

# Appendix

## 1 CONDITIONAL STRATEGIES WITHOUT PUNISHMENT

The following equations define the respective individual payoff whenever an agent employs a cooperation (defection) strategy:

$$\begin{aligned} X_c &= 1 + b \frac{T_c(n-1) + 1}{n} - c \\ X_d &= 1 + b \frac{T_c(n-1)}{n} \end{aligned} \quad (1)$$

In this case, the conditional expected payoffs for other group members are:

$$\begin{aligned} \overline{X}_c &= T_c X_c + (1 - T_c) \left( X_d + \frac{b}{n} \right) \\ \overline{X}_d &= T_c \left( X_c - \frac{b}{n} \right) + (1 - T_c) X_d \end{aligned} \quad (2)$$

Here, note that the following inequalities hold:

$$\begin{aligned} \overline{X}_c &= T_c X_c + (1 - T_c)(X_c + c) \\ &= X_c + (1 - T_c)c \\ &\geq X_c \end{aligned} \quad (3)$$

$$\begin{aligned} \overline{X}_d &= T_c(X_d - c) + (1 - T_c)X_d \\ &= X_d - cT_c \\ &\leq X_d \end{aligned} \quad (4)$$

Thus, the expected utility of cooperation and defection take the form:

$$\begin{aligned} E[U_c] &= \beta \overline{X}_c + (1 - \beta) X_c \\ E[U_d] &= \alpha \overline{X}_d + (1 - \alpha) X_d \end{aligned}$$

where  $\alpha$  and  $\beta$  are restricted to  $\Omega$ :

$$\Omega = \left\{ (\alpha, \beta) : -1 \leq \alpha \leq 1, -1 \leq \beta \leq \alpha \right\}.$$

The preceding expectations can be re-written as:

$$\begin{aligned}
E[U_c] &= \beta \left[ X_c + (1 - T_c)c \right] + (1 - \beta)X_c \\
&= X_c + \beta(1 - T_c)c \\
E[U_d] &= \alpha \left[ X_d - cT_c \right] + (1 - \alpha)X_d \\
&= X_d - \alpha cT_c
\end{aligned} \tag{5}$$

Thus, cooperation evolves whenever:

$$\begin{aligned}
&E[U_c] > E[U_d] \\
\implies &\alpha T_c + \beta(1 - T_c) > 1 - \frac{b}{nc} \\
\implies &\beta > -\left( \frac{T_c}{1 - T_c} \right) \alpha + \frac{1 - \frac{b}{nc}}{1 - T_c}.
\end{aligned} \tag{6}$$

## 2 CONDITIONAL STRATEGIES WITH PUNISHMENT

The following equations define the respective individual payoff whenever an agent employs a cooperation (or defection) strategy. Now, groups are also comprised of ( $nT_p$ ) punishers who reduce the earnings of defectors by ( $p$ ) at a personal cost ( $k$ ):

$$\begin{aligned}
X_c &= 1 + b \frac{T_c(n-1) + 1}{n} - c \\
X_d &= 1 + b \frac{T_c(n-1)}{n} - pT_pT_c \\
X_{cp} &= 1 + b \frac{T_c(n-1) + 1}{n} - c - k(1 - T_c).
\end{aligned} \tag{7}$$

The conditional expected payoffs for other group members are:

$$\begin{aligned}
\overline{X}_c &= T_c \left[ (1 - T_p)X_c + T_pX_{cp} \right] + (1 - T_c) \left[ X_d + \frac{b}{n} \right] \\
\overline{X}_d &= T_c \left[ (1 - T_p) \left( X_c - \frac{b}{n} \right) + T_p \left( X_{cp} - \frac{b}{n} - k \right) \right] + (1 - T_c)X_d \\
\overline{X}_{cp} &= T_c \left[ (1 - T_p)X_c + T_pX_{cp} \right] + (1 - T_c) \left[ X_d + \frac{b}{n} - p \right].
\end{aligned} \tag{8}$$

We can re-write the equations in (8) as:

$$\begin{aligned}
\overline{X}_c &= X_c + T_c \left[ T_p(X_{cp} - X_c) \right] + (1 - T_c) \left[ c - pT_pT_c \right] \\
&= X_c + (1 - T_c) \left[ c - T_pT_c(k + p) \right] \\
&\geq X_c
\end{aligned} \tag{9}$$

where the last inequality holds if:  $T_pT_c(k + p) < c$ .

$$\begin{aligned}
\overline{X}_d &= X_d + T_c \left[ T_p(pT_pT_c - c - k - k(1 - T_c)) + (1 - T_p)(pT_pT_c - c) \right] \\
&= X_d + T_c \left[ pT_pT_c - T_pk(2 - T_c) - c \right] \\
&\leq X_d
\end{aligned} \tag{10}$$

where the last inequality holds if:  $T_p((k + p)T_c - 2k) < c$ .

$$\begin{aligned}
\overline{X}_{cp} &= X_{cp} + T_c \left[ (1 - T_p)k(1 - T_c) \right] + (1 - T_c) \left[ c - pT_pT_c - p + k(1 - T_c) \right] \\
&= X_{cp} + (1 - T_c) \left[ c + k - p - T_pT_c(k + p) \right] \\
&\geq X_{cp}
\end{aligned} \tag{11}$$

where the last inequality holds if:  $T_pT_c(k + p) < c + k - p$ .

Hence, the expected utility of cooperation, defection, and punishment take the form:

$$\begin{aligned}
E[U_c] &= \beta \overline{X}_c + (1 - \beta)X_c; & \overline{X}_c &> X_c \\
E[U_d] &= \alpha \overline{X}_d + (1 - \alpha)X_d & \overline{X}_d &< X_d \\
E[U_p] &= \beta \overline{X}_{cp} + (1 - \beta)X_{cp}; & \overline{X}_{cp} &> X_{cp}
\end{aligned} \tag{12}$$

where  $(\alpha, \beta) \in \Omega$ . The preceding expectations may written as:

$$\begin{aligned}
E[U_c] &= X_c + \beta(1 - T_c) \left[ c - T_pT_c(k + p) \right] \\
E[U_d] &= X_d + \alpha T_c \left[ (pT_pT_c - T_pk(2 - T_c) - c) \right] \\
E[U_p] &= X_{cp} + \beta(1 - T_c) \left[ c + k - p - T_pT_c(k + p) \right]
\end{aligned}$$

**Notes:**

\*  $\bar{X}_c \geq X_c \implies \bar{X}_d \leq X_d$  whenever  $T_c > \frac{2k}{k+p}$ . The converse is not necessarily true.

\*  $X_{cp} - X_c = -k(1 - T_c) < 0$ . Also:

$$X_c - X_d = \frac{b}{n} - c + pT_pT_c$$

## 2.1 EVOLUTION OF COOPERATION: $E[U_c] > E[U_d]$

**Case 1:**  $\bar{X}_c > X_c$ ;  $\bar{X}_d < X_d$

$$\begin{aligned} & E[U_c] > E[U_d] \\ \implies & X_c - X_d + \beta(1 - T_c) \left[ c - T_pT_c(k + p) \right] > \alpha T_c \left[ pT_pT_c - T_pk(2 - T_c) - c \right] \\ \implies & \beta > \frac{T_c [pT_pT_c - T_pk(2 - T_c) - c] \alpha + c - \frac{b}{n} - pT_pT_c}{(1 - T_c)[c - T_pT_c(k + p)]} \end{aligned} \quad (13)$$

Note here that setting  $T_p = 0$  in condition (13) recovers condition (7).

**Case 2:**  $\bar{X}_c < X_c$ ;  $\bar{X}_d < X_d$

$$\begin{aligned} & E[U_c] > E[U_d] \\ \implies & X_c - X_d + \alpha(1 - T_c) \left[ c - T_pT_c(k + p) \right] > \alpha T_c \left[ pT_pT_c - T_pk(2 - T_c) - c \right] \\ \implies & \alpha > \frac{c - \frac{b}{n} - pT_pT_c}{(1 - T_c)[c - T_pT_c(k + p)] - T_c [pT_pT_c - T_pk(2 - T_c) - c]} \end{aligned} \quad (14)$$

provided  $(1 - T_c)(c - T_pT_c(k + p)) > T_c(pT_pT_c - T_pk(2 - T_c) - c)$ .

**Case 3:**  $\bar{X}_c < X_c$ ;  $\bar{X}_d > X_d$

$$\begin{aligned} & E[U_c] > E[U_d] \\ \implies & \beta < \frac{(1 - T_c)[c - T_pT_c(k + p)] \alpha - c + \frac{b}{n} + pT_pT_c(n - 1)}{T_c[pT_pT_c - T_pk(2 - T_c) - c]}. \end{aligned} \quad (15)$$

**Case 4:**  $\bar{X}_c > X_c$ ;  $\bar{X}_d > X_d$

$$E[U_c] > E[U_d]$$

$$\implies \beta > \frac{c - \frac{b}{n} - pT_pT_c}{(1 - T_c)[c - T_pT_c(k + p)] - T_c[pT_pT_c - T_pk(2 - T_c) - c]}. \quad (16)$$

## 2.2 EVOLUTION OF PUNISHMENT: $E[U_p] > E[U_c]$

Costly punishment evolves whenever it yields a larger expected utility than cooperation alone. This expectation is influenced by within-group strategy distribution and the social welfare preferences of agents.

**Case 1:**  $\bar{X}_c > X_c$ ;  $\bar{X}_{cp} > X_{cp}$

$$E[U_p] > E[U_c]$$

$$\implies X_{cp} + \beta(1 - T_c) \left[ c + k - p - T_pT_c(k + p) \right] > X_c + \beta(1 - T_c) \left[ c - T_pT_c(k + p) \right]$$

$$X_{cp} - X_c + \beta(1 - T_c)(k - p) > 0$$

$$\implies \beta > \frac{k}{k - p} \quad (17)$$

provided  $k > p$ .

**Case 2:**  $\bar{X}_c < X_c$ ;  $\bar{X}_{cp} < X_{cp}$

$$E[U_p] > E[U_c]$$

$$\implies \alpha > \frac{k}{k - p}. \quad (18)$$

**Case 3:**  $\bar{X}_c > X_c$ ;  $\bar{X}_{cp} < X_{cp}$

$$E[U_p] > E[U_c]$$

$$\implies X_{cp} + \alpha(1 - T_c) \left[ c + k - p - T_pT_c(k + p) \right] > X_c + \beta(1 - T_c) \left[ c - T_pT_c(k + p) \right]$$

$$\implies \beta < \frac{[c + k - p - T_pT_c(k + p)] \alpha - k}{c - T_pT_c(k + p)} \quad (19)$$

**Case 4:**  $\overline{X}_c < X_c$ ;  $\overline{X}_{cp} > X_{cp}$

$$\begin{aligned} E[U_p] &> E[U_c] \\ \implies \beta &> \frac{[c - T_p T_c(k + p)] \alpha + k}{c + k - p - T_p T_c(k + p)} \end{aligned} \quad (20)$$