G^{D0})$  and  $\pi_T^D(GS; \rho_G^{D0})$  are the same function, as all guilty defendants are type GS in the base model. So the plea offer is higher in the model with weak types.