

# Supplementary Material

## Costly incentives design from an institutional perspective: cooperation, sustainability, and affluence

Xin Zhou\*, Adam Belloum, Michael H. Lees, Tom van Engers, Cees de Laat  
*University of Amsterdam, Amsterdam, 1098XH, Netherlands*

### Supplementary part 1. Population equilibrium

To predict the effect of incentive mechanism, **evolutionary game theory** which describes the population of agents engaging in pairwise interaction, has been generally accepted as a common framework to model and interpret the evolution of cooperation in a social dilemma. This part shows the derivation process of the analytical results of the population in equilibrium.

Assume a finite population of size  $N$ ,  $M$  of them participant a market play PDG, two types of strategies for cooperation  $C$  and defection  $D$  are well mixed. Let  $\pi(C)$  denote the expected payoff of strategy  $C$ ,  $\pi(D)$  denote that of strategy  $D$ , and  $\bar{\pi}(\mathbf{x})$  denote the average payoff of the whole population:

$$\pi(C) = (1 - c_0 + R_{CC})x + (-T - c_0 + R_{CD})y \quad (1)$$

$$\pi(D) = (T - c_0 - F_{CD})x + (0 - c_0 - F_{DD})y \quad (2)$$

$$\bar{\pi}(\mathbf{x}) = \pi(C)x + \pi(D)y \quad (3)$$

The replicator dynamics of cooperators and defectors is:

$$\dot{x} = x(1-x)[\pi(C) - \pi(D)]. \quad (4)$$

With the equations (1), (2), (3), and (4), let  $\dot{x} = 0$ ,

$$(x^2 - x)[x + xR_{CC} + (1-x)R_{CD} + xF_{CD} + (1-x)F_{DD} - T] = 0 \quad (5)$$

Thus, the fixed points are  $x^* = 0$ ,  $x^* = 1$ , and  $x^* = \frac{-F_{DD} + T - R_{CD}}{1 + R_{CC} - R_{CD} + F_{CD} - F_{DD}}$ .

We next discuss if the fixed point is the Nash equilibrium (NE) and satisfy the requirement of being as the evolutionary stable strategy (ESS).

Let  $\mathbf{s}^* = (x^*, 1 - x^*)$ ,  $\mathbf{s} = (p, 1 - p)$ ,  $\mathbf{s} \neq \mathbf{s}^*$ , an ESS can be defined as a mixed strategy  $\mathbf{s}^*$ , such that for any strategy  $\mathbf{s}$  and any sufficient small  $\varepsilon > 0$ ,

$$\pi[\mathbf{s}^*, (1 - \varepsilon)\mathbf{s}^* + \varepsilon\mathbf{s}] > \pi[\mathbf{s}, (1 - \varepsilon)\mathbf{s}^* + \varepsilon\mathbf{s}]. \quad (6)$$

Using the linearity in probability of expected payoffs reduces 6 to:

$$(1 - \varepsilon)\pi(\mathbf{s}^*, \mathbf{s}^*) + \varepsilon\pi(\mathbf{s}^*, \mathbf{s}) > (1 - \varepsilon)\pi(\mathbf{s}, \mathbf{s}^*) + \varepsilon\pi(\mathbf{s}, \mathbf{s}). \quad (7)$$

If for all small  $\varepsilon > 0$  and for all  $\mathbf{s}$ ,

$$\pi(\mathbf{s}^*, \mathbf{s}^*) \geq \pi(\mathbf{s}, \mathbf{s}^*), \quad (8)$$

then  $\mathbf{s}^*$  is a symmetric NE. Further, if

$$\pi(\mathbf{s}^*, \mathbf{s}) > \pi(\mathbf{s}, \mathbf{s}) \quad (9)$$

whenever  $\pi(\mathbf{s}^*, \mathbf{s}^*) = \pi(\mathbf{s}, \mathbf{s}^*)$ ,  $\mathbf{s}^*$  is an ESS [1].

1) fixed point at  $x^* = 0$ ,  $\mathbf{s}^* = (0, 1)$ ,  $\mathbf{s} = (p, 1 - p)$  ( $0 < p \leq 1$ ), we have:

$$\pi(\mathbf{s}^*, \mathbf{s}^*) = -c_0 - F_{DD} \quad (8a)$$

$$\pi(\mathbf{s}, \mathbf{s}^*) = p(-T - c_0 + R_{CD}) + (1 - p)(0 - c_0 - F_{DD}) \quad (8b)$$

$$\pi(\mathbf{s}^*, \mathbf{s}) = p(T - c_0 - F_{CD}) + (1 - p)(0 - c_0 - F_{DD}) \quad (9a)$$

$$\begin{aligned} \pi(\mathbf{s}, \mathbf{s}) = & p^2(1 - c_0 + R_{CC}) + p(1 - p)(-T - c_0 + R_{CD}) \\ & + p(1 - p)(T - c_0 - F_{CD}) + (1 - p)^2(0 - c_0 - F_{DD}) \end{aligned} \quad (9b)$$

By 8a - 8b, we have

$$\begin{aligned} 8a - 8b &= -p(-T - c_0 + R_{CD}) + p(0 - c_0 - F_{DD}) \\ &= p(T - F_{DD} - R_{CD}), \end{aligned} \quad (8c)$$

and by 9a - 9b, we have

$$9a - 9b = p^2(T - F_{CD} - R_{CC}) + p(1 - p)(-F_{DD} + T - R_{CD}). \quad (9c)$$

The Table 1 shows the requirements for  $x^* = 0$  being a NE or an ESS under different incentive policies.

Supplementary table 1. NE and ESS analysis under different conditions when  $x^* = 0$

Incentive Policy	Scale of parameters	NE	ESS
<b>Reward</b>	$0 < R_{CC} + R_{CD}, F_{CD} = F_{DD} = 0$	When $T \geq R_{CD}$ , $8c \geq 0$	When $T > R_{CD}$ , $9c > 0$
<b>Punishment</b>	$R_{CD} = R_{DD} = 0, 0 < F_{CD} + F_{DD}$	When $T \geq F_{DD}$ , $8c \geq 0$	When $T \geq F_{CD}$ , $9c > 0$
<b>Mixed incentives</b>	$R_{CD} + R_{CC} \neq 0, F_{CD} + F_{DD} \neq 0$	When $T \geq R_{CD} + F_{DD}$ , $8c \geq 0$	When $T \geq \max[R_{CD} + F_{DD}, R_{CC} + F_{CD}]$ , $9c > 0$

2) Consider the fixed point at  $x^* = 1$ ,  $\mathbf{s}^* = (1, 0)$ ,  $\mathbf{s} = (1 - p, p)$  ( $0 < p \leq 1$ ), we have:

$$\pi(\mathbf{s}^*, \mathbf{s}^*) = 1 - c_0 + R_{CC} \quad (8d)$$

$$\pi(\mathbf{s}, \mathbf{s}^*) = (1 - p)(1 - c_0 + R_{CC}) + p(T - c_0 - F_{CD}) \quad (8e)$$

$$\pi(\mathbf{s}^*, \mathbf{s}) = (1 - p)(1 - c_0 + R_{CC}) + p(-T - c_0 + R_{CD}) \quad (9d)$$

$$\begin{aligned} \pi(\mathbf{s}, \mathbf{s}) = & (1 - p)^2(1 - c_0 + R_{CC}) + p(1 - p)(-T - c_0 + R_{CD}) \\ & + p(1 - p)(T - c_0 - F_{CD}) + p^2(0 - c_0 - F_{DD}) \end{aligned} \quad (9e)$$

By 8d - 8e, we have

$$\begin{aligned} 8d - 8e &= p(1 - c_0 + R_{CC}) - p(T - c_0 - F_{CD}) \\ &= p(1 - T + R_{CC} + F_{CD}), \end{aligned} \quad (8f)$$

and by 9d - 9e, we have

$$\begin{aligned} 9d - 9e &= p(1 - p)(1 - c_0 + R_{CC}) - p(1 - p)(T - c_0 - F_{CD}) \\ &+ p^2(-T - c_0 + R_{CD}) - p^2(0 - c_0 - F_{DD}) \\ &= p(1 - p)(1 - T + R_{CC} + F_{CD}) + p^2(-T + R_{CD} + F_{DD}). \end{aligned} \quad (9f)$$

The Table 2 shows the requirements for  $x^* = 1$  being a NE or an ESS under different incentive policies.

Supplementary table 2. NE and ESS analysis under different conditions when  $x^* = 1$

Incentive Policy	Scale of parameters	NE	ESS
<b>Reward</b>	$0 < R_{CC} + R_{CD}, F_{CD} = F_{DD} = 0$	When $R_{CC} \geq T - 1, 8f \geq 0$	When $R_{CC} \geq T, 9f > 0$
<b>Punishment</b>	$R_{CD} = R_{DD} = 0, 0 < F_{CD} + F_{DD}$	When $F_{CD} \geq T - 1, 8f \geq 0$	When $F_{DD} > T, 9f > 0$
<b>Mixed incentives</b>	$R_{CD} + R_{CC} \neq 0, F_{CD} + F_{DD} \neq 0$	When $R_{CC} + F_{CD} \geq T - 1, 8f \geq 0$	When $\min[R_{CD} + F_{DD}, R_{CC} + F_{CD}] \geq T, 9f > 0$

- 3) Consider the fixed point at  $x^* = \frac{T - R_{CD} - F_{DD}}{1 + R_{CC} - R_{CD} + F_{CD} - F_{DD}} = \frac{B}{A}$ . Let  $q = \frac{B}{A}, \mathbf{s}^* = (q, 1 - q), \mathbf{s} = (p, 1 - p)$  ( $0 \leq p \leq 1, p \neq q$ ), we have:

$$\begin{aligned} \pi(\mathbf{s}^*, \mathbf{s}^*) &= q^2(1 - c_0 + R_{CC}) + q(1 - q)(-T - c_0 + R_{CD}) \\ &\quad + q(1 - q)(T - c_0 - F_{CD}) + (1 - q)^2(0 - c_0 - F_{DD}) \end{aligned} \quad (8g)$$

$$\begin{aligned} \pi(\mathbf{s}, \mathbf{s}^*) &= pq(1 - c_0 + R_{CC}) + p(1 - q)(-T - c_0 + R_{CD}) \\ &\quad + (1 - p)q(T - c_0 - F_{CD}) + (1 - p)(1 - q)(0 - c_0 - F_{DD}) \end{aligned} \quad (8h)$$

$$\begin{aligned} \pi(\mathbf{s}^*, \mathbf{s}) &= pq(1 - c_0 + R_{CC}) + (1 - p)q(-T - c_0 + R_{CD}) \\ &\quad + p(1 - q)(T - c_0 - F_{CD}) + (1 - p)(1 - q)(0 - c_0 - F_{DD}) \end{aligned} \quad (9g)$$

$$\begin{aligned} \pi(\mathbf{s}, \mathbf{s}) &= p^2(1 - c_0 + R_{CC}) + p(1 - p)(-T - c_0 + R_{CD}) \\ &\quad + p(1 - p)(T - c_0 - F_{CD}) + (1 - p)^2(0 - c_0 - F_{DD}) \end{aligned} \quad (9h)$$

By 8g - 8h, we have

$$\begin{aligned} 8g - 8h &= q(q - p)(1 - c_0 + R_{CC} - T + c_0 + F_{CD}) + (1 - q)(q - p)(-T - c_0 + R_{CD} + c_0 + F_{DD}) \\ &= q(q - p)(1 - T + R_{CC} + F_{CD}) + (1 - q)(q - p)(-T + R_{CD} + F_{DD}) \\ &= (q - p)[q(1 + R_{CC} + F_{CD} - R_{CD} - F_{DD}) - T + R_{CD} + F_{DD}] \\ &= (q - p)[q(1 + R_{CC} + F_{CD}) + (1 - q)(R_{CD} + F_{DD})], \end{aligned} \quad (8i)$$

and by 9g - 9h, we have

$$\begin{aligned} 9g - 9h &= p(q - p)(1 - c_0 + R_{CC} - T + c_0 + F_{CD}) + (1 - p)(q - p)(-T - c_0 + R_{CD} + c_0 + F_{DD}) \\ &= p(q - p)(1 - T + R_{CC} + F_{CD}) + (1 - p)(q - p)(-T + R_{CD} + F_{DD}) \\ &= (q - p)[p(1 + R_{CC} + F_{CD} - R_{CD} - F_{DD}) - T + R_{CD} + F_{DD}] \\ &= (q - p)[p(1 + R_{CC} + F_{CD}) + (1 - p)(R_{CD} + F_{DD}) - T]. \end{aligned} \quad (9i)$$

The Table 3 shows the requirements for  $x^* = q$  being a NE or an ESS under different incentive policies.

Supplementary table 3. NE and ESS analysis under different conditions when  $x^* = q$

Incentive Policy	Scale of parameters	NE	ESS
<b>Reward</b>	$0 < R_{CC} + R_{CD}, F_{CD} = F_{DD} = 0$	$8i = 0$ , thus always hold	$9i < 0$ , thus the ESS is not hold
<b>Punishment</b>	$R_{CD} = R_{DD} = 0, 0 < F_{CD} + F_{DD}$	$8i = 0$ , thus always hold	$9i < 0$ , thus the ESS is not hold
<b>Mixed incentives</b>	$R_{CD} + R_{CC} \neq 0, F_{CD} + F_{DD} \neq 0$	$8i = 0$ , thus always hold	$9i < 0$ , thus the ESS is not hold

## Supplementary part 2. Rate of $R_{CC}$ ( $F_{DD}$ ) in $R_{CC} + F_{CD}$ ( $R_{CD} + F_{DD}$ )

Let  $\alpha$  be the rate,  $k = R_{CC} + F_{CD}$ , to minimize the difference between reward ( $R_{CD} + R_{CC}$ ) and punishment ( $F_{CD} + F_{DD}$ ), meanwhile satisfy the constraints are  $R_{CD} \geq R_{CC}$ ,  $F_{CD} \geq F_{DD}$ , we have:

$$\begin{aligned} \min_{\alpha} \quad & k\alpha + (k+1)\alpha - [k(1-\alpha) + (1+k)\alpha] \\ \text{s.t.} \quad & (1-\alpha)2 \geq 3\alpha \\ & (1-\alpha)1 \geq 4\alpha \end{aligned} \tag{10}$$

The solution is  $\alpha = 0.2$ . Consequently, these four parameters are set as shown in Table 4.

Supplementary table 4. Mixed incentives setup

$R_{CC} + F_{CD}$	$R_{CC}$	$F_{CD}$	$R_{CD} + F_{DD}$	$R_{CD}$	$F_{DD}$
1	0.2	0.8	2	1.6	0.4
1.25	0.25	1	2.25	1.8	0.45
1.5	0.3	1.2	2.5	2	0.5
1.75	0.35	1.4	2.75	2.2	0.55
2	0.4	1.6	3	2.4	0.6
2.25	0.45	1.8	3.25	2.6	0.65
2.5	0.5	2	3.5	2.8	0.7
2.75	0.55	2.2	3.75	3	0.75
3	0.6	2.4	4	3.2	0.8

## Supplementary part 3. Accumulated wealth of the third-party

From equation 4, we have

$$\begin{aligned} \dot{x} = \frac{dx}{dt} = & (1 + R_{CC} - R_{CD} + F_{CD} - F_{DD})x^3 \\ & + (2R_{CD} + 2F_{DD} - F_{CD} - R_{CC} - T - 1)x^2 \\ & - (R_{CD} + F_{DD} - T)x \end{aligned} \tag{11}$$

Let  $a = 1 + R_{CC} - R_{CD} + F_{CD} - F_{DD}$ ,  $b = 2R_{CD} + 2F_{DD} - F_{CD} - R_{CC} - T - 1$ ,  $c = T - R_{CD} - F_{DD}$ , we have:

$$\frac{dx}{dt} = ax^3 + bx^2 + cx \tag{12}$$

We can then solve the differential equation:

$$\begin{aligned} \int \frac{1}{ax^3 + bx^2 + cx} dx &= \int \frac{1}{x(ax^2 + bx + c)} dx \\ &= \int \frac{1}{x(ax^2 + bx + c)} dx \\ &= \int \frac{1}{cx} + \frac{-ax - b}{c(ax^2 + bx + c)} dx \\ &= \frac{1}{c} \int \frac{1}{x} dx - \frac{a}{c} \int \frac{x}{ax^2 + bx + c} dx - \frac{b}{c} \int \frac{1}{ax^2 + bx + c} dx \\ &= \frac{1}{c} \ln x - \frac{a}{c} \left( \frac{\ln(ax^2 + bx + c)}{2a} - \frac{b}{a\sqrt{4ac - b^2}} \tan^{-1} \left( \frac{2ax + b}{\sqrt{4ac - b^2}} \right) \right) \\ &\quad - \frac{b}{c} \frac{2}{a\sqrt{4ac - b^2}} \tan^{-1} \left( \frac{2ax + b}{\sqrt{4ac - b^2}} \right) + C \\ &= t \end{aligned} \tag{13}$$

$x^{(t)}$  is the inverse function of 13, let  $x^{(t)} := F(t, a, b, c)$ . The wealth of the third party at time step  $t$  is shown as 2.3

$$W_T^{(t)} = M^{(t)} \left( c_0 + x^{(t)}(1-x^{(t)})F_{CD} + (1-x^{(t)})^2 F_{DD} - (x^{(t)})^2 R_{CC} - x^{(t)}(1-x^{(t)})R_{CD} - \alpha \left( x^{(t)}(1-x^{(t)})F_{CD} + (1-x^{(t)})^2 F_{DD} \right) \right) \quad (14)$$

$$= M^{(t)} \left( (x^{(t)})^2 ((1-\alpha)(F_{DD} - F_{CD}) - R_{CC} + R_{CD}) + x^{(t)}((1-\alpha)(F_{CD} - 2F_{DD}) - R_{CD}) + c_0 + (1-\alpha)F_{DD} \right) \quad (15)$$

$$= M^{(t)} \left( F^2(t, a, b, c) ((1-\alpha)(F_{DD} - F_{CD}) - R_{CC} + R_{CD}) + F(t, a, b, c) ((1-\alpha)(F_{CD} - 2F_{DD}) - R_{CD}) + c_0 + (1-\alpha)F_{DD} \right) \quad (16)$$

$$(17)$$

Thus,

$$W_T = \int W_T^{(t)} dt \quad (18)$$

$$= ((1-\alpha)(F_{DD} - F_{CD}) - R_{CC} + R_{CD}) \int M^{(t)} F^2(t, a, b, c) dt + ((1-\alpha)(F_{CD} - 2F_{DD}) - R_{CD}) \int M^{(t)} F(t, a, b, c) dt + (c_0 + (1-\alpha)F_{DD}) \int M^{(t)} dt \quad (19)$$

## References

1. Crawford, V.P. Learning and mixed-strategy equilibria in evolutionary games. *J.Theor.Biol.* **140**, 537–550 (1989).