

# When is Reputation Bad?<sup>1</sup>

Jeffrey Ely

Drew Fudenberg

David K. Levine<sup>2</sup>

First Version: April 22, 2002

This Version: October 23, 2002

**Abstract:** In traditional reputation theory, reputation is good for the long-run player. In "Bad Reputation," Ely and Valimaki give an example in which reputation is unambiguously bad. This paper characterizes a more general class of games in which that insight holds, and presents some examples to illustrate when the bad reputation effect does and does not play a role. The key properties are that participation is optional for the short-run players, and that every action of the long-run player that makes the short-run players want to participate has a chance of being interpreted as a signal that the long-run player is "bad." We also broaden the set of commitment types, allowing many types, including the "Stackelberg type" used to prove positive results on reputation. Although reputation need not be bad if the probability of the Stackelberg type is too high, the relative probability of the Stackelberg type can be high when all commitment types are unlikely.

---

<sup>1</sup> We are grateful to Adam Szeidl for careful proofreading, to Juuso Valimaki for helpful conversations, and to National Science Foundation Grants SES-9730181, SES99-86170, SES99-85462, and SES-0112018

<sup>2</sup> Departments of Economics, Northwestern University, Harvard University and UCLA.

## 1. Introduction

A long-run player playing against a sequence of short-lived opponents can build a reputation for playing in a specific way and so obtain the benefits of commitment power. To model these “reputation effects,” the literature following Kreps and Wilson [1982] and Milgrom and Roberts [1982] has supposed that there is positive prior probability that the long-run player is a “commitment type” who always plays a specific strategy.<sup>3</sup> In “Bad Reputation,” Ely and Valimaki [2001] (henceforth EV) construct an example in which introducing a particular commitment type hurts the long-run player. When the game is played only once and there are no commitment types, the unique sequential equilibrium is good for the long-run player. This remains an equilibrium when the game is repeated without commitment types, regardless of the player’s discount factor. However, when a particular “bad” commitment type is introduced, the only Nash equilibria are “bad” for a patient long-run player.<sup>4</sup>

What is not clear from EV is when reputation is bad. This paper extends the ideas in EV to a more general class of games in an effort to find the demarcation between “bad” and “good” reputation. In addition, we try to relate the EV conclusions to past work on reputation effects.

Reputation effects are most powerful when the long-run player is very patient, and Fudenberg and Levine [1992] (FL) provided upper and lower bounds on the limiting values of the equilibrium payoff of the long-run player as that player’s discount factor tends to 1. The upper bound

---

<sup>3</sup> See Sorin [1999] for a recent survey of the reputation effects literature, and its relationship to the literature on merging of opinions.

<sup>4</sup> It is obvious that incomplete information about the long-run player’s type can be harmful when the long-run player is impatient, since incomplete information can be harmful in one-shot games. Fudenberg-Kreps [1987] argue that a better measure of the “power of reputation effects” is to hold fixed the prior distribution over the reputation-builder’s types, and compare the reputation-building scenario to one in which the reputation builder’s opponents do not observe how the reputation builder has played against other opponents. They discuss why reputation effects might be detrimental in the somewhat different setting of a large long-run player facing many simultaneous small but long run opponents.

corresponds to the usual notion of the “Stackelberg payoff.” The lower bound, called the “generalized Stackelberg payoff,” weakens this notion to allow the short-run players to have incorrect beliefs about the long-run player’s strategy, so long as the beliefs are not disconfirmed by the information that the short-run players get to observe. When the stage game is a one-shot simultaneous-move game, actions are observed, payoffs are generic, and commitment types have full support, these two bounds coincide, so that the limit of the Nash equilibrium payoffs as the long-run player’s discount factor tends to one is the single point corresponding to the Stackelberg payoff. For extensive-move stage games, with public outcomes corresponding to terminal nodes, the bounds can differ. However, although FL provided examples in which the lower bound is attained, in those examples the upper bound was attained as well, and we are not aware of past work that determines the range of possible limiting values for a fairly general class of games.

Here we examine the upper bound more closely for a specific class of games designed to capture the insight of EV. Specifically, we define a class of “bad reputation” games, in which the long-run player can do no better than if the short-run players choose not to participate. This extends the EV example in a number of ways. We allow a broad class of stage games in which participation by the short-run players is optional; allowing for many actions, many signals, many short-run players, and a wide variety of payoffs. Especially important, we allow for a broad range of types, including types that are committed to “good” actions, as well as types that are committed to “bad” actions. Earlier research suggests that to attain the upper bound on the long-run player’s payoff, it can be important to include the “Stackelberg type” that is committed to the stage-game action the long-run player would choose in a Stackelberg equilibrium.<sup>5</sup> We

---

<sup>5</sup> EV consider two specifications for the bad type, either “committed” (to playing the bad action) or “strategic” (willing to play a different action occasionally to increase entry and the future payoff from playing “bad.”) In a related model, Mailath and Samuelson [1998] argue that “bad” types – and specifically strategic bad types – are more plausible than Stackelberg types. We are sympathetic to the argument that strategic bad types may be

find that the EV result fails if the probability of this Stackelberg type is too high, but extends to the case where the probability of the Stackelberg type is sufficiently low, but nonzero. This shows that it is not essential to rule out the types that support “good” reputation effects in order to derive the bad reputation result. Moreover, the *relative* probability of the Stackelberg type can be high when all commitment types are unlikely; in this sense our conditions hold for “almost all” sufficiently small commitment-type perturbations of the complete information limit.

By extending the EV example to a broad class of stage games we are able to more clearly identify the types of assumptions key to a bad reputation. There are several such properties, notably that the short-run players can either individually or collectively choose not to participate. However, most of the assumptions on the structure of the game seem to involve little loss of economic applicability. The key substantive assumption seems to be that every action of the long-run player that makes the short-run players want to participate in the game has a chance of being misconstrued as a signal of a “bad reputation.”

EV motivate their example by considering an automobile mechanic who has specialized knowledge of the work that needs to be done to repair the car. We think that we have identified a broader class of bad reputation games that can be interpreted as “expert advice.” This includes consulting a doctor or stockbroker, or in the macroeconomics context, can be the decision whether or not to turn to the IMF for assistance. In EV, the short-run players observe only the advice, but not the consequences of the advice. Here we explicitly consider what happens when the short-run players observe the consequences as well. We also show that there are other distinct classes of games with rather different observation structures that are bad reputation games, such as our “teaching evaluation” game, where “advice” is not an issue because the long-run player does not

---

more likely than commitment types, but this does not imply that the probability of commitment types should be zero. Instead, we would argue that it is preferable for models to allow for a wide range of types, especially those with fairly simple behavior rules.

privately observe anything that is payoff-relevant for the short-run player. Finally, we illustrate the boundaries of bad reputation by giving a number of examples and classes of participation games that are not bad reputation games.

## 2. The Model

### 2.1. The Dynamic Game

There are  $N + 1$  players, a long run-player 1, and  $N$  short-run players  $2 \dots N + 1$ . The game begins at time  $t = 1$  and is infinitely repeated. Each period, each player  $i$  chooses from a finite action space  $A^i$ . We denote individual actions  $a^i$ , and action profiles by  $\mathbf{a}$ . We also use  $\mathbf{a}^{-i}$  to denote the play of all players except player  $i$  and  $\mathbf{a}^{-i-j}$  to denote the play of all players except players  $i$  and  $j$ .

The long-run player discounts the future with discount factor  $\mathbf{d}$ . Each short-run player plays only in one period, and is replaced by an identical short-run player in the next period. There is a set  $\Theta$  of types of long-run player. There are two sorts of types: type  $\mathbf{0} \in \Theta$  is called the “rational type,” and is the focus of our interest, with utility described below. For each pure action  $a^1$ , type  $\mathbf{q}(a^1)$  is a “committed type,” that is constrained to play  $a^1$ . These are the only possible types in  $\Theta$ . Note that we do not require that every pure action commitment type has positive probability. The stage game utility functions are  $u^i(\mathbf{a})$ , where  $u^1(\mathbf{a})$  corresponds to the long-run player of type  $\mathbf{q} = \mathbf{0}$ . The common prior distribution over long-run player types is denoted  $\mathbf{m}(\mathbf{0})$ .

There is a finite public signal space  $Y$  with signal probabilities  $\mathbf{r}(y | \mathbf{a})$ . All players observe the history of the public signals. Short-run players observe only the history of the public signals, and in particular observe neither the past actions of the long-run player, nor of previous short-run players. We do not assume that the payoffs depend on the actions only through the signals, so the short-run players at date  $t$  need not

know the realized payoffs of the previous generations of short-run players.<sup>6</sup>

We let  $h_t = (y_1, y_2, \dots, y_t)$  denote the public history through the end of period  $t$ . We denote the null history by  $\mathbf{0}$ . We let  $h_t^1$  denote the private history known only to the long-run player. This includes his own actions, and may or may not include the actions of the short-run players he has faced in the past. A strategy for the long-run player is a sequence of maps  $\mathbf{s}^1(h_t, h_t^1, \mathbf{q}) \in \text{conhull } A^1 \equiv \mathcal{A}^1$ ; a strategy profile for the short-run players is a sequence of maps  $\mathbf{s}^j(h_t) \in \text{conhull } A^j \equiv \mathcal{A}^j$ . (Note that  $\mathcal{A}^{-1}$  denotes the product of the  $\mathcal{A}^j$ 's, not the convex hull of the product.) A short-run profile  $\mathbf{a}^{-1}$  is a Nash response to  $\mathbf{a}^1$  if  $u^i(\mathbf{a}^1, \mathbf{a}^i, \mathbf{a}^{-1-i}) \geq u^i(\mathbf{a}^1, \mathbf{a}^i, \mathbf{a}^{-1-i})$  for all  $\mathbf{a}^i \in A^i$ . We denote the set of short-run Nash responses to  $\mathbf{a}^1$  by  $B(\mathbf{a}^1)$ .

Given strategy profiles  $\mathbf{s}$ , the prior distribution over types  $\mathbf{m}(\mathbf{0})$  and a public history  $h_t$  that has positive probability under  $\mathbf{s}$ , we can calculate from  $\mathbf{s}^1$  the conditional probability of long-run player actions  $\bar{\mathbf{a}}^1(h_t)$  given the public history. A *Nash Equilibrium* is a strategy profile  $\mathbf{s}$  such that for each positive probability history

- 1)  $\mathbf{s}^{-1}(h_t) \in B(\bar{\mathbf{a}}^1(h_t))$  [short-run players optimize]
- 2)  $\mathbf{s}^1(h_t, h_t^1, \mathbf{q}(\mathbf{a}^1)) = \mathbf{a}^1$  [committed types play accordingly]
- 3)  $\mathbf{s}^1(\cdot, \cdot, \mathbf{0})$  is a best-response to  $\mathbf{s}^{-1}$  [rational type optimizes].

## 2.2 The Ely-Valimaki Example

We will use the EV example to illustrate our assumptions and definitions. In EV, the long-run player is a mechanic, her action is a map from the privately observed state of the customer's car  $\mathbf{w} \in \{E, T\}$  to

---

<sup>6</sup> Fudenberg and Levine [1992] assumed that a player's payoff was determined by his own action and the realized signal, but that assumption was not used in the analysis. The assumption is used in models with more than one long-run player to justify the restriction to public equilibria, but it is not needed here.

announcements  $\{e, t\}$ , where  $E$  means the car needs a new engine,  $T$  means it needs at tune-up, and the announcements, which are what the mechanic says the car needs, determine what is actually done to the car. Thus the long-run player's action space is a map from privately observed states to announcements,  $A^1 = \{ee, et, te, tt\}$ , where the first component indicates the announcement in response to the signal  $E$  and the second to  $T$ . There is one short-run player each period who chooses an element of  $A^2 = \{In, Out\}$ . The public signal takes on the values  $Y = \{e, t, Out\}$ . If the short-run player chooses  $Out$  the signal is  $Out$ , that is  $r(Out | a^1, Out) = 1$ ; otherwise the signal is the announcement of the long-run player. The two states of the car are assumed to be i.i.d. and equally likely, so  $r(e | (et, In)) = r(e | (te, In)) = 1/2$ ,  $r(e | (ee, In)) = 1$ , and  $r(e | (tt, In)) = 0$ . If the short-run player chooses  $Out$ , each player gets utility 0. If he plays  $In$  and the long-run player's announcement is truthful (that is, matches the state), the short-run player receives  $u$ ; if it is untruthful, it is  $-w$  where  $w > u > 0$ . The "rational type" of long-run player has exactly the same stage-game payoff function as the short run players. Thus when the long-run player is certain to be the rational type, the strategic form of the stage game is

	<i>In</i>	<i>Out</i>
<i>ee</i>	$(u - w) / 2, (u - w) / 2$	0, 0
<i>et</i>	$u, u$	0, 0
<i>te</i>	$-w, -w$	0, 0
<i>tt</i>	$(u - w) / 2, (u - w) / 2$	0, 0

Figure 1

When the rational type is the only type in the model, there is an equilibrium where he chooses the action that matches the state, all short-run players enter, and the rational type's payoff is  $u$ . However, EV show that when there is also a probability that the long-run player is a "bad

$ee$ , the long-amount that converges to 0 as the discount factor goes to 1. The intuition for this result has four steps. First of all, the rational type must  $et$  if its equilibrium  $p$  is 0, because the short-run player is too  $ee$ . Second, from Bayes' rule it follows that there is some number  $K$  such that  $E$  in periods where the rational type is playing honestly will make the posterior probability of the bad type so high that all subsequent short-run players play out. Third, when there have been  $K - 1$  successive observations of  $E$ , the rational type of long run player is tempted to play  $tt$  instead of  $et$ , even though this lowers his short-run payoff, to avoid driving out the short-run players with another observation of  $E$ . Thus, the long-run player is tempted to take an action that is worse for both himself and the short-run players in order to avoid being incorrectly tagged as a "bad type;" an induction argument shows that honest play by the rational type unravels.

### **2.3. Participation Games and Bad Reputation Games**

We consider "participation games" in which the short-run players may choose not to participate. The crucial aspect of non-participation by the short-run players is that it conceals the action taken by the long-run player from subsequent short-run players; this is what allows the lower bound on the long-run player's Nash equilibrium payoff in the EV example to be lower than Stackelberg payoff. We will then define "bad reputation" games as a subclass of participation games that have the additional features needed for the bad reputation result; the following is a brief summary of the key features of these games. First of all, there will be a set of "friendly" actions that must receive sufficiently high probability to induce the short-run players to participate, such as  $et$  in the EV example. Next, there are "bad signals" that are most strongly linked to some "unfriendly" actions that deter participation, but which also have



positive probability under friendly play; in EV the bad signal is  $E$ . Finally, there are some actions that are not friendly, but reduce the probability of the bad signals, such as  $tt$  in EV; we call these actions “temptations.” If there is a positive prior probability that the long-run player is a “bad type” that is committed to one of the unfriendly actions, some histories of play will induce the short-run players to exit, so to avoid these histories the rational type of long-run player may choose to play one of the temptations; foreseeing this, the short-run players will chose not to enter. Our main result shows that this leads to a “bad reputation” result as long as the prior does not assign too much probability to types that are committed to play friendly actions.

To model the option to not participate, we assume that certain public signals  $y^e \in Y^e$  are *exit signals*. Associated with these exit signals are *exit profiles*, which are pure action profiles  $e^{-1} \in E^{-1} \subseteq A^{-1}$  for the short run players.

For each such  $e$ ,  $\mathbf{r}(y^e | a^1, e^{-1}) = \mathbf{r}(y^e | e^{-1})$  for all  $a^1$ , and  $\mathbf{r}(Y^e | e^{-1}) = 1$ . In other words, if an exit profile is chosen, an exit signal must occur, and the distribution of exit signals is independent of the long-run player's action. Moreover, if  $a^{-1} \notin E^{-1}$  then  $\mathbf{r}(y^e | a^1, a^{-1}) = 0$  for all  $a^1 \in A^1, y^e \in Y^e$ . We refer to  $A^{-1} - E^{-1}$  as the entry profiles. Note that an entry profile cannot give rise to an exit signal. A *participation game* is a game in which  $E^{-1} \neq \emptyset$ . The remainder of the paper specializes to participation games.

We begin by distinguishing actions by the long-run player that cause the short-run players to exit (unfriendly actions), and those that are needed to get them to enter (friendly actions).

**Definition 1:** A non-empty finite set of pure actions for the long-run player  $N^1$  is unfriendly if there is a number  $\mathbf{y} < 1$  such that  $\mathbf{a}^1(N^1) \geq \mathbf{y}$  implies  $B(\mathbf{a}^1) \subseteq \text{conhull } E^{-1}$ .

*Remark:* This definition says that unfriendly actions induce exit, in the strong sense that exit is the only best response if the probability of the

unfriendly actions is sufficiently high. There will often be many sets of unfriendly actions. In the EV example the set  $\{ee, tt, te\}$  is unfriendly, and so is any subset.

**Definition 2:** A non-empty finite set of mixed actions  $F^1$  for the long run player is friendly if there is a number  $\mathbf{g} > \mathbf{0}$  such that  $B(\mathbf{a}^1) \cap [\mathcal{A}^{-1} - \text{conhull}(E^{-1})] \neq \emptyset$  implies  $\mathbf{a}^1 \geq \mathbf{g}f^1$  for some  $f^1 \in F^1$ . The number  $\mathbf{g}$  is called the size of the friendly set.

*Remark:* This definition says that the probabilities given to every pure action must be bounded below by a scale factor times some friendly mixture if the short-run players are not to exit. Note that weight on a friendly action is necessary for entry, but need not be sufficient for entry, and that a friendly set must be non-empty. There may also be many different friendly sets. Suppose that  $F^1$  is friendly of size  $\hat{\mathbf{g}}$ , and let  $\underline{F}^1 \equiv \min\{f(\mathbf{a}) > 0 \mid f \in F^1\}$ . Then if  $f^1 \in F^1$  it may be replaced by any mixture over the support of  $f^1$ , and the resulting set will be friendly of size  $\hat{\mathbf{g}}\underline{F}_1$ . Similarly, if we have a friendly set and we eliminate mixtures  $\tilde{f}^1 \in F^1$  whose support contained in the support of some different  $f^1 \in F^1$ , we get a new friendly set with a smaller value of  $\underline{\mathbf{a}}$ . In the EV example, the action  $et$  is friendly, with

$$\underline{\mathbf{a}} = \frac{w - u}{w + u/2}.$$

Finally, consider a pure action  $\mathbf{a}^1$  such that  $B(\mathbf{a}^1) \cap [\mathcal{A}^{-1} - \text{conhull}(E^{-1})] \neq \emptyset$ . Since  $\mathbf{a}^1$  is pure,  $\mathbf{a}^1 \geq \mathbf{g}f^1$  is possible only if  $f^1 = \mathbf{a}^1$ . In other words, any pure action that permits short-run entry (such as  $et$  in the EV example) *must* be in *every* friendly set.

**Definition 3:** The support  $A^1(F^1)$  of a friendly set  $F^1$  is the actions that are played with positive probability:

$$A^1(F^1) \equiv \{\mathbf{a}^1 \in A^1 \mid f^1(\mathbf{a}^1) > 0, f^1 \in F^1\}$$

We say that a friendly set  $F^1$  is *orthogonal* to an unfriendly set  $N^1$  if  $N^1 \cap A^1(F^1) = \emptyset$ .

Next we consider what signals may reveal about actions.

**Definition 4:** We say that a set of signals  $\hat{Y}$  is unambiguous for a set of actions  $N^1$  if for all  $a^{-1} \notin E^{-1}, \hat{y} \in \hat{Y}, n^1 \in N^1, a^1 \notin N^1$  we have  $r(\hat{y} | n^1, a^{-1}) > r(\hat{y} | a^1, a^{-1})$ .

This is a strong condition: every action in  $N^1$  must assign a higher probability to each signal in  $\hat{Y}$  than any action not in  $N^1$ . A given set of actions may not have signals that are unambiguous; in the case of the EV example,  $E$  is an unambiguous signal for the unfriendly set  $\{ee\}$ .

**Definition 5:** An action  $a^1$  is vulnerable to temptation relative to a set of signals  $\hat{Y}$  if there exist numbers  $\underline{r}, \tilde{r} > 0$  and an action  $b^1$  such that

- 1) If  $a^{-1} \notin E^{-1}, \hat{y} \in \hat{Y}$ , then  $r(\hat{y} | b^1, a^{-1}) \leq r(\hat{y} | a^1, a^{-1}) - \underline{r}$ .
- 2) If  $a^{-1} \notin E^{-1}$  and  $y \notin \hat{Y} \cup Y^e$  then  $r(y | b^1, a^{-1}) \geq (1 + \tilde{r})r(y | a^1, a^{-1})$ .
- 3) For all  $e^{-1} \in E^{-1}$ ,  $u^1(b^1, e^{-1}) \geq u^1(a^1, e^{-1})$ .

The action  $b^1$  is called a *temptation*, and the parameters  $\underline{r}, \tilde{r}$  are the *temptation bounds*.

In other words, an action is vulnerable if it is possible to lower the probability of all of the signals in  $\hat{Y}$  by at least  $\underline{r}$  while increasing the probability of each other signal by at least the multiple  $(1 + \tilde{r})$ . Notice that for an action to be vulnerable to a temptation, it must place at least weight  $\underline{r}$  on each signal in  $\hat{Y}$ . Notice also that the definition does not control the payoff to the vulnerable action when the short-run players participate – the temptation here is not to increase short-run payoff, but rather to decrease the probability of the signals in  $\hat{Y}$ . In the EV example, the action  $et$  is vulnerable relative to  $\{E\}$ . The temptation  $b^1$  is  $tt$ , which sends the

probability of the signal  $E$  to zero. (Since there is one other signal, condition 2 of the definition is immediate.)

Notice that if an action  $a^1$  is vulnerable, it cannot be the case that if  $a^{-1} \notin \text{conhull} E^{-1}$  then  $r(\cdot | a^1, a^{-1}) = r(\cdot | a^{-1})$ ; the distribution of signals must be in some way dependent on the long-run player's action if the short-run players do not exit. This is related to the notion of an action being identified, as in Fudenberg, Levine and Maskin [1994]. Here we allow the possibility that there are strategies such as  $et$  and  $te$  from the EV example that are not identified, but do not allow complete lack of identification unless the short-run players play in  $E^{-1}$  with probability one.

**Definition 6:** A mixed action  $a^1$  for the long run player is enforceable if there does not exist another action  $\tilde{a}^1$  such that for all  $a^{-1} \in E^{-1}$ ,  $u^1(\tilde{a}^1, a^{-1}) \geq u^1(a^1, a^{-1})$  and for all  $a^{-1} \in A^{-1} - E^{-1}$ ,  $u^1(\tilde{a}^1, a^{-1}) > u^1(a^1, a^{-1})$  and  $r(\cdot | \tilde{a}^1, a^{-1}) = r(\cdot | a^1, a^{-1})$ . When  $a^1$  is not enforceable, we say that the action  $\tilde{a}^1$  defeats  $a^1$ .

If an action is not enforceable then there is necessarily lack of identification, since  $a^1$  and  $\tilde{a}^1$  induce exactly the same distribution over signals. The key point is that if the short-run players enter with positive probability, the rational type cannot play an action that is not enforceable: by switching to  $\tilde{a}^1$  he would strictly increase his current payoff, while maintaining the same distribution over signals, and so the same future utility. Note also that a mixed action that assigns positive probability to unenforceable actions is not enforceable: if  $a^1$  assigns probability  $p$  to unenforceable action  $a^1$ , then  $a^1$  is defeated by the mixed action  $\hat{a}^1$  formed by replacing the probability on  $a^1$  with the action  $\tilde{a}^1$  that defeats  $a^1$ .

**Definition 7:** A participation game has an exit minmax if

$$\begin{aligned} & \max_{\mathbf{a}^{-1} \in E^{-1} \cap \text{range}(B)} \max_{\mathbf{a}^1} u^1(\mathbf{a}^1, \mathbf{a}^{-1}) = \\ & \min_{\mathbf{a}^{-1} \in \text{range}(B)} \max_{\mathbf{a}^1} u^1(\mathbf{a}^1, \mathbf{a}^{-1}) \end{aligned}$$

In other words, any exit strategy forces the long-run player to the minmax payoff, where the relevant notion of minmax incorporates the restriction that the action profile chosen by the short-run players must lie in the range of  $B$ .<sup>7</sup> It is convenient in this case to normalize the minmax payoff to 0. We are now in a position to define a class of games we call *bad reputation games*.

**Definition 8:** *A participation game is a bad reputation game if it has an exit minmax, there is an unfriendly set  $N^1$ , a friendly set  $F^1$  that is orthogonal to  $N^1$ , and a set of signals  $\hat{Y}$  that are unambiguous for  $N^1$ , and such that every enforceable  $f^1 \in F^1$  is vulnerable to temptation relative to  $\hat{Y}$ . The signals  $\hat{Y}$  are called the bad signals.*

In particular, the EV game is a bad reputation game. We take the friendly set to be  $\{et\}$ , the unfriendly set to be  $\{ee\}$  and the unfriendly signals to be  $\{E\}$ . We have already observed that  $\{et\}$  is a friendly set and  $\{ee\}$  unfriendly. The two are obviously orthogonal, and  $\{E\}$  is unambiguous for  $\{ee\}$ .

In a bad reputation game, the relevant temptations are those relative to  $\hat{Y}$ . For the remainder of the paper when we examine a bad reputation game and refer to a temptation, we will always mean relative to the set  $\hat{Y}$ .

---

<sup>7</sup> When there is a single short-run player this restriction collapses to the constraint of not playing strictly dominated strategies, but when there are multiple short-run players it involves additional restrictions. It is clear that no equilibrium could give the long-run player a lower payoff than the minmax level defined in definition 7. Conversely, in complete-information games, any long-run player payoff above this level can be supported by a perfect Bayesian equilibrium if actions are identified and the public observations have a “product structure” (Fudenberg and Levine [1994]). This is true in particular when actions are publicly observed as shown in Fudenberg, Kreps and Maskin [1990].

For any bad reputation game, it is useful to define several constants describing the game. Recall that  $\mathbf{y}$  is the probability in the definition of an unfriendly set and that  $\mathbf{g}$  is the scale factor in the definition of a friendly set. Since the friendly set is finite, we may define  $\mathbf{j} > \mathbf{0}$  to be the minimum, taken over elements of the friendly set, of the values  $\underline{\mathbf{r}}$  in the definition of temptation. Define

$$r = \min_{n^1 \in N^1, a^1 \notin N^1, a^{-1} \notin \text{conhull}(E^{-1}), \hat{\mathbf{y}} \in \hat{Y}} \frac{r(\hat{\mathbf{y}} \mid n^1, \mathbf{a}^{-1})}{r(\hat{\mathbf{y}} \mid a^1, \mathbf{a}^{-1})}.$$

Since the friendly set is non-empty and orthogonal to the unfriendly set, the denominator of this expression is well defined, and since  $\hat{Y}$  is unambiguous for the unfriendly set,  $r > 1$ .

Also define

$$\mathbf{h} = -\log(\mathbf{g}\mathbf{j}) / \log r$$

which is positive, and

$$k_0 = -\frac{\log(\mathbf{y})}{\log\left(\mathbf{y} + (1 - \mathbf{y})\frac{1}{r}\right)}.$$

### 3. The Theorem

We now prove our main result: In a bad reputation game with a sufficiently patient long-run player and likely enough unfriendly types, in any Nash equilibrium, the long-run player gets approximately the exit payoff. The proof uses several Lemmas proven in the Appendix.

We begin by describing what it means for unfriendly types to be likely “enough.” Let  $\Theta(F^1)$  be the commitment types corresponding to actions in the support of  $F^1$ ; we will call these the *friendly commitment types*. Let  $\hat{\Theta}$  be the *unfriendly commitment types* corresponding to the unfriendly set  $N^1$ .

**Definition 9:** A bad reputation game has commitment size  $\mathbf{e}, \mathbf{f}$  if

$$\mathbf{m}(\mathbf{0})[\Theta(F^1)] \leq \mathbf{e} \left( \frac{\mathbf{m}(\mathbf{0})[\hat{\Theta}]}{\mathbf{m}(\mathbf{0})[\Theta(F^1)]} \right)^{\mathbf{f}}$$

where  $\mathbf{f} > 0$ .

This notion of commitment size places a bound on the prior probability of friendly commitment types that depends on the prior probability of the unfriendly types. Since  $\mathbf{f}$  is positive, the larger the prior probability of  $\hat{\Theta}$ , the larger the probability of the friendly commitment types is allowed to be. The hypothesis that the priors have commitment size  $\mathbf{e}, \mathbf{f}$  for sufficiently small  $\mathbf{e}$  is a key assumption driving our main results.

Note that the assumption of a given commitment size does not place any restrictions on the relative probabilities of commitment types. In particular, let  $\tilde{\mathbf{m}}$  be a fixed prior distribution over the commitment types, and consider priors of the form  $\mathbf{I} \tilde{\mathbf{m}}$ , where the remaining probability is assigned to the rational type. Then the right-hand side of the inequality defining commitment size depends only on  $\tilde{\mathbf{m}}$ , and not on  $\mathbf{I}$ , while the left-hand side has the form  $\mathbf{I} \tilde{\mathbf{m}}$ . Hence for sufficiently small  $\mathbf{I}$  the assumption of commitment size  $\mathbf{e}, \mathbf{h}$  is satisfied. Note that the EV example has commitment size  $\mathbf{0}, \mathbf{h}$  since the only types are the rational type and the commitment type who plays  $ee$ .

Define  $U^1 = \max_a u^1(a) - \min\{0, u^1\}$ , and let  $\tilde{\mathbf{r}}_{\min} = \min_{f^1 \in F^1} \tilde{\mathbf{r}}$ . Recall that  $\underline{F}^1 \equiv \min\{f(a) > 0 \mid f \in F^1\}$  and that  $\mathbf{g}$  is the scale factor in the definition of a friendly set.

**Theorem 1:** *In a bad reputation game of commitment size  $((\mathbf{g}\underline{F}^1/2)^{(1+\mathbf{h})}, \mathbf{h})$  let  $\bar{\mathbf{v}}^1$  be the supremum of the payoff of the rational type in any Nash equilibrium. Then*

$$\bar{v}^1 \leq (1 - d)k^* \left( \frac{1}{\tilde{r}_{\min}} \right)^{k^*} \left( 1 + \frac{1}{\tilde{r}_{\min}} \right) U^1,$$

where  $k^* = k_0 + \log(\mathbf{m}(\mathbf{0})[\hat{\Theta}]) / \log\left(\mathbf{y} + (1 - \mathbf{y})\frac{1}{r}\right)$ . In particular,  $\lim_{d \rightarrow 1} \bar{v}^1 \leq 0$ .

To prove this we use a series of Lemmas proven in the Appendix. For the rest of this section, we fix an arbitrary Nash equilibrium. Given this equilibrium, let  $v^1(h_t)$  denote the expected continuation value to the rational long-run player, and let  $\mathbf{m}(h_t)$ ,  $\mathbf{a}^{-1}(h_t)$  be the posterior beliefs and strategy of the short-run players at history  $h_t$ . Notice that the expected present value to the rational long-run player conditional on a positive probability public/private history pair must not depend on the private history  $h_t^1$ , or the rational long-run type would be failing to optimize. If  $\mathbf{a}^1$  has positive probability under  $\bar{\mathbf{a}}^1(h_t)$ , and  $\mathbf{a}^{-1}$  positive probability under  $\mathbf{a}^{-1}(h_t)$ , then we define

$$v^1(h_t, \mathbf{a}) \equiv (1 - d)u^1(\mathbf{a}) + d \sum_y \mathbf{r}(y | \mathbf{a})v^1(h_t, y).$$

When mixed actions  $\mathbf{a}^1$  and  $\mathbf{a}^{-1}$  put weight only on such positive-probability  $\mathbf{a}^1, \mathbf{a}^{-1}$ , it is convenient also to define  $v^1(h_t, \mathbf{a})$  in the natural way.

**Lemma 1:** *In a participation game, if  $h_t$  is a positive probability history in which  $\hat{y} \in \hat{Y}$  occurs in period  $t$  and  $\mathbf{m}(h_{t-1})[\Theta(F^1)] \leq \mathbf{g}\underline{F}^1/2$  then  $\mathbf{a}_0^1(h_t) \geq (\mathbf{g}/2)f^1$  for some friendly  $f^1$ .*

In other words, when the prior on friendly types is sufficiently low, entry can occur only if the rational type is playing a friendly strategy with appreciable probability. This is a consequence of the definition of friendly strategies: entry requires that the overall strategy assigns some minimum probability to a friendly strategy, and if the friendly types are unlikely,



then a non-negligible part of this probability must come from the play of the rational type.

Recall that  $\mathbf{h} = -\log(\mathbf{g}\mathbf{j}) / \log r$ .

**Lemma 2:** *In a bad reputation game, if  $\mathbf{h}_t$  is a positive probability history, and the signals in  $\mathbf{h}_t$  all lie in  $Y^e \cup \widehat{Y}$ , then*

a) *At most*

$$\mathbf{k} = \mathbf{k}_0 + \frac{\log(\mathbf{m}(0)[\widehat{\Theta}])}{\log\left(\mathbf{y} + (1 - \mathbf{y})\frac{1}{r}\right)}$$

*of the signals are in  $\widehat{Y}$ .*

b) *If the commitment size is  $((\mathbf{g}/2)\underline{F}^1)^{1+\mathbf{h}}$ ,  $\mathbf{h}$  then  $\mathbf{m}(\mathbf{h}_t)[\Theta(F^1)] \leq (\mathbf{g}/2)\underline{F}^1$ .*

*Remark:* The intuition for part *a* is simple, and closely related to the argument about the deterministic evolution of beliefs in FL: The short-run players exit if they think it is likely that entry will lead to the observation of a bad signal. Hence each observation of a bad signal is a “surprise” that increases the posterior probability of the bad type by (at least) a fixed ratio greater than 1, so along a history that consists of only bad signals and exit signals, the posterior probability of the bad type eventually gets high enough that all subsequent short-run players exit. This argument holds no matter what other types have positive probability, and it is the only part of this lemma that would be needed when there are only two types, one rational and one bad, as in EV.

However, as we will show by example below, we cannot expect the “bad reputation” result to hold when there is sufficiently high probability of the Stackelberg type. Part *b* of the lemma says that if the

initial probability of the friendly types is sufficiently low compared to the prior probability of the bad types, then along any history on the path of play which consists only of exit outcomes and bad outcomes, the probability of the Stackelberg type remains low. The intuition for this result is that because of the assumption that the unfriendly and friendly sets are orthogonal,  $r > 1$ , so each observation of a bad signal not only increases the probability of the bad type, it increases the relative probability of this type compared to any friendly commitment type, and this bounds the rate of growth of the probability of friendly types.<sup>8</sup>

Define

$$\bar{u}^1(y, \tilde{\mathbf{r}}) = \begin{cases} \left(1 + \frac{1}{\tilde{\mathbf{r}}}\right)U^1 & y \in \hat{Y} \\ 0 & \text{otherwise} \end{cases}$$

$$\bar{\mathbf{d}}^1(y, \tilde{\mathbf{r}}) = \begin{cases} \frac{\mathbf{d}}{\tilde{\mathbf{r}}} + 1 & y \in \hat{Y} \\ \mathbf{d} & \text{otherwise} \end{cases}$$

and  $Y(h_t) = \{y \in Y^e \cup \hat{Y} \mid \mathbf{r}(y \mid \bar{\mathbf{a}}^1(h_t), \mathbf{a}^{-1}(h_t)) > 0\}$ .

**Lemma 3:** *In a participation game if  $\mathbf{a}^{-1}(h_t) \in \text{conhull}(E^{-1})$ , or  $\mathbf{a}^{-1}(h_t) \notin \text{conhull} E^{-1}$  and  $\mathbf{a}_0^1(h_t) \geq c f^1$  for some  $c > 0$  and vulnerable friendly action  $f^1$  with temptation bounds  $\underline{\mathbf{r}}, \tilde{\mathbf{r}}$  then*

$$v^1(h_t) \leq \max_{y \in Y(h_t)} (1 - \mathbf{d})\bar{u}^1(y, \tilde{\mathbf{r}}) + \bar{\mathbf{d}}(y, \tilde{\mathbf{r}})v^1(h_t, y).$$

*Remark:* This lemma says that if the rational type is playing a friendly strategy, his payoff is bounded by a one-period gain and the continuation payoff conditional on a bad signal. This follows from the assumption that for every entry-inducing strategy it is possible to lower the probability of all of the signals in  $\hat{Y}$  by at least  $\mathbf{j}$  while increasing the probability of

---

<sup>8</sup> If there were a type with a history-dependent strategy, this part of the lemma would need to be modified.

each other signal by at least the multiple  $\tilde{\mathbf{r}}_{\min}$ . The fact that the rational type chooses not to reduce the probability of the bad signal means that the continuation payoff after the bad signal cannot be much worse than the overall continuation payoff.

*Proof of Theorem 1:* Given an equilibrium, we begin by constructing a positive probability sequence of histories beginning with an initial history at date 0. Given  $h_t$  already constructed, we define  $h_{t+1} = (h_t, y_{t+1})$  where

$$y_{t+1} \in \arg \max_{y \in Y(h_t)} (1 - \mathbf{d})\bar{u}^1(y, \tilde{\mathbf{r}}) + \bar{\mathbf{d}}(y, \tilde{\mathbf{r}})v^1(h_t, y).$$

We know that  $Y(h_t)$  is not empty because either  $\mathbf{a}^{-1}(h_t) \in \text{conhull } E^{-1}$ , or  $\mathbf{a}^{-1}(h_t) \notin \text{conhull } E^{-1}$ . This latter case implies that  $\bar{\mathbf{a}}^1(h_t) \geq \mathbf{g}f^1$  for some friendly  $f^1$ , and since only enforceable actions can be played in equilibrium, this  $f^1$  must be vulnerable to temptation, so  $\mathbf{r}(\hat{Y} | \bar{\mathbf{a}}^1(h_t), \mathbf{a}^{-1}(h_t)) \geq \mathbf{g}\mathbf{r}(\hat{Y} | f^1, \mathbf{a}^{-1}(h_t)) > 0$ .

Now apply Lemma 2 to conclude that for each  $h_t$  at most  $k^*$  of the signals are in  $\hat{Y}$  and  $\mathbf{m}(h_t)[\Theta(F^1)] \leq \mathbf{g}\underline{F}^1/2$ . Consider an  $h_t$  such that  $\mathbf{a}^{-1}(h_t) \notin \text{conhull } E^{-1}$ . From the definition of a friendly action, we know that  $\bar{\mathbf{a}}^1(h_t) \geq \mathbf{g}f^1$  for some friendly  $f^1$ , so  $\mathbf{m}(h_t)[\Theta(F^1)] \leq \mathbf{g}\underline{F}^1/2$  and Lemma 1 implies that  $\mathbf{a}_0^1(h_t) \geq \mathbf{g}f^1/2$ . Now apply Lemma 3 to conclude that for each  $h_t$

$$v^1(h_t) \leq (1 - \mathbf{d})\bar{u}^1(y_{t+1}, \tilde{\mathbf{r}}_{\min}) + \bar{\mathbf{d}}(y_{t+1}, \tilde{\mathbf{r}}_{\min})v^1(h_{t+1})$$

Since  $v^1(h_t) \leq U^1$ , it follows that

$$v^1(0) \leq (1 - \mathbf{d})\sum_{t=1}^{\infty} \prod_{t=2}^t \bar{\mathbf{d}}(y_t, \tilde{\mathbf{r}}_{\min})\bar{u}^1(y_t, \tilde{\mathbf{r}}_{\min}).$$

Since  $\bar{u}^1(y^e, \tilde{\mathbf{r}}_{\min}) = 0$ , and  $y_t \in \hat{Y}$  at most  $k^*$  times, this gives the desired bound. Notice that the fact that  $\bar{u}^1(y^e, \tilde{\mathbf{r}}_{\min}) = 0$  follows from the assumption of exit minmax: it is here that we make use of the fact that exit gives the long-run player no more than the minmax.

☑

## 4. Examples

We now consider a number of examples to illustrate the scope of Theorem 1, and also the extent to which the assumptions are necessary as well as sufficient. To begin, Example 4.1 illustrates what happens when the prior puts too much weight on some committed types for the hypothesis of commitment size  $\underline{g}F^1/2$  to be satisfied. Example 4.2 shows that the EV conclusion is not robust to the addition of an observed action that makes the short-run players want to enter. Example 4.3 examines participation games that are not bad reputation games, and example 4.4 illustrates the role of the exit-minmax assumption. In all of the examples but 4.1, we assume that the hypothesis of commitment size  $\underline{g}F^1/2$  is satisfied, and investigate whether the game is a bad reputation game. The following section considers a class of bad-reputation principal-agent games.

### 4.1: EV With Stackelberg Type

We have relaxed the original assumptions of EV in a number of ways. One important extension is that we allow for positive probabilities of all commitment types. In particular, we allow a positive probability of a “Stackelberg type” committed to the honest strategy  $et$ , which is the optimal commitment. However, a hypothesis of the theorem is that the prior satisfy the commitment size assumption.

Here we illustrate that assumption in the context of the EV example. Suppose in particular that there are 3 types, rational, bad, and Stackelberg. The set of possible priors can be represented by the simplex in figure 2.

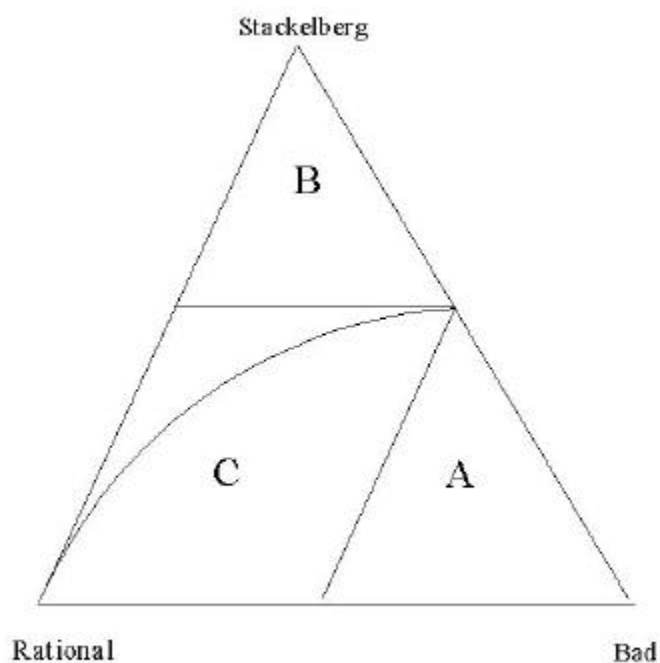


Figure 2

When the prior falls into the region A, the probability of the bad type is too high, and the short run players refuse to enter regardless of the behavior of the rational type. Bad reputation arises because the long-run player tries to prevent the posterior from moving into this region. In EV the prior assigned probability zero to the Stackelberg type. Thus the prior and all posteriors on the equilibrium path belong to the lower boundary of the simplex. When there is a sufficiently high probability of the Stackelberg type, the short-run players will enter regardless of the behavior of the rational type; this is region B. Note that the boundaries of these regions intersect on the right edge of the simplex: this point represents the mixture between  $ee$  and  $et$  which makes the short-run player indifferent between entry and exit.

When the prior is in region A, there will be no entry and the long-run player obtains the minmax payoff of zero. On the other hand, when the prior is in region B, there is a Nash (and indeed sequential) equilibrium in which the long run player receives the best commitment payoff, which is “ $u$ ” in the notation of EV. The equilibrium is constructed as follows.

Consider the game in which the posterior probability of the bad type is zero. In this game there exists a sequential equilibrium in which the long-run player gets  $u$ . Suppose that we assume that this is the continuation payoff in the original game in any subform in which the long-run player played  $t$  at least once in the past. A sequential equilibrium of this modified game is clearly a sequential equilibrium of the original game, and by standard arguments, this modified game has a sequential equilibrium. How much does the rational long-run player get in this sequential equilibrium? One option is to play  $tt$  in the first period. Since the short-run player is entering regardless, this means that beginning in period 2 the rational type gets  $u$ . In the first period he gets  $(u - w)/2$ . Hence in equilibrium he gets at least  $(1 - d)(u - w)/2 + du$ , which converges to  $u$  as  $d \rightarrow 1$ .

Our theorem is about the set of equilibrium payoffs for priors outside of these two regions. The theorem states that there is a curve, whose shape is represented in the figure, such that when the prior falls below this curve (region C), the set of equilibrium payoffs for the long-run player is bounded above by a value that approaches the minmax value as the discount factor converges to 1. The diagram shows that the left boundary of the simplex is an asymptote for this curve as it approaches the complete information prior (i.e.  $m(0)(q_0) = 1$ ) in the lower left corner. This illustrates the important aspect of the commitment size restriction: it is satisfied for “almost all” sufficiently small perturbations of the complete information limit.

### 4.2: Adding an Observed Action to EV

We now modify the EV game by adding a new observable action "g" for the long-run player called "give away money." This action induces the short-run players to participate ( $B(g) \cap \text{conhull}(E^{-1}) = \emptyset$ ) so it must be in every friendly set. Since the action is observable, it is not vulnerable to temptation with respect to any signals that are unambiguous for the unfriendly actions, so this is not a bad reputation game. Moreover, even without a Stackelberg type the EV conclusion fails in this game: there is an equilibrium where the rational type plays g in the first period. This reveals that he is the rational type, and there is entry in all subsequent periods, while playing anything else reveals him to be the bad type so that all subsequent short run players exit. Thus the assumption that every friendly action is vulnerable to temptation is seen to be both important and economically restrictive.

### 4.3. Orthogonality Issues

Suppose friendly actions send the bad signal by putting positive weight on unfriendly actions. An important class of games in which this is the case are those in which, conditional on entry, the long-run players' actions are observed. In this case the bad signals correspond to unfriendly actions, and bad signals can only have positive probability when the unfriendly action is played with positive probability. Moreover, in some games, the only friendly strategies involve randomizing in this way.

**Proposition 1:** *If for every friendly set and unfriendly set there is a friendly action whose restriction to the complement of the unfriendly set has probability 0 of generating a signal that is unambiguous for the unfriendly set, the game is not a bad reputation game.*

*Proof:* The assumption that the friendly and unfriendly sets are orthogonal is violated.

☑

To see that this makes a difference, consider the following two-person game with observed actions:

	L	M	R
U	0,4	1,3	0,0
D	0,0	1,3	0,4

Figure 3

where  $L$  and  $R$  correspond to exit and  $M$  to entry.<sup>9</sup> In this case entry can be induced only by mixing with probability of  $U$  between  $\frac{1}{4}$  and  $\frac{3}{4}$ .

We will first explain why the bad reputation theorem does not hold here, and then show that its conclusion fails as well.

Both  $U$  and  $D$  are unfriendly, and we need to choose either one or both of them to be in the unfriendly set. If we include both actions in the unfriendly set, then it is impossible to find an orthogonal friendly set. If we include only one of the actions in the unfriendly set, and chose a friendly set that includes mixed actions, then orthogonality is again violated, while if we specify that the friendly set is the singleton corresponding to the other action, then the friendly action is not vulnerable to temptation.

The conclusion of the theorem fails here as well: suppose that the only commitment type with positive probability is  $D$ , and that the probability of the bad type is less than  $\frac{1}{4}$ . Consider the following strategies: For any current probability  $\mathbf{m}(h_t)[D]$  less than  $\frac{1}{4}$  the rational type mixes so that the overall probability of  $D$  is exactly  $\frac{3}{4}$ . (In particular, this is true when the long-run player has been revealed to be rational, so that  $\mathbf{m}(h_t)[D] = 0$ .) The short-run player always enters. If  $U$  is observed,

---

<sup>9</sup> In this example, the short-run player has several exit actions, and his payoff depends on the long-run player's action. This is a necessary feature of two-player games where the only friendly strategies are mixed, but it is not necessary in three-player games – think of a game where player 3 has veto power, 3 only plays In if 2 plays M, and 2's payoffs to M are as in the payoff matrix of this example.



the type is revealed to be rational. If  $D$  is observed, the probability of the bad type increases by a factor of  $4/3$ , so when it first exceeds  $1/4$  it is at most equal to  $1/3$ . At this point, the rational type may reveal himself by playing  $U$  with probability 1, while preserving the incentive of the short-run player to enter. It is easy to see that this is a Nash equilibrium for any discount factor of the long-run player, yet in this equilibrium, the long-run player's payoff is 1.

We say that an action  $f^*$  is *sufficient for entry* if, for some  $\underline{a} < 1$ ,  $\mathbf{a}^1 \geq \underline{a}f^*$  implies that there is  $\mathbf{a}^2 \in B(\mathbf{a}^1)$  with positive probability of entry. In the example above the friendly action is sufficient for entry, and sends the bad signal only because of mixing onto the unfriendly action. That is, the sufficient action mixes between a pure action that does not send the bad signal, and an unfriendly action. If there is a friendly action that does not send the bad signal at all, then the game is not a bad reputation game since such an action cannot admit a temptation. More strongly, if an action sufficient for entry does not send the bad signal at all, then a patient rational player can do almost as well as in the absence of bad reputation effects.

**Proposition 2:** *If there is an action  $f^*$  that is sufficient for entry and does not send any bad signal, the only commitment types are unfriendly types, and the probability of commitment types is sufficiently low, then as  $\mathbf{d} \rightarrow 1$  there are sequential equilibrium payoffs for the rational type that approach the highest sequential equilibrium payoff without committed types.*

*Proof:* Suppose that the prior probability of committed types is sufficiently low that the short-run players will enter when the rational type plays  $f^*$ . Then it is a sequential equilibrium for the rational type to play  $f^*$  in the initial period with entry by the short-run players. Subsequently, if a bad signal was observed, the short-run players stay out. If a bad signal

was not observed, the probability of committed types is zero, and the continuation equilibrium is the best possible without committed types. On the equilibrium path, the rational type payoff clearly approaches that of the highest payoff without committed types, since he gets that amount beginning in period 2, and payoffs in period 1 are bounded below.

☑

#### 4.4: Exit Minmax

In participation games, reputation plays a role because the short run players will guard against unfriendly types by exiting. This is “bad” for the long-run player only if exit is worse than the payoff he otherwise would receive, and the exit minmax assumption ensures that this is the case.

In participation games without exit minmax, there are outcomes that are even worse for the long-run player than obtaining a bad reputation. In this case it is possible that there exist equilibria in which the long-run player is deterred from his temptation to avoid exit by the even stronger threat of a minmaxing punishment. For example consider the game in Figure 4, where the first matrix represents the payoffs, and the second represents the distribution of signals conditional on entry.

	<i>In</i>	<i>Out</i> <sup>1</sup>	<i>Out</i> <sup>2</sup>
<i>F</i>	1,1	0,0	-2,0
<i>U</i>	1,0	0,1	-2,0
<i>T</i>	1,0	0,0	-2,1

	<i>g</i>	<i>b</i>	<i>r</i>
<i>F</i>	½	½	0
<i>U</i>	0	1	0
<i>T</i>	½	0	½

Figure 4

This game is a participation game with exit actions  $Out_1$  and  $Out_2$ , unfriendly action  $U$  and friendly action  $F$  vulnerable to temptation  $T$ . There are only two types, the rational type and a bad type that plays  $U$ . Exit minmax fails because the maximum exit payoff exceeds the minmax payoff, and we claim that there are good equilibria in this game because

the threat of exiting with  $Out_2$  is worse than the fear of obtaining a reputation for playing  $U$  which would only lead to exit with  $Out_1$ .

To see this, consider the following strategy profile. The rational type plays  $F$  at every history unless the signal  $r$  has appeared at least once; in that case the rational type plays  $T$ . The short run player plays  $Out_2$  if a signal of  $r$  has ever appeared. Otherwise, the short run player plays  $Out_1$  if the posterior probability of the bad type exceeds  $\frac{1}{2}$  and  $In$  if this probability is less than  $\frac{1}{2}$ . Observations of  $r$  are interpreted as signals that the long-run player is rational.

Since  $(T, Out_2)$  is a Nash equilibrium of the stage game, the continuation play after a signal of  $r$  is a sequential equilibrium. When  $r$  has not appeared, the long run player optimally plays  $F$ . Playing  $U$  gives no short-run gain and hastens the onset of  $Out_1$ , and playing  $T$  shifts probability from the bad signal  $b$  to the signal  $r$  which is even worse.<sup>10</sup> The short-run players are playing short-run best responses. In this equilibrium, the long run player does not give in to the temptation to play  $T$ . As a result, with positive probability, the short-run players never become sufficiently pessimistic to begin exiting, and so the long run player achieves his best payoff.

In the above example there were two exit actions. The next proposition states that when there is only one exit action and the long-run player's exit payoff is independent of his own action, the worst Nash equilibrium payoff for the long run player is (not much worse than) his exit payoff. Note that this condition is satisfied in the principal-agent applications discussed in section 5. The proposition is a consequence of FL (1992).

**Proposition 3:** *Consider a participation game with a single short-run player and a unique exit action. If*

---

<sup>10</sup> Playing  $T$  gives probability  $\frac{1}{2}$  of shifting to the absorbing state where payoffs are  $-2$ . Playing the equilibrium action of  $F$  has probability at most  $\frac{1}{2}$  of switching to the state where payoffs are 0.

(i) there exists a pure action<sup>11</sup>  $\hat{a}^1$ , such that  $B(\hat{a}^1) = \{exit\}$ ,

(ii) the prior distribution assigns positive probability to a type that is committed to  $\hat{a}^1$ ,

and

(iii) the long-run player's action is identified conditional on entry then there is a lower bound on the payoffs to the rational type which converges to the exit payoff, as the discount factor approaches 1.

*Proof:* FL (1992) established<sup>12</sup> that for any game there exists a lower bound  $b(\mathbf{d})$  on the set of Nash equilibrium payoffs for the rational type, and that as  $\mathbf{d} \rightarrow 1$ ,  $b(\mathbf{d})$  converges to a limit that is at least

$$\max_{a^1 \in C} \min_{a^2 \in \tilde{B}(a^1)} u^1(\mathbf{a}^1, \mathbf{a}^2)$$

where  $\tilde{B}(a^1)$  is the set of self-confirmed best-responses to for the short-run player to  $a^1$ , and  $C$  is the set of actions corresponding to the support of the prior distribution over commitment types. Because the long-run player's action is identified conditional on entry and  $B(\hat{a}^1) = \{exit\}$ , we have  $\tilde{B}(\hat{a}^1) = \{exit\}$ , and because the type that plays  $\hat{a}^1$  has positive prior probability, the FL (1992) bound is at least  $u^1(\hat{a}^1, exit)$ .

☑

For games satisfying the conditions of the proposition, the exit minmax condition is not necessary for bad reputation. The worst

---

<sup>11</sup> The assumption that this is a pure action is not necessary here; we state the result this way for consistency with the rest of the paper.

<sup>12</sup> The statement of the FL theorem requires that commitment types including mixing types have full support, in which case the set  $C$  is the space of all (mixed) actions, but the proof given there also shows that the version of the lower bound given here is correct.

equilibrium continuation value that the short-run players could inflict is arbitrarily close to the exit payoff and hence a patient long run player could not be deterred from his temptation to avoid a bad reputation.

## 5. Poor Reputation Games and Strong Temptations

Recall that an action is vulnerable to a temptation if when the short-run players participate, the temptation lowers the probability of all bad signals, and increases the probability of all others. In this case the bad reputation result requires the exit minmax condition, as demonstrated by the example in Section 4.4. Notice, however, that in the example the relative probability of  $g$  and  $r$  is changed by the temptation. If the temptation satisfies the stronger property that the relative probability of the other signals remains constant, then we can weaken the assumption of exit minmax. In this section we develop this result, and give an application to games with two actions.

First we give a formal definition of a strong temptation:

**Definition 10:** *An action  $a^1$  is vulnerable to a strong temptation relative to a set of signals  $\hat{Y}$  if there exists a number  $\underline{r} > 0$  and an action  $b^1$  such that*

$$1) \text{ If } a^{-1} \notin E^{-1}, \hat{y} \in \hat{Y} \text{ then } \mathbf{r}(\hat{y} | b^1, a^{-1}) \leq \mathbf{r}(\hat{y} | a^1, a^{-1}) - \underline{r}$$

$$2) \text{ If } a^{-1} \notin E^{-1} \text{ and } y, y' \notin \hat{Y} \cup Y^e \text{ then } \frac{\mathbf{r}(y | b^1, a^{-1})}{\mathbf{r}(y' | b^1, a^{-1})} = \frac{\mathbf{r}(y | a^1, a^{-1})}{\mathbf{r}(y' | a^1, a^{-1})}.$$

$$3) \text{ For all } e^{-1} \in E^{-1}, u^1(b^1, e^{-1}) \geq u^1(a^1, e^{-1}).$$

The action  $b^1$  is called a strong temptation.

The first and third parts of this definition are the same as in the definition of a temptation; the additional strength comes from part (2), which requires that the temptation not merely increase the probability of all of the good signals, but leave their relative probabilities unchanged.

Note that strong temptation is equivalent to temptation in games in which the set  $Y \setminus (\widehat{Y} \cup Y^e)$  has a single element, for example games in which there are only two entry signals; in particular applies when the game of Section 4.4 is modified so that the only signals when entry occurs are  $g$  and  $r$ .

This condition lets us prove an analog of lemma 3: and  $\mathbf{d}$ :

**Lemma 4:** *In a participation game, if  $\mathbf{a}^{-1}(h_t) \in \text{conhull } E^{-1}$  or  $\mathbf{a}^{-1}(h_t) \notin \text{conhull } E^{-1}$  and  $\mathbf{a}_0^1(h_t) \geq \mathbf{g}f^1$  for some  $\mathbf{g} > \mathbf{0}$  and friendly action  $f^1$  that is vulnerable to a strong temptation size  $\underline{\mathbf{r}}$ , then*

$$v^1(h_t) \leq \max_{y \in Y(h_t)} (1 - \mathbf{d})\bar{u}^1(y, \underline{\mathbf{r}}) + \bar{\mathbf{d}}^1(y, \underline{\mathbf{r}})dv^1(h_t, y),$$

$$\text{where } \bar{u}^1(y, \mathbf{r}) = \begin{cases} \left(1 + \frac{1}{|\widehat{Y}| \mathbf{r}}\right) U^1 & \text{if } y \in \widehat{Y} \\ \hat{u}^1 & \text{otherwise} \end{cases}$$

$$\text{and } \bar{\mathbf{d}}^1(y, \mathbf{r}) = \begin{cases} \mathbf{d} \left(1 + \frac{1}{|\widehat{Y}| \mathbf{r}}\right) & y \in \widehat{Y} \\ \mathbf{d} & \text{otherwise} \end{cases}.$$

The proof, which is in the Appendix, follows that of lemma 3, but takes advantage of the fact that the long-run player's continuation expected value, conditional on a friendly action, a non-exit profile, and a signal not in  $\widehat{Y} \cup Y^e$ , is the same for the equilibrium action and the strong temptation  $b^1$ .

Define

$$\hat{u}^1 = \max_{\mathbf{a}^1, \mathbf{a}^{-1} \in \text{conhull}(E^{-1}) \cap \text{image}(B)} u^1(\mathbf{a}^1, \mathbf{a}^{-1})$$

This is a bound on the long-run player's payoff when the short-run players play exit actions that are a best response to some (possibly incorrect) conjectures.

**Definition 11:** A participation game is a poor reputation game if there is an unfriendly set  $N^1$ , a friendly set  $F^1$  that is orthogonal to  $N^1$ , and a set of signals  $\hat{Y}$  that are unambiguous for  $N^1$ , and such that every enforceable  $f^1 \in F^1$  is vulnerable to strong temptation relative to  $\hat{Y}$ .

The next result says that poor reputation games have much the same consequences as bad reputation games.

**Theorem 2:** In a poor reputation game of commitment size  $\mathbf{g}F^1/2, \mathbf{h}$ ,

$$\lim_{d \rightarrow 1} \bar{v}^1 \leq \hat{u}_1.$$

With lemma 4 in hand, the proof of Theorem 2 is very close to that of Theorem 1, and is omitted. Notice that it is possible for a game to be both a bad reputation game and a poor reputation game, and, since strong and ordinary temptation are equivalent when  $Y \setminus (\hat{Y} \cup Y^e)$  is a singleton, the two are necessarily equivalent in this case. The original EV game is such an example. Notice also that example 4.4 in which we construct a non-bad equilibrium has three signals rather than two. With two signals, the game would still fail the exit minmax condition and fail to be a bad reputation game, but it would still be a poor reputation game, and would not admit a good equilibrium. Finally, observe the proofs of both Lemma 3 and 4 can be generalized, so that the difference between the best equilibrium payoff (in the limit as  $d \rightarrow 1$ ) and the most favorable outcome with exit is bounded by a scale factor times the the product of two terms, namely (i) the change in relative probabilities induced by a temptation and (ii) the excess of the best result given exit over the minmax. In particular, the bound on the difference is continuous in the each of these terms, so that if either is small the best equilibrium payoff for a patient long-run player can only exceed the best exit payoff by a small amount.

We turn now to the special case of two-player participation games where there is only one signal in  $\hat{Y}$  and short-run player payoffs depend only on the signal. We focus on the case where one signal in  $Y \setminus (\hat{Y} \cup Y^e)$ , so that bad reputation implies poor reputation. We show

that these games are not poor reputation games (and by implication not bad reputation games either).

**Proposition 4:** *In a two-player participation game suppose there are only two “entry signals” (that is two elements of  $Y - Y^e$ ), that the short-run player has only two actions, and that the short-run player’s realized payoff is determined by the signal. Then the game is not a poor reputation game.*

*Proof:* Notice that since the short-run player has only two actions, they correspond to “entry” and “exit” respectively. Consequently, the short-run player payoff conditional on entry depends only on the distribution over signals induced by the long-run player action. If we normalize the short-run player’s payoff function so that his exit payoff is 0, and suppose that both the friendly and unfriendly sets are non-empty, then one signal yields a negative payoff and the other signal’s payoff is positive; call these the “bad” and “good” signals respectively. If the game has no non-empty unfriendly set, it is not a poor reputation game; so we can suppose there is at least one non-empty unfriendly set. Any unfriendly set  $\widehat{A}^1$  consists of actions with a sufficiently high probability of sending the bad signal, and the bad signal (as a singleton set) is the only set  $\widehat{Y}$  that can be unambiguous for  $\widehat{A}^1$ . Let  $f^1$  be the friendly action in the (finite) friendly set that maximizes the short-run player’s payoff. The payoff to this action, conditional on it not generating the bad signal with the negative payoff, is positive, and since any temptation relative to  $\widehat{Y}$  must reduce the probability of the bad signal, a temptation must give the short-run player a higher payoff than this “friendliest” friendly action. For this to be true, there must be a pure strategy  $\widehat{b}^1$  that gives the short-run player at least this same utility. Clearly  $\widehat{b}^1$  induces entry, and since it is a pure strategy, it must be in the friendly set. This contradicts the fact that  $f^1$  was assumed to maximize short-run player utility in the friendly set.

☑



We believe that the assumptions of this proposition imply that there is an equilibrium where the rational type's payoff is bounded below by a positive number as  $d \rightarrow 1$  but we have not been able to show this.

## 6. Principal-Agent Entry Games

In this section we consider a class of applications which have the nature of an agency relationship. The long-run player (the agent) takes an action that affects the payoffs of both a principal (that period's short run player) and herself. When the principal's and the agent's preferences differ over the action set, and the action is not perfectly observed, we have a classical problem of incentives. A repeated interaction can often substitute for explicit contracts in alleviating this incentive problem, as the long run agent's objective of establishing a good reputation can provide an incentive for efficient behavior. In this section we classify agency environments in which the repeated interaction has the opposite effect: Bad reputation can intensify rather than mitigate the agent's incentive problem.

There is a single short-run player (the principal) whose only choice is whether to enter or to exit. If the principal enters, then the long-run player (the agent) chooses a payoff-relevant action, otherwise both players receive a reservation value which is normalized to zero. Formally  $A^2 = \{exit, enter\}$  and  $u^2(a^1, exit) = 0$  for each  $a^1 \in A^1$ . For simplicity we write  $u^2(a^1, enter) = u^2(a^1)$ . We assume there is an action  $a^1 \in A^1$  for which  $u^1(a^1) \geq 0$ , so that the exit minmax assumption is satisfied. (Note that this assumption will hold whenever the principal has the option to refuse to participate. Note also that from Theorem 2 this assumption is not necessary for games with two signals.)

For these games we can immediately identify the relevant friendly set. Define

$$F^1 = \{a^1 \in A^1 : u^2(a^1) \geq 0\}$$

which is the set of pure friendly actions. We know that  $F^1 \subset \hat{F}^1$  for any friendly set  $\hat{F}^1$ . In fact, within the class of principal-agent games, any bad

reputation game is a bad reputation game with friendly set  $F^1$ . To see this note that if  $\mathbf{a}^1(F^1) = 0$  then  $u^2(\mathbf{a}^1) < 0$ , i.e. exit is the unique best reply to  $\mathbf{a}^1$ . Thus  $F^1$  is itself a friendly set.<sup>13</sup> Furthermore  $\text{supp}(F^1) \subset \text{supp}(\hat{F}^1)$  so that orthogonality is preserved, and if every  $f^1 \in \hat{F}^1$  is vulnerable then every  $f^1 \in F^1$  is vulnerable. Thus we can restrict attention to  $F^1$ .

To show that these are bad reputation games, it suffices to find an unfriendly set orthogonal to  $F^1$  with unambiguous signals, such that every enforceable point in  $F^1$  is vulnerable to a temptation.

### 6.1 Games with Hidden Information

In these games the principal has some private information that is relevant for a decision affecting both principal and agent. Each period, nature draws a state  $\mathbf{w} \in \Omega$  independently from a probability distribution that we denote by  $p$ .<sup>14</sup> The agent privately observes the state and then selects a decision  $d \in D$ . Conditional on the realized state and the decision of the agent, a signal  $z \in Z$  is drawn from the distribution  $m(z|\mathbf{w}, d)$  where we assume that  $m(z|\mathbf{w}, d) > 0$  for each  $z, \mathbf{w}$ , and  $d$ . Future short run players observe both  $z$  and the decision  $d$ . Each player  $j$  has state-dependent utility function  $\mathbf{p}^j(\mathbf{w}, d, z)$  and evaluates stage payoffs according to expected utility with respect to the distributions  $p(\mathbf{w})$  and  $m(z|\mathbf{w}, d)$ .

To apply Theorem 1, we find conditions under which this defines a bad reputation game. The set of actions  $A^1$  for the long-run player is the set of maps  $a^1 : \Omega \rightarrow D$ . The stage-game utility function is

$$u^j(a^1) = \sum_{\mathbf{w} \in \Omega} p(\mathbf{w}) \sum_{z \in Z} m(z|\mathbf{w}, d) \mathbf{p}^j(\mathbf{w}, a^1(\mathbf{w}), z).$$

Finally  $Y - Y^e = Z \times D$  and

---

<sup>13</sup> When there is more than one principal, this conclusion does not follow, and mixed friendly sets will generally have to be considered. See the discussion in section 4.3 and footnote 7.

<sup>14</sup> This is a slight abuse of notation, as  $p$  also denotes the probability distribution over types in the incomplete-information games, but no ambiguity should result.

$$r((z,d)|a^1,entry) = \sum_{\{\mathbf{w}:a^1(\mathbf{w})=d\}} p(\mathbf{w})m(z|\mathbf{w},d)$$

**Proposition 5:** *The hidden information game is a bad reputation game if there exists a decision  $d$  such that  $a^1 \in F^1$  implies  $\emptyset \neq \{\mathbf{w}:a^1(\mathbf{w})=d\} \neq \Omega$ .*

*Proof:* Let  $a(d)$  denote the constant action that chooses  $d$  regardless of the state  $\mathbf{w}$ , and take  $\hat{A} = \{a(d)\}$ . Because  $m(z|\mathbf{w},d) > 0$ , the set of signals  $Y^d = Z \times \{d\}$  is unambiguous for  $\hat{A}$ . If  $a^1$  is friendly, then  $\Omega_d = \{\mathbf{w}:a^1(\mathbf{w})=d\} \neq \emptyset$ . For each  $b \neq d$  let  $b^1$  be the action defined by  $b^1(\mathbf{w}) = b$  and  $b^1(\mathbf{w}) = a^1(\mathbf{w}), \mathbf{w} \notin \Omega_d$ . Then any mixed strategy that puts positive weight on every  $b^1$  will decrease the probability of signals in  $Y^d$  and increase the probability of all other non-exit signals, so  $a^1$  is vulnerable to such mixed strategy.

☑

Many examples can be found that meet the condition of the proposition, and the EV example is a special case; the theorem extends the example to allow for public signals  $z$  about the short run players' realized payoffs.. For illustration, consider the following extension of the EV example. If the correct repair is chosen, then the car works, otherwise it does not, and this outcome is observed by future motorists. Proposition 5 implies that as long as the mechanic's diagnosis is not perfect, the game is a bad reputation game. Formally, let  $Z = \{work, not\}$ . Suppose that for each motorist (independently and with equal probability), nature selects a necessary repair from the set  $\{tuneup, engine\}$ . Conditional on this, the mechanic observes state  $\mathbf{w} \in \{tuneup, engine\}$  representing the mechanic's diagnosis, and selects a repair  $d \in \{tuneup, engine\}$ . Suppose that with positive probability (but as small as desired)  $\mathbf{w}$  is incorrect, i.e. not equal to the repair that is truly necessary. Then it follows that  $m(z|\mathbf{w},d) > 0$  for each  $z, \mathbf{w}$ , and  $d$ . We can now define the motorists' payoff as a function of  $z, \mathbf{w}$ , and  $d$  to yield expected payoff function  $u^2(a^1)$  identical to the original EV stage-game payoffs. Thus the

friendly set (which is defined only in terms of the short-run player's payoffs) from EV remains so, and we can apply Proposition 5 using the decision  $d = engine$  to show that the modified game is a bad reputation game.

The strength of Proposition 5 raises the question of whether people typically do employ advisors, and if so, why they do. One set of responses is that in some cases the advisor is not really used for advice, but is either employed as a costly signal, or is hired not for the advice *per se* but for its implementation. For example, a country might know the advice that the IMF will recommend, but find it useful to delegate the implementation of the advice so that it can avoid taking full responsibility for the resulting hardships.

A second set of responses explains why Proposition 5 does not apply even though the advisor really is relied on for advice. Some advice games aren't participation games because the advisor can make "speeches" even without any "customers;" this may describe political advisors and investment columnists. Alternatively, perhaps some of the short-run players are "naïve" and enter even when entry is not a best response to equilibrium play. These sorts of "noise players" ensure that the long player can build a track record; formally, the game would not once again not be a participation game because the long-run player's action would always have some effect on the expected distribution of signals given the play of the "rational" short-run players.

The third type of response builds on the fact that Theorems 1 and 2 do not pin down play, as opposed to payoffs, and are most powerful for discount factors that are close to 1. Thus, our results would be consistent with a game where the rational type played "honestly" (or Stackelberg) for a long time for a given fixed discount factor.

## 6.2 Rules Rather than Discretion

We can build on the analysis of hidden information games to discuss the emergence of rules over discretion in agency relationships. To motivate the idea, consider college admissions. The university (the long-run agent) receives an application. The applicant is described by a set of characteristics  $\mathbf{w} \in \Omega = \Omega^o \times \Omega^n$ . Some ( $\Omega^o$ ) of these characteristics are publicly observable (for example race and SAT scores) and others ( $\Omega^n$ ) are observed only by the university. This may include information that is truly private (like an interview) or information that require the expertise of the agent to interpret (for example, the strength of the applicant's high school.) A pure strategy for the university is a map from characteristics to the decision space  $D = \{\text{admit, deny}\}$ . The probability of drawing characteristics  $\mathbf{w}$  is  $p(\mathbf{w}) > 0$ . The university's preferences over applicants are summarized by the payoff function  $\mathbf{p}^1(\mathbf{w})$  if the student is admitted, and  $R$  if the student is denied.

The short-run principal, player 2, is the state governor who chooses between allowing the university discretion in admissions, or imposing a rigid admission rule based on observable characteristics. There are many possible rules that the principal might use, but since she is a short-run player we can restrict attention to the rule that maximizes the principal's expected short-run payoff. This rule is a mapping  $g : \Omega^o \rightarrow D$  that mandates admission if and only if  $\mathbf{w}^o \in g^{-1}(\text{admit})$ . The imposition of a rigid admission rule represents "exit." The public signal at date  $t$  is  $y_t = (d_t, a_t^2, \mathbf{w}^o)$ , where  $d_t \in D$  is the date- $t$  decision, and any signal with  $a_t^2 = \text{rule}$  is an exit signal. The governor shares the same preferences as the university, receiving a utility of  $\mathbf{p}^1(\mathbf{w})$  for admits and  $R$  for rejects.

Because the university can always implement  $g$  its own, exit minmax condition is satisfied. In order for discretion to improve upon,  $g$  for some set of verifiable characteristics, the admission decision should depend on the unverifiable characteristics. That is  $a^1 \in F^1$  only if  $a^1(\mathbf{w}^o, \mathbf{w}^n) = \text{admit}$  and  $a^1(\mathbf{w}^o, \hat{\mathbf{w}}^n) = \text{deny}$  for some  $\mathbf{w}^n, \hat{\mathbf{w}}^n$  and  $\mathbf{w}^o$ . Then

by essentially the same argument as in Proposition 5, the game is a bad reputation game with unfriendly set

$$\{a^1 : a^1(\mathbf{w}^o, \cdot) = \text{deny}\}$$

For example,  $\mathbf{w}^o$  may be racial characteristics, and the types associated with this unfriendly set represent the governor's fear that the university admissions are biased against members of the race in question.

## 6.2. Games with Hidden Actions

On the other hand, agency games with hidden actions, or moral hazard, tend not to be susceptible to bad reputation effects. The problem is that the second part of the definition of temptation typically fails because deviations will generally lower the probability of some good signals. However, a special case in which a hidden action game is a bad reputation game occurs when there is only one short-run player and only two signals.

The following proposition is an immediate application of the definition of a bad reputation game in this setting.

**Proposition 6:** *Suppose in a principal-agent entry game that  $Y - Y^e = \{y^L, y^H\}$  and that  $\hat{a}^1$  strictly maximizes the probability of  $y^L$  with  $u^2(\hat{a}^1) < 0$ . If for every friendly enforceable  $a^1$  there is a  $b^1$  such that  $\mathbf{r}(y^L | b^1) < \mathbf{r}(y^L | a^1)$  the game is a bad reputation game.*

To apply this result, suppose that the agent chooses an action from a one-dimensional set ordered so that higher actions are more likely to give rise to the high signal. Specifically, we let  $A^1 = \{\underline{a}^1, \dots, \bar{a}^1\} \subset \mathfrak{R}$  and  $\mathbf{r}(y^H | a^1)$  are an increasing function of  $a^1$ . We assume that  $u^2(a^1)$  is concave so that  $F^1$  is an interval. Whether or not the game is a bad reputation game then depends on whether the principal prefers extreme or interior actions.

**Proposition 7:** *The hidden action game with two non-exit signals is a bad reputation game if and only if  $\{\underline{a}^1, \bar{a}^1\} \subset A^1 - F^1$ .*

*Proof:* Suppose  $\{\underline{a}^1, \bar{a}^1\} \subset A^1 - F^1$ . Then  $\hat{a}^1 = \underline{a}^1$  and every friendly action  $a^1 > \hat{a}^1$  is vulnerable to the (unfriendly) temptation  $\bar{a}^1$ . On the other hand if,  $\bar{a}^1 \in F^1$ , then the only candidate set of bad signals is  $\hat{Y} = \{y^L\}$ , meaning that  $\bar{a}^1$  is not vulnerable to a temptation. In case  $\underline{a}^1 \in F^1$ , we simply reverse the role of the signals.

☑

In these two-signal games, as in the hidden information games, short-run player utility depends on aspects of the long-run player strategy that is unobserved by subsequent short-run players. Proposition 4 shows that this must be the case for a game with two entry signals to be a bad-reputation game.

## 7. Multilateral Entry Games

We now consider games with multiple principals. In these “multilateral entry” games, the short-run players choose only whether to participate or exit. If any short-run player chooses to exit, that player receives the reservation payoff of 0, but play between the agent and other principals is unaffected. That is,  $A^j = \{exit, enter\}$  for each  $j > 1$ , and the unique exit profile is  $a_e^{-1} \equiv (exit, \dots, exit)$ . The payoff of the short-run players who enter depends only on the action of the principal, and not on how many other short-run players chose to enter; to simplify notation we denote this “entry payoff” as  $u^j(a^1)$ . If all principals exit, the long-run player’s payoff is 0; if  $m$  of them choose to enter, the long-run player’s payoff is  $u^1(a^1, m)$ . We assume that the agent cannot be forced to participate, so that there exists an action  $a^1$  such that for all  $m$ ,  $u^1(a^1, m) \geq 0$ .

We do not require that  $u^1(a^1, m)$  is linear in  $m$ , so this class of games includes those in which the agent has the opportunity to take a costly action prior to the entry decision of the short-run players. Consider for example, a game in which the long-run player is an expert advisor, and the decision of the short-run player is whether or not to pay the long-run player for advice. One example of this is the EV example of car repairs, where the long-run player is able to determine the type of repair the car needs. Other examples include stockbrokers advising clients on portfolio choices, doctors advising patients on treatments, and the IMF advising countries on economic policies. In the EV example, the private information emerges as a consequence of the decision of the short-run player to consult the long-run player, so the advice is specific to the short-run player. In another cases, at least some part of the information is not specific to the short-run player. The advisor receives a report about the



general desirability of various actions, and then meets with each of his  $n$  short-run customers, possibly learning about their individual needs. Here the advisor receives the signal regardless of whether or not he is consulted by any particular short-run player, and he may incur costs ahead of time for doing so. That is, the long-run player's payoff may depend on his action even if the short-run players decline to participate.

Costs incurred on exit are consistent with a bad reputation game provided that conditional on exit the temptations are less costly than the friendly actions. For example, the long-run player might be a stockbroker, and the general non-client specific information might be something about general economic conditions, acquired in advance in the form of economic reports that will be presented to the client. The friendly actions in this case are to report truthfully; the bad action might be to always claim that times are good. In this case the temptation is to announce that times are bad when they are actually good, to avoid being mistaken for the type that always announces good times. If it is costly to put together a persuasive package of economic data indicating that times are bad when in fact they are good this would not be a bad reputation game. If it is more costly to put together an honest report, then it would be a candidate for a bad reputation game.

We have the following obvious extension of Proposition 5.

**Proposition 8:** *Suppose in a multilateral entry game that  $Y - Y^e = \{y^L, y^H\}$  and that  $\hat{a}^1$  strictly maximizes the probability of  $y^L$  with  $u^j(\hat{a}^1) < 0$ . If for every friendly enforceable  $a^1$  there is a  $b^1$  such that  $r(y^L | b^1) < r(y^L | a^1)$  the game is a bad reputation game.*

As an illustration, suppose that the short-run players are students, the long-run player a teacher, and the signals are teaching evaluations. (This model could apply equally well to the decision to attend a particular college, graduate school, or take a particular job.) Each period, each short-run player decides whether to enter - that is, take the class, or not. The long run player has a pair of binary choices: he can either teach well

or teach poorly, and he can either administer teaching evaluations honestly or manipulate them. The public signals are whether the evaluations (averaged over respondents) are good,  $y^H$  or poor,  $y^L$ . If the evaluations are administered honestly and the class is taught well, there is probability .9 of a good evaluation. If evaluations are administered honestly and the class is taught poorly, the probability of good evaluations is only .1. Manipulating the evaluations is certain to lead to a good evaluation. All players get 0 if no students decide not to take the class. For a short-run player who enters, the short run player's payoffs are +1 for good teaching and -1 for bad. Let  $m$  denote the number of students who take the class. The rational type of long-run player pays a cost of  $m$  to teach well; good evaluations are worth  $2m$ , while manipulating evaluations costs  $3m$ . Hence in the one-shot game with only the rational type, the unique sequential equilibrium is for the rational type to teach well and not manipulate the evaluations, for an expected payoff of .8.

However, when there is a small probability that the instructor is a bad type, and the instructor faces a sequence of short-run students, Proposition 7 applies. To see this, we see that teaching poorly and administering the evaluations honestly is the unfriendly action  $\hat{a}^1$ . The friendly set consists of the pure actions “teach well, administer honest evaluations” and “teach well, manipulate.” Crucially, the action “teach well, manipulate” is unenforceable: teach poorly and manipulate yields a higher stage game payoff and the same distribution over signals. So the only enforceable action in the friendly set is “teach well, administer honestly.” This admits the temptation “teach poorly, manipulate.” Here the short-run player recognizes that if the long-run player chooses not to send the signal honestly, he loses his incentive to teach well, and so there is no reason to enter

## Appendix: Proofs

**Lemma 1:** *If  $h_t$  is a positive probability history in which  $\hat{y} \in \hat{Y}$  occurs in period  $t$  and  $\mathbf{m}(h_{t-1})[\Theta(F^1)] \leq \mathbf{g}\underline{F}^1/2$  then  $\mathbf{a}_0^1(h_t) \geq (\mathbf{g}/2)f^1$  for some friendly  $f^1$ .*

*Proof:* Given  $h_{t-1}$  the short-run players' profile has positive probability on a profile that does not exit. At such profiles  $\bar{\mathbf{a}}^1(h_t) \geq \mathbf{g}f^1$  for some friendly  $f^1$ . In particular, we must have  $\mathbf{m}(h_{t-1})[a] + \mathbf{a}_0^1(h_t)(a) \geq \mathbf{g}f^1(a)$  for each action  $a$  in the support of  $f^1$ . Since  $\mathbf{m}(h_{t-1})[\Theta(F^1)] \leq \mathbf{g}\underline{F}^1/2$  we see that  $\mathbf{a}_0^1(h_t)(a) \geq (\mathbf{g}/2)f^1(a) - (\mathbf{g}/2)\underline{F}^1 \geq (\mathbf{g}/2)f^1(a)$ , where the last inequality follows from the definition of  $\underline{F}^1$ . Since this holds for each  $a$  in the support of  $f^1$ , the conclusion follows. □

We will let  $p(\cdot | h_{t-1}) = \mathbf{r}(\cdot | \bar{\mathbf{a}}^1(h_{t-1}), \mathbf{a}^{-1}(h_{t-1}))$  denote probability distributions over signals induced by the equilibrium strategies at history  $h_{t-1}$ ; similarly

$$p(\cdot | \hat{\Theta}, h_{t-1}) = \sum_{q \in \hat{\Theta}} \mathbf{m}(h_t)[q] \mathbf{r}(\cdot | \mathbf{a}^1(q), \mathbf{a}^{-1}(h_{t-1}))$$

denotes the equilibrium distribution on signals conditional on  $q$  being in the set  $\hat{\Theta}$  of types that are committed to actions in  $N^1$ . The probability distribution  $A^{-1}$  induced by a mixed profile  $\mathbf{a}^{-1}$  can be written as a convex combination of the component distributions  $\mathbf{a}_{-e}^{-1}$ , which has support entirely in  $A^{-1} - E^{-1}$ , and  $\mathbf{a}_e^{-1}$ , which has support entirely in  $E^{-1}$ . Then  $I\mathbf{a}_{-e}^{-1} + (1 - I)\mathbf{a}_e^{-1}$  induces the same distribution over  $A^{-1}$  as  $\mathbf{a}^{-1}(h_t)$ , where  $\mathbf{a}^{-1}(h_t) \notin \text{conhull } E^{-1}$ ,  $I > 0$ . Although,  $\mathbf{a}_{-e}^{-1}$  and  $\mathbf{a}_e^{-1}$  need not correspond to mixed strategy profiles (since they can have correlation), we may still write  $u^1(\mathbf{a}^1, \mathbf{a}_e^{-1}), v^1(h_t, \mathbf{a}^1, \mathbf{a}_e^{-1}), \mathbf{r}(\cdot | \mathbf{a}^1, \mathbf{a}_e^{-1})$  and so forth for the expected values of  $u^1(\mathbf{a}^1, \mathbf{a}^{-1}), v^1(h_t, \mathbf{a}^1, \mathbf{a}^{-1}), \mathbf{r}(\cdot | \mathbf{a}^1, \mathbf{a}^{-1})$  with respect to the weights  $\mathbf{a}_e^{-1}$ ,

and similarly for  $\mathbf{a}_e^{-1}$ . With that in mind, let  $p(\cdot | \text{entry}, h_{t-1}) = \mathbf{r}(\cdot | \bar{\mathbf{a}}^1(h_{t-1}), \mathbf{a}_e^{-1}(h_{t-1}))$  be the distribution of signals after history  $h_{t-1}$ , given that the realization of the short-run players' (equilibrium) action is an entry profile, and let  $p(\cdot | \mathbf{a}^1, \text{entry}, h_{t-1}) = \mathbf{r}(\cdot | \mathbf{a}^1, \mathbf{a}_e^{-1}(h_{t-1}))$ .

**Lemma 2:** *In a bad reputation game, if  $h_t$  is a positive probability history with respect to a Nash equilibrium, and the signals in  $h_t$  all lie in  $Y^e \cup \hat{Y}$*

a) At most  $k^* = k_0 + \log(\mathbf{m}(0)[\hat{\Theta}]) / \log\left(\mathbf{y} + (1 - \mathbf{y}) \frac{1}{r}\right)$  of the signals are in  $\hat{Y}$ .

b) If the commitment size is  $((\mathbf{g}F^1 / 2)^{1+h}, \mathbf{h})$  then  $\mathbf{m}(h_t)[\Theta(F^1)] \leq \mathbf{g}F^1 / 2$ .

*Proof:* First observe that if  $\mathbf{m}(h_{t-1})[\hat{\Theta}] \geq \mathbf{y}$ , then the short-run players must exit in period  $t$ , so  $\mathbf{m}(h_t) = \mathbf{m}(h_{t-1})$ . Suppose that  $h_t$  is a positive probability history in which  $\hat{y}$  occurs in period  $t$ . From Bayes' rule

$$\mathbf{m}(h_t)[\hat{\Theta}] = \frac{p(\hat{y} | \hat{\Theta}, h_{t-1})\mathbf{m}(h_{t-1})[\hat{\Theta}]}{p(\hat{y} | h_{t-1})}$$

Since  $\hat{y}$  has positive probability at time  $t$  conditional on  $h_{t-1}$ , it must be that  $\mathbf{a}^{-1}(h_{t-1})$  has positive probability of entry. It follows that  $\bar{\mathbf{a}}^1(h_{t-1})$  puts weight less than  $\mathbf{y}$  on  $N^1$ , and that for any  $\mathbf{a}^1 \notin N^1$

$$\frac{p(\hat{y} | \hat{\Theta}, h_{t-1})}{p(\hat{y} | \mathbf{a}^1, h_{t-1})} \geq r$$

Consequently

$$\mathbf{m}(h_t)[\hat{\Theta}] \geq \frac{\mathbf{m}(h_{t-1})[\hat{\Theta}]}{\mathbf{y} + (1 - \mathbf{y})\frac{1}{r}}.$$

Since signals in  $Y^e$  convey no information about the long-run player's type, it follows that if all signals lie in  $Y^e \cup \hat{Y}$ , and signals in  $\hat{Y}$  occur  $k$  times, then

$$\mathbf{m}(h_t)[\hat{\Theta}] \geq \left( \frac{1}{\mathbf{y} + (1 - \mathbf{y})\frac{1}{r}} \right)^k \mathbf{m}(0)[\hat{\Theta}]$$

Hence if

$$k \geq -\frac{\log(\mathbf{y}) - \log(\mathbf{m}(0)[\hat{\Theta}])}{\log\left(\mathbf{y} + (1 - \mathbf{y})\frac{1}{r}\right)}$$

then  $\mathbf{m}(h_t)[\hat{\Theta}] \geq \mathbf{y}$ , so in all subsequent periods the signal must be an exit signal. (Recall that in a bad reputation game,  $r > 1$ ; this implies that the denominator above is not zero.) Again because  $r > 1$ , it is sufficient that

$$k \geq \frac{-\log(\mathbf{y}) + \log(\mathbf{m}(0)[\hat{\Theta}])}{\log\left(\mathbf{y} + (1 - \mathbf{y})\frac{1}{r}\right)}$$

which is the condition in part *a*).

We now turn to part *b*). For any history  $h$  on the equilibrium path at which entry occurs with positive probability, we must have  $\bar{\mathbf{a}}^1(h) \geq \mathbf{g}f^1$  for some friendly  $f^1$ . By assumption every enforceable friendly action is vulnerable to temptation, so that conditional on entry, the total probability of each bad signal  $\hat{y} \in \hat{Y}$  is at least  $\mathbf{jg}$  at such a history  $h$ . In particular, consider any history  $h_{t-1}$  after which a bad signal  $\hat{y}$  actually occurs. Since bad signals are entry signals, entry must have had positive probability at  $h_{t-1}$ , and hence conditional on entry,  $\hat{y}$  had total probability at least  $\mathbf{jg}$ . Bayes' rule then implies

$\mathbf{m}(h_t)[\mathbf{q}] \leq (1/j \mathbf{g})\mathbf{m}(h_{t-1})[\mathbf{q}]$  for any type  $\mathbf{q}$ . Thus  $\mathbf{m}(h_t)[\Theta(F^1)] \leq (1/j \mathbf{g})\mathbf{m}(h_{t-1})[\Theta(F^1)]$  for each  $t$  in which a bad signal occurs. In particular

$$\mathbf{m}(h_t)[\Theta(F^1)] \leq (1/j \mathbf{g})^k \mathbf{m}(0)[\Theta(F^1)] \quad (1.1)$$

where  $k$  is the number of bad signals in  $h_t$ .

Bayes' rule gives us the following inequalities

$$\mathbf{m}(h_t)[\hat{\Theta}] \geq \frac{\mathbf{m}(h_{t-1})[\hat{\Theta}] \min_{n^1 \in N^1} p(\hat{y} | n^1, \text{entry}, h_{t-1})}{p(\hat{y} | \text{entry}, h_{t-1})}$$

and

$$\begin{aligned} \mathbf{m}(h_t)[\Theta(F^1)] &\leq \frac{\mathbf{m}(h_{t-1})[\Theta(F^1)] \max_{a^1 \in F^1} p(\hat{y} | a^1, \text{entry}, h_{t-1})}{p(\hat{y} | \text{entry}, h_{t-1})} \\ &\leq \frac{\mathbf{m}(h_{t-1})[\Theta(F^1)] \max_{a^1 \notin N^1} p(\hat{y} | a^1, \text{entry}, h_{t-1})}{p(\hat{y} | \text{entry}, h_{t-1})} \end{aligned}$$

where the last inequality comes from the assumption that  $F^1$  and  $N^1$  are orthogonal.

Divide the first inequality by the second and apply the definition of  $r$  to get

$$\frac{\mathbf{m}(h_t)[\hat{\Theta}]}{\mathbf{m}(h_t)[\Theta(F^1)]} \geq r \frac{\mathbf{m}(h_{t-1})[\hat{\Theta}]}{\mathbf{m}(h_{t-1})[\Theta(F^1)]}$$

for each  $t \in \{1, \dots, t\}$  such that  $y_t \in \hat{Y}$ . Since there are  $k$  such  $t$ , we obtain

$$\frac{\mathbf{m}(h_t)[\hat{\Theta}]}{\mathbf{m}(h_t)[\Theta(F^1)]} \geq r^k \frac{\mathbf{m}(0)[\hat{\Theta}]}{\mathbf{m}(0)[\Theta(F^1)]}. \quad (1.2)$$

Finally, we define  $\mathbf{k}$  by the following equation

$$\left( \frac{\mathbf{m}(\mathbf{0})[\hat{\Theta}]}{\mathbf{m}(\mathbf{0})[\Theta(F^1)]} \right) r^k = \left( \frac{\underline{\mathbf{g}}}{2} \underline{F} \right)^{-1}. \quad (1.3)$$

Our commitment size assumption<sup>15</sup> is that

$$(1.4)$$

$$\begin{aligned} \mathbf{m}(\mathbf{0})[\Theta(F^1)] &\leq \left( \frac{\underline{\mathbf{g}}}{2} \underline{F} \right)^{(1-\log(\underline{\mathbf{g}}j)/\log r)} \left[ \frac{\mathbf{m}(\mathbf{0})[\hat{\Theta}]}{\mathbf{m}(\mathbf{0})[\Theta(F^1)]} \right]^{\left( \log\left(\frac{1}{\underline{\mathbf{g}}j}\right)/\log(r) \right)} \\ &= \left( \frac{\underline{\mathbf{g}}}{2} \underline{F} \right)^{\left( \frac{1}{\underline{\mathbf{g}}j} \right)^{\frac{\log\left(\frac{\mathbf{m}(\mathbf{0})[\hat{\Theta}]}{\mathbf{m}(\mathbf{0})[\Theta(F^1)]} + \log\left(\frac{\underline{\mathbf{g}}}{2} \underline{F}\right)}{\log r}}} \right)} \\ &= \left( \frac{\underline{\mathbf{g}}}{2} \underline{F} \right)^{\left( \frac{1}{\underline{\mathbf{g}}j} \right)^{-k}} \end{aligned}$$

With these preliminaries in hand, we can conclude the proof. First suppose that

$$\frac{\mathbf{m}(h_t)[\hat{\Theta}]}{\mathbf{m}(h_t)[\Theta(F^1)]} > \left( \frac{\underline{\mathbf{g}}}{2} \underline{F} \right)^{-1}.$$

Then it follows immediately that

$$\mathbf{m}(h_t)[\Theta(F^1)] < \frac{\underline{\mathbf{g}}}{2} \underline{F}$$

and we are done. On the other hand, the opposite inequality implies by (1.2) and (1.3) that  $k \leq \mathbf{k}$ . Consequently, by (1.1) and (1.4) we have

$$\mathbf{m}(h_t)[\Theta(F^1)] \leq \frac{\underline{\mathbf{g}}}{2} \underline{F}.$$

□

Recall that

$$\bar{u}^1(y, \tilde{\mathbf{r}}) = \begin{cases} \left( 1 + \frac{1}{\tilde{\mathbf{r}}} \right) U^1 & y \in \hat{Y} \\ 0 & \text{otherwise} \end{cases}$$

<sup>15</sup> The first line in this derivation follows from the definition of  $\mathbf{h}$ . The second line is an

application of the rule:  $a^{\frac{\log b}{c}} = e^{\frac{\log b \log a}{c}} = \left( e^{\log b} \right)^{\frac{\log a}{c}} = b^{\frac{\log a}{c}}$ .

$$\bar{\mathbf{d}}^1(\mathbf{y}, \tilde{\mathbf{r}}) = \begin{cases} \frac{\mathbf{d}}{\tilde{\mathbf{r}}} + 1 & \mathbf{y} \in \hat{Y} \\ \mathbf{d} & \text{otherwise} \end{cases}$$

and  $Y(\mathbf{h}_t) = \{\mathbf{y} \in Y^e \cup \hat{Y} \mid \mathbf{r}(\mathbf{y} \mid \bar{\mathbf{a}}^1(\mathbf{h}_t), \mathbf{a}^{-1}(\mathbf{h}_t)) > 0\}$ .

**Lemma 3:** *In a participation game if  $\mathbf{a}^{-1}(\mathbf{h}_t) \in \text{conhull } E^{-1}$  or  $\mathbf{a}^{-1}(\mathbf{h}_t) \notin \text{conhull } E^{-1}$  and  $\mathbf{a}_0^1(\mathbf{h}_t) \geq c\mathbf{f}^1$  for some  $c > 0$  and vulnerable friendly action  $f^1$  with temptation bounds  $\underline{\mathbf{r}}, \tilde{\mathbf{r}}$  then*

$$v^1(\mathbf{h}_t) \leq \max_{\mathbf{y} \in Y(\mathbf{h}_t)} (1 - \mathbf{d})\bar{u}^1(\mathbf{y}, \tilde{\mathbf{r}}) + \bar{\mathbf{d}}(\mathbf{y}, \tilde{\mathbf{r}})v^1(\mathbf{h}_t, \mathbf{y}).$$

Proof: We need to calculate the long-run player payoff separately as a function of whether the short-run players exit or not. Using the decomposition of the short-run players' profile that we introduced before the proof of lemma 2, we write

$$v^1(\mathbf{h}_t) = \mathbf{I}v^1(\mathbf{h}_t, f^1, \mathbf{a}_e^{-1}) + (1 - \mathbf{I})v^1(\mathbf{h}_t, f^1, \mathbf{a}_{-e}^{-1}).$$

First assume  $\mathbf{I} \geq 0$  and consider the value  $v^1(\mathbf{h}_t, f^1, \mathbf{a}_e^{-1})$  conditional on exit. By exit minmax, this value is no more than  $\max_{\mathbf{y} \in Y(\mathbf{h}_t)} \mathbf{d}v^1(\mathbf{h}_t, \mathbf{y})$  and thus from the definition of  $\bar{u}^1$  we can conclude that

$$v^1(\mathbf{h}_t, f^1, \mathbf{a}_e^{-1}) \leq \max_{\mathbf{y} \in Y(\mathbf{h}_t)} (1 - \mathbf{d})\bar{u}^1(\mathbf{y}, \tilde{\mathbf{r}}) + \bar{\mathbf{d}}(\mathbf{y}, \tilde{\mathbf{r}})v^1(\mathbf{h}_t, \mathbf{y}).$$

If  $\mathbf{a}^{-1}(\mathbf{h}_t) \in \text{conhull } E^{-1}$  then  $\mathbf{I} = 1$  and we are done with the first case in the statement.

Now consider now the second case in the claim of the lemma,  $\mathbf{a}^{-1}(\mathbf{h}_t) \notin \text{conhull } E^{-1}$ , so that entry has positive probability and  $\mathbf{I} < 1$ . In this case,  $f^1$  must be enforceable, since otherwise the rational type could do better by changing from his equilibrium action to one that replaces the probability on the support of  $f^1$  with probability on the improving and observationally equivalent mixed strategy that defeats it,  $f^1$  and get a strictly higher utility. Every enforceable friendly action is



vulnerable to temptation, so let  $b^1$  be a temptation for  $f^1$ . Since  $f^1$  is played in equilibrium, it earns at least as much as  $b^1$ , so that

$$\begin{aligned} & v^1(h_t, f^1, \mathbf{a}^{-1}(h_t)) - v^1(h_t, b^1, \mathbf{a}^{-1}(h_t)) \\ &= \mathbf{I} [v^1(h_t, f^1, \mathbf{a}_e^{-1}) - v^1(h_t, b^1, \mathbf{a}_e^{-1})] + (1 - \mathbf{I}) [v^1(h_t, f^1, \mathbf{a}_{-e}^{-1}) - v^1(h_t, b^1, \mathbf{a}_{-e}^{-1})] \\ &\geq 0 \end{aligned}$$

The expression  $\mathbf{I} [v^1(h_t, f^1, \mathbf{a}_e^{-1}) - v^1(h_t, b^1, \mathbf{a}_e^{-1})]$  is non-positive because given exit,  $b^1$  and  $f^1$  induce the same distribution over signals, and hence earn identical continuation values, and by the definition of a temptation,  $u^1(b^1, \mathbf{a}_e^{-1}) \geq u^1(f^1, \mathbf{a}_e^{-1})$ , so that  $b^1$  does at least as well in the current period..

Thus,  $v^1(h_t, f^1, \mathbf{a}_e^{-1}) - v^1(h_t, b^1, \mathbf{a}_e^{-1}) \geq 0$ . Expanding this inequality and using the fact that  $u^1(f^1, \mathbf{a}_{-e}^{-1}) - u^1(b^1, \mathbf{a}_{-e}^{-1}) \leq U^1$ ,

$$(1.5) \quad \begin{aligned} & (1 - \mathbf{d})U^1 + \mathbf{d} \sum_{\hat{y} \in \hat{Y}} [\mathbf{r}(\hat{y} | f^1, \mathbf{a}_e^{-1}) - \mathbf{r}(\hat{y} | b^1, \mathbf{a}_e^{-1})] v^1(h_t, \hat{y}) \geq \\ & \mathbf{d} \left[ \sum_{y \in Y/\hat{Y}} [\mathbf{r}(y | b^1, \mathbf{a}_e^{-1}) - \mathbf{r}(y | f^1, \mathbf{a}_e^{-1})] v^1(h_t, y) \right] \end{aligned}$$

Define

$$v^1(h_t, \hat{Y}) = \max_{y \in \hat{Y}} v^1(h_t, y) \geq 0.$$

The inequality holds because continuation values for histories on the equilibrium path of a Nash equilibrium must exceed the minmax value, which we have normalized to zero. We will use this fact repeatedly in the remainder of the proof. By the definition of a temptation,  $\mathbf{r}(\hat{y} | b^1, \mathbf{a}_e^{-1}) < \mathbf{r}(\hat{y} | f^1, \mathbf{a}_e^{-1})$  for each  $\hat{y}$ . Thus, inequality (1.5) can be reduced to

$$\begin{aligned}
(1 - \mathbf{d})U^1 + \mathbf{d}v^1(h_t, \widehat{Y}) &\geq \mathbf{d} \left\{ \sum_{y \in Y/\widehat{Y}} [\mathbf{r}(y | b^1, \mathbf{a}_e^{-1}) - \mathbf{r}(y | f^1, \mathbf{a}_e^{-1})] v^1(h_t, y) \right\} \\
&\geq \widehat{\mathbf{d}} \sum_{y \in Y/\widehat{Y}} \mathbf{r}(y | f^1, \mathbf{a}_e^{-1}) v^1(h_t, y)
\end{aligned}$$

where the second inequality uses part 2 of the definition of a temptation. We can now expand the definition of  $v^1(h_t, f^1, \mathbf{a}_e^{-1})$  and bound it as follows.

$$\begin{aligned}
&(1 - \mathbf{d})u^1(f^1, \mathbf{a}_e^{-1}) + \mathbf{d} \left[ \sum_{\widehat{y} \in \widehat{Y}} \mathbf{r}(\widehat{y} | f^1, \mathbf{a}_e^{-1}) v^1(h_t, \widehat{y}) + \sum_{y \notin \widehat{Y}} \mathbf{r}(y | f^1, \mathbf{a}_e^{-1}) v^1(h_t, y) \right] \\
&\leq (1 - \mathbf{d})U^1 + \mathbf{d}v^1(h_t, \widehat{Y}) + \frac{(1 - \mathbf{d})U^1}{\widehat{\mathbf{r}}} + \frac{\mathbf{d}}{\widehat{\mathbf{r}}} v^1(h_t, \widehat{Y}) \\
&= \max_{\widehat{y} \in \widehat{Y}} (1 - \mathbf{d})\bar{u}^1(\widehat{y}, \widehat{\mathbf{r}}) + \left( \frac{\mathbf{d}}{\widehat{\mathbf{r}}} + 1 \right) v^1(h_t, \widehat{y}) \\
&\leq \max_{y \in Y(h_t)} (1 - \mathbf{d})\bar{u}^1(y, \widehat{\mathbf{r}}) + \bar{\mathbf{d}}^1(y, \widehat{\mathbf{r}}) v^1(h_t, y)
\end{aligned}$$

where the last inequality follows because when  $I < 1$ ,  $\widehat{Y} \subseteq Y(h_t)$  (Recall that the vulnerability of  $f^1$  implies  $\mathbf{r}(\widehat{y} | f^1, \mathbf{a}_e^{-1}) > 0$  for each  $\widehat{y} \in \widehat{Y}$  and each entry profile  $\mathbf{a}^{-1}$ .)

This concludes the proof because if  $I = 0$  then  $v^1(h_t) = v^1(h_t, f^1, \mathbf{a}_e^{-1})$  and if  $I \in (0, 1)$  then  $v^1(h_t) \leq \max \{ v^1(h_t, f^1, \mathbf{a}_e^{-1}), v^1(h_t, f^1, \mathbf{a}_e^{-1}) \}$ .

□

**Lemma 4:** *In a participation game, if  $\mathbf{a}^{-1}(h_t) \in \text{conhull } E^{-1}$  or  $\mathbf{a}^{-1}(h_t) \notin \text{conhull } E^{-1}$  and  $\mathbf{a}_0^1(h_t) \geq c f^1$  for some  $c > 0$  and friendly action  $f^1$  that is vulnerable to a strong temptation size  $\underline{\mathbf{r}}$ , then*

$$v^1(h_t) \leq \max_{y \in Y(h_t)} (1 - \mathbf{d})\bar{u}^1(y, \underline{\mathbf{r}}) + \bar{\mathbf{d}}^1(y, \underline{\mathbf{r}}) \mathbf{d}v^1(h_t, y)$$

$$\text{where } \bar{u}^1(y, \underline{\mathbf{r}}) = \begin{cases} \left( 1 + \frac{1}{|\widehat{Y}| \underline{\mathbf{r}}} \right) U^1 & \text{if } y \in \widehat{Y} \\ \widehat{u}^1 & \text{otherwise} \end{cases}$$

$$\text{and } \bar{\mathbf{d}}^1(\mathbf{y}, \mathbf{r}) = \begin{cases} \mathbf{d} \left( 1 + \frac{1}{|\widehat{Y}|} \mathbf{r} \right) & y \in \widehat{Y} \\ \mathbf{d} & \text{otherwise} \end{cases}$$

*Proof:* The proof of Lemma 4 is similar to that of lemma 3; indeed the case  $\mathbf{l} = \mathbf{0}$  is identical, and when  $\mathbf{l} > \mathbf{0}$  the derivation of equation (1.5) is unchanged

$$(1 - \mathbf{d})U^1 + \mathbf{d} \sum_{\widehat{y} \in \widehat{Y}} [\mathbf{r}(\widehat{y} | f^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(\widehat{y} | b^1, \mathbf{a}_{-e}^{-1})] v^1(\mathbf{h}_t, \widehat{y}) \geq \mathbf{d} \left[ \sum_{y \in Y/\widehat{Y}} [\mathbf{r}(y | b^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(y | f^1, \mathbf{a}_{-e}^{-1})] v^1(\mathbf{h}_t, y) \right] \quad (1.5)$$

Since the good signals are changed proportionately by the temptation, it

follows that  $\mathbf{r}(y | b^1, \mathbf{a}_{-e}^{-1}) = \frac{\mathbf{r}(Y/\widehat{Y} | b^1, \mathbf{a}_{-e}^{-1})}{\mathbf{r}(Y/\widehat{Y} | f^1, \mathbf{a}_{-e}^{-1})} \mathbf{r}(y | f^1, \mathbf{a}_{-e}^{-1})$  for each

$y \in Y/\widehat{Y}$ . Thus,

$$\begin{aligned} & \sum_{y \in Y/\widehat{Y}} [\mathbf{r}(y | b^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(y | f^1, \mathbf{a}_{-e}^{-1})] v^1(\mathbf{h}_t, y) \\ &= \sum_{y \in Y/\widehat{Y}} \left[ \left( \frac{\mathbf{r}(Y/\widehat{Y} | b^1, \mathbf{a}_{-e}^{-1})}{\mathbf{r}(Y/\widehat{Y} | f^1, \mathbf{a}_{-e}^{-1})} - 1 \right) \mathbf{r}(y | f^1, \mathbf{a}_{-e}^{-1}) \right] v^1(\mathbf{h}_t, y) \\ &= \left( \frac{\mathbf{r}(Y/\widehat{Y} | b^1, \mathbf{a}_{-e}^{-1})}{\mathbf{r}(Y/\widehat{Y} | f^1, \mathbf{a}_{-e}^{-1})} - 1 \right) \sum_{y \in Y/\widehat{Y}} \mathbf{r}(y | f^1, \mathbf{a}_{-e}^{-1}) v^1(\mathbf{h}_t, y) \\ &= \left( \frac{\mathbf{r}(Y/\widehat{Y} | b^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(Y/\widehat{Y} | f^1, \mathbf{a}_{-e}^{-1})}{\mathbf{r}(Y/\widehat{Y} | f^1, \mathbf{a}_{-e}^{-1})} \right) v^1(\mathbf{h}_t, Y/\widehat{Y}) \end{aligned}$$

where  $v^1(\mathbf{h}_t, Y/\widehat{Y})$  is the expected continuation value after playing  $f^1$  and observing a signal in  $Y/\widehat{Y}$ . Substituting into (1.5),

$$\begin{aligned} & (1 - \mathbf{d})U^1 + \mathbf{d} \sum_{\widehat{y} \in \widehat{Y}} [\mathbf{r}(\widehat{y} | f^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(\widehat{y} | b^1, \mathbf{a}_{-e}^{-1})] v^1(\mathbf{h}_t, \widehat{y}) \\ & \geq \mathbf{d} \left[ \frac{\mathbf{r}(Y/\widehat{Y} | b^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(Y/\widehat{Y} | f^1, \mathbf{a}_{-e}^{-1})}{\mathbf{r}(Y/\widehat{Y} | f^1, \mathbf{a}_{-e}^{-1})} \right] v^1(\mathbf{h}_t, Y/\widehat{Y}) \\ & \geq \mathbf{d} [\mathbf{r}(Y/\widehat{Y} | b^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(Y/\widehat{Y} | f^1, \mathbf{a}_{-e}^{-1})] v^1(\mathbf{h}_t, Y/\widehat{Y}) \end{aligned}$$

Set

$$v^1(h_t, \widehat{Y}) = \max_{y \in \widehat{Y}} v^1(h_t, y).$$

From the fact that  $b^1$  reduces the probability of every bad signal by a positive amount,

$$\begin{aligned} & [\mathbf{r}(\widehat{y} \mid f^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(\widehat{y} \mid b^1, \mathbf{a}_{-e}^{-1})] v^1(h_t, \widehat{y}) \\ & \leq [\mathbf{r}(\widehat{Y} \mid f^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(\widehat{Y} \mid b^1, \mathbf{a}_{-e}^{-1})] v^1(h_t, \widehat{Y}) \end{aligned}$$

for each  $\widehat{y} \in \widehat{Y}$ . Thus,

$$(1.6) \quad \begin{aligned} & (1 - d)U^1 + d[\mathbf{r}(\widehat{Y} \mid f^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(\widehat{Y} \mid b^1, \mathbf{a}_{-e}^{-1})] v^1(h_t, \widehat{Y}) \geq \\ & d[\mathbf{r}(Y / \widehat{Y} \mid b^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(Y / \widehat{Y} \mid f^1, \mathbf{a}_{-e}^{-1})] v^1(h_t, Y / \widehat{Y}) \end{aligned}$$

which implies

$$\begin{aligned} v^1(h_t, Y / \widehat{Y}) & \leq \frac{(1 - d)U^1 + d[\mathbf{r}(\widehat{Y} \mid f^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(\widehat{Y} \mid b^1, \mathbf{a}_{-e}^{-1})] v^1(h_t, \widehat{Y})}{d[\mathbf{r}(Y / \widehat{Y} \mid b^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(Y / \widehat{Y} \mid f^1, \mathbf{a}_{-e}^{-1})]} \\ & \leq \frac{(1 - d)U^1 + d v^1(h_t, \widehat{Y})}{d[\mathbf{r}(Y / \widehat{Y} \mid b^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(Y / \widehat{Y} \mid f^1, \mathbf{a}_{-e}^{-1})]} \end{aligned}$$

where the second line uses the fact that  $v^1(h_t, \widehat{Y}) \geq 0$ .

Because  $b^1$  lowers the probability of all bad signals by at least  $\underline{\mathbf{r}}$ , it raises the total probability of the remaining signals by at least  $|\widehat{Y}| \underline{\mathbf{r}}$ , i.e.  $[\mathbf{r}(Y / \widehat{Y} \mid b^1, \mathbf{a}_{-e}^{-1}) - \mathbf{r}(Y / \widehat{Y} \mid f^1, \mathbf{a}_{-e}^{-1})] \geq |\widehat{Y}| \underline{\mathbf{r}}$ . This and the fact that the numerator on the right hand side of the previous inequality is non-negative gives

$$v^1(Y / \widehat{Y}) \leq \frac{(1 - d)U^1 + d v^1(h_t, \widehat{Y})}{d |\widehat{Y}| \underline{\mathbf{r}}}$$

Finally, we conclude:

$$\begin{aligned}
v_{\sim e}^1 &\equiv (1 - \mathbf{d})u^1(f^1, \mathbf{a}_{\sim e}^{-1}) + \mathbf{d}\left(\sum_{\hat{y} \in \hat{Y}} \mathbf{r}(\hat{y} \mid f^1, \mathbf{a}_{\sim e}^{-1})v^1(h_t, \hat{y}) + \mathbf{r}(Y/\hat{Y} \mid f^1, \mathbf{a}_{\sim e}^{-1})v^1(h_t, Y/\hat{Y})\right) \\
&\leq (1 - \mathbf{d})u^1(f^1, \mathbf{a}_{\sim e}^{-1}) + \mathbf{d}\left(\mathbf{r}(\hat{Y} \mid f^1, \mathbf{a}_{\sim e}^{-1})v^1(h_t, \hat{Y}) + \mathbf{r}(Y/\hat{Y} \mid f^1, \mathbf{a}_{\sim e}^{-1})v^1(h_t, Y/\hat{Y})\right) \\
&\leq (1 - \mathbf{d})u^1(f^1, \mathbf{a}_{\sim e}^{-1}) + \mathbf{d}\left(v^1(h_t, \hat{Y}) + v^1(h_t, Y/\hat{Y})\right) \\
&\leq (1 - \mathbf{d})U^1 + \mathbf{d}v^1(h_t, \hat{Y}) + \mathbf{d}\left(\frac{(1 - \mathbf{d})U^1 + \mathbf{d}v^1(h_t, \hat{Y})}{\mathbf{d}|\hat{Y}|_{\underline{\mathbf{r}}}}\right) \\
&= (1 - \mathbf{d})\left(1 + \frac{1}{|\hat{Y}|_{\underline{\mathbf{r}}}}\right)U^1 + \mathbf{d}\left(1 + \frac{1}{|\hat{Y}|_{\underline{\mathbf{r}}}}\right)v^1(h_t, \hat{Y}) \\
&= \max_{y \in \hat{Y}} (1 - \mathbf{d})\bar{u}^1(y, \underline{\mathbf{r}}) + \bar{\mathbf{d}}(y, \underline{\mathbf{r}})v^1(h_t, y)
\end{aligned}$$

The conclusion of the proof is now identical to that of Lemma 3: if  $\mathbf{I} = 0$  then  $v^1(h_t) = v^1(h_t, f^1, \mathbf{a}_{\sim e}^{-1})$  and if  $\mathbf{I} \in (0, 1)$  then  $v^1(h_t) \leq \max\{v^1(h_t, f^1, \mathbf{a}_{\sim e}^{-1}), v^1(h_t, f^1, \mathbf{a}_{\sim e}^{-1})\}$ .

☑

## References

- Ely, J. and J. Valimaki [2002] "Bad Reputation," *NAJ Economics*, 4: 2, <http://www.najecon.org/v4.htm> and forthcoming *Quarterly Journal of Economics*
- Fudenberg, D. and D. Kreps [1987] "Reputation and Simultaneous Opponents" *Review of Economic Studies*, 54: 541-568
- Fudenberg, D., D. Kreps, and E. Maskin [1990] "Repeated Games with Long-run and Short-run Players," *Review of Economic Studies*, 57, 555-573.
- Fudenberg, D. and D. K. Levine [1994] "Efficiency and Observability in Games with Long-Run and Short-Run Players," *Journal of Economic Theory*, 62 , 103-135
- Fudenberg, D. and D. K. Levine [1992] "Maintaining a Reputation when Strategies are Imperfectly Observed," *Review of Economic Studies*, 59: 561-579.
- Fudenberg, D. and D. K. Levine [1989] "Reputation and Equilibrium Selection in Games with a Single Long-Run Player" *Econometrica*, 57: 759-778.
- Fudenberg, D., E. Maskin, and D.K. Levine [1994] "The Folk Theorem in Repeated Games with Imperfect Public Information," *Econometrica* 62, 997-1039.
- Kreps, D. and R. Wilson [1982] "Reputation and Imperfect Information," *Journal of Economic Theory*, 27:253-279.
- Mailath, G. and L. Samuelson [1998] "Your Reputation is Who You're Not, Not Who You'd like to Be," mimeo.
- Milgrom, P. and J. Roberts [1982] "Predation, Reputation, and Entry Deterrence," *Journal of Economic Theory*, 27:280-213.
- Sorin, S. [1999] Merging, Reputation, and Repeated Games with Incomplete Information," *Games and Economic Behavior* 29, 274-308