

Supplementary Appendix to: When is Reputation Bad?

Jeffrey Ely
Drew Fudenberg
David K. Levine¹

November 22, 2005

In Ely, Fudenberg and Levine [2005], hereafter EFL, we defined an action to vulnerable to a temptation if conditional on participation by the short-run players, the temptation lowers the probability of all “bad” signals, and increases the probability of all other signals. EFL’s first bad reputation result requires an exit minmax condition. If the temptation satisfies the stronger property that the relative probability of the other signals remains constant, then the assumption of exit minmax can be weakened. We define a strong temptation:

Definition 3S: An action a is vulnerable to a strong temptation relative to a set of signals \hat{Y} if there exists a number $\underline{\rho} > 0$ and an action d such that

- 1) If $b \notin E$, $\hat{y} \in \hat{Y}$ then $\rho(\hat{y}|d,b) \leq \rho(\hat{y}|a,b) - \underline{\rho}$
- 2) If $b \notin E$ and $y, y' \notin \hat{Y} \cup Y^E$ then $\frac{\rho(y|d,b)}{\rho(y'|d,b)} = \frac{\rho(y|a,b)}{\rho(y'|a,b)}$.
- 3) For all $e \in E$, $u^L(d,e) \geq u^L(a,e)$.

The action d is called a strong temptation.

This condition lets us prove an analog of lemma 3 in EFL which we prove here.

¹ Departments of Economics, Northwestern University, Harvard University and UCLA.

Lemma 3S: *In a participation game, if $\beta(h_t)\{E\} = 1$ or $\beta(h_t)\{E\} < 1$ and $\alpha_0(h_t)(f) > 0$ for some friendly action f that is vulnerable to a strong temptation of size $\underline{\rho}$, then*

$$v(h_t) \leq \max_{y \in Y(h_t)} (1 - \delta) \bar{u}(y, \tilde{\rho}) + \delta \bar{\delta}(y, \tilde{\rho}) v(h_t, y)$$

where

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

$$\text{and } \bar{\delta}(y, \rho) = \begin{cases} \delta \left(1 + \frac{1}{|\hat{Y}| \tilde{\rho}}\right) & y \in \hat{Y} \\ \delta & \text{otherwise} \end{cases}$$

Proof: The derivation of equation (1.5) is unchanged

$$\begin{aligned} & (1 - \delta)U^L + \delta \sum_{\hat{y} \in \hat{Y}} [\rho(\hat{y} | f, \beta_{-E}) - \rho(\hat{y} | d, \beta)] v(h_t, \hat{y}) \\ & \geq \delta \left[\sum_{y \in Y \setminus \hat{Y}} [\rho(y | d, \beta_{-E}) - \rho(y | f, \beta_{-E})] v(h_t, y) \right] \end{aligned} \quad (1.5)$$

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

$$\rho(y | b, \beta_{-E}) = \frac{\rho(Y \setminus \hat{Y} | d, \beta_{-E})}{\rho(Y \setminus \hat{Y} | f, \beta_{-E})} \rho(y | f, \beta_{-E})$$

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

$$\begin{aligned} & \sum_{y \in Y \setminus \hat{Y}} [\rho(y | d, \beta_{-E}) - \rho(y | f, \beta_{-E})] v(h_t, y) \\ & = \sum_{y \in Y \setminus \hat{Y}} \left[\left(\frac{\rho(Y \setminus \hat{Y} | d, \beta_{-E})}{\rho(Y \setminus \hat{Y} | f, \beta_{-E})} - 1 \right) \rho(y | f, \beta_{-E}) \right] v(h_t, y) \\ & = \left(\frac{\rho(Y \setminus \hat{Y} | d, \beta_{-E})}{\rho(Y \setminus \hat{Y} | f, \beta_{-E})} - 1 \right) \sum_{y \in Y \setminus \hat{Y}} \rho(y | f, \beta_{-E}) v(h_t, y) \\ & = \left(\frac{\rho(Y \setminus \hat{Y} | d, \beta_{-E}) - \rho(Y \setminus \hat{Y} | f, \beta_{-E})}{\rho(Y \setminus \hat{Y} | f, \beta_{-E})} \right) v(h_t, Y \setminus \hat{Y}) \end{aligned}$$

where $v(h_t, Y \setminus \hat{Y})$ is the expected continuation value after playing f and observing a signal in $Y \setminus \hat{Y}$. Substituting into (1.5),

$$\begin{aligned}
& (1 - \delta)U^L + \delta \sum_{\hat{y} \in \hat{Y}} [\rho(\hat{y} | f, \beta_{-E}) - \rho(\hat{y} | d, \beta_{-E})] v(h_t, \hat{y}) \\
& \geq \delta \left[\frac{\rho(Y \setminus \hat{Y} | d, \beta_{-E}) - \rho(Y \setminus \hat{Y} | f, \beta_{-E})}{\rho(Y \setminus \hat{Y} | f, \beta_{-E})} \right] v(h_t, Y \setminus \hat{Y}) \\
& \geq \delta [\rho(Y \setminus \hat{Y} | d, \beta_{-E}) - \rho(Y \setminus \hat{Y} | f, \beta_{-E})] v(h_t, Y \setminus \hat{Y})
\end{aligned}$$

Set

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

From the fact that d reduces the probability of every bad signal by a positive amount,

$$\begin{aligned}
& [\rho(\hat{y} | f, \beta_{-E}) - \rho(\hat{y} | d, \beta_{-E})] v(h_t, \hat{y}) \\
& \leq [\rho(\hat{Y} | f, \beta_{-E}) - \rho(\hat{Y} | d, \beta_{-E})] v(h_t, \hat{Y})
\end{aligned}$$

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

$$\begin{aligned}
& (1 - \delta)U^L + \delta [\rho(\hat{Y} | f, \beta_{-E}) - \rho(\hat{Y} | d, \beta_{-E})] v(h_t, \hat{Y}) \\
& \geq \delta [\rho(Y \setminus \hat{Y} | d, \beta_{-E}) - \rho(Y \setminus \hat{Y} | f, \beta_{-E})] v(h_t, Y \setminus \hat{Y})
\end{aligned} \tag{1.6}$$

which implies

$$\begin{aligned}
v(h_t, Y \setminus \hat{Y}) & \leq \frac{(1 - \delta)U^L + \delta [\rho(\hat{Y} | f, \beta_{-E}) - \rho(\hat{Y} | d, \beta_{-E})] v(h_t, \hat{Y})}{\delta [\rho(Y \setminus \hat{Y} | d, \beta_{-E}) - \rho(Y \setminus \hat{Y} | f, \beta_{-E})]} \\
& \leq \frac{(1 - \delta)U^L + \delta v(h_t, \hat{Y})}{\delta [\rho(Y \setminus \hat{Y} | d, \beta_{-E}) - \rho(Y \setminus \hat{Y} | f, \beta_{-E})]}
\end{aligned}$$

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

Because d lowers the probability of all bad signals by at least $\underline{\rho}$, it raises the total probability of the remaining signals by at least $|\hat{Y}| \underline{\rho}$, that is,

$$[\rho(Y \setminus \hat{Y} | d, \beta_{-E}) - \rho(Y \setminus \hat{Y} | f, \beta_{-E})] \geq |\hat{Y}| \underline{\rho}.$$

This and the fact that the numerator on the right hand side of the previous inequality is non-negative gives

$$v(Y \setminus \hat{Y}) \leq \frac{(1 - \delta)U^L + \delta v(h_t, \hat{Y})}{\delta |\hat{Y}| \underline{\rho}}$$

Finally, we conclude:

$$\begin{aligned}
& (1 - \delta)u^L(f, \beta_{-E}) + \delta \left(\sum_{\hat{y} \in \hat{Y}} \rho(\hat{y} \mid f, \beta_{-E}) v(h_t, \hat{y}) + \rho(Y \setminus \hat{Y} \mid f, \beta_{-E}) v(h_t, Y \setminus \hat{Y}) \right) \\
& \leq (1 - \delta)u^L(f, \beta_{-E}) + \delta \left(\rho(\hat{Y} \mid f, \beta_{-E}) v(h_t, \hat{Y}) + \rho(Y \setminus \hat{Y} \mid f, \beta_{-E}) v(h_t, Y \setminus \hat{Y}) \right) \\
& \leq (1 - \delta)u^L(f, \beta_{-E}) + \delta \left(v(h_t, \hat{Y}) + v(h_t, Y \setminus \hat{Y}) \right) \\
& \leq (1 - \delta)U^L + \delta v(h_t, \hat{Y}) + \delta \left(\frac{(1 - \delta)U^L + \delta v(h_t, \hat{Y})}{\delta |\hat{Y}|_{\underline{\rho}}} \right) \\
& = (1 - \delta) \left(1 + \frac{1}{|\hat{Y}|_{\underline{\rho}}} \right) U^L + \delta \left(1 + \frac{1}{|\hat{Y}|_{\underline{\rho}}} \right) v(h_t, \hat{Y}) \\
& = \max_{y \in \hat{Y}} (1 - \delta) \bar{u}(y, \underline{\rho}) + \bar{\delta}(y, \underline{\rho}) v(h_t, y)
\end{aligned}$$

The conclusion of the proof is now identical to that of Lemma 3: if $\lambda = 0$ then $v(h_t) = v(h_t, f, \beta_{-E})$ and if $\lambda \in (0, 1)$ then $v(h_t) \leq \max \{ v(h_t, f, \beta_E), v(h_t, f, \beta_{-E}) \}$.

□