

Poisoning Attacks in Games Example

-

June 20, 2021

1 Correlated Equilibrium Example

1.1 Numerical Example

From Roth 2017,

-	A	B	C
A	(1, 1)	(-1, -1)	(0, 0)
B	(-1, -1)	(1, 1)	(0, 0)
C	(0, 0)	(0, 0)	(-1.1, -1.1)

- Dominant strategy: none
- Pure strategy Nash: $(A, A), (B, B)$
- Mixed strategy Nash: $\left(A \left(\frac{1}{2} \right) B \left(\frac{1}{2} \right), A \left(\frac{1}{2} \right) B \left(\frac{1}{2} \right) \right)$

Check indifference conditions:

$$M \left(A, A \left(\frac{1}{2} \right) B \left(\frac{1}{2} \right) \right) = \frac{1}{2} (1) + \frac{1}{2} (-1) = 0$$
$$M \left(B, A \left(\frac{1}{2} \right) B \left(\frac{1}{2} \right) \right) = \frac{1}{2} (-1) + \frac{1}{2} (1) = 0$$

Check support conditions:

$$M \left(C, A \left(\frac{1}{2} \right) B \left(\frac{1}{2} \right) \right) = 0 \geq \frac{1}{2} (0) + \frac{1}{2} (0)$$

- Correlated equilibrium: $\left((A, A) \left(\frac{1}{2} \right) (B, B) \left(\frac{1}{2} \right) \right)$

The equilibrium payoff is:

$$M \left((A, A) \left(\frac{1}{2} \right) (B, B) \left(\frac{1}{2} \right) \right) = \frac{1}{2} (1) + \frac{1}{2} (1) = 1$$

Check best response conditioned on the player receiving the signal A (i.e. the player infers that the strategy is (A, A) so the other player will play A):

$$M \left(B, (A, A) \left(\frac{1}{2} \right) (B, B) \left(\frac{1}{2} \right) | A \right) = M(B, A) = -1 < 1$$

$$M \left(C, (A, A) \left(\frac{1}{2} \right) (B, B) \left(\frac{1}{2} \right) | A \right) = M(C, A) = 0 < 1$$

Check best response conditioned on the player receiving the signal B (i.e. the player infers that the strategy is (B, B) so the other player will play B):

$$M \left(A, (A, A) \left(\frac{1}{2} \right) (B, B) \left(\frac{1}{2} \right) | B \right) = M(A, B) = -1 < 1$$

$$M \left(C, (A, A) \left(\frac{1}{2} \right) (B, B) \left(\frac{1}{2} \right) | B \right) = M(C, B) = 0 < 1$$

This cannot be implemented as a mixed strategy equilibrium since the actions are not independent.

- Coarse correlated equilibrium: $\left((A, A) \left(\frac{1}{3} \right) (B, B) \left(\frac{1}{3} \right) (C, C) \left(\frac{1}{3} \right) \right)$

The equilibrium payoff is:

$$M \left((A, A) \left(\frac{1}{3} \right), (B, B) \left(\frac{1}{3} \right), (C, C) \left(\frac{1}{3} \right) \right) = \frac{1}{3} (1) + \frac{1}{3} (1) + \frac{1}{3} (-1.1) = 0.3$$

Check best response not conditioned on the signal (i.e. the other player receives each of the signals A, B, C with probability $\frac{1}{3}$):

$$M \left(A, (A, A) \left(\frac{1}{3} \right) (B, B) \left(\frac{1}{3} \right) (C, C) \left(\frac{1}{3} \right) \right) = M \left(A, A \left(\frac{1}{3} \right) B \left(\frac{1}{3} \right) C \left(\frac{1}{3} \right) \right) = \frac{1}{3} (1) + \frac{1}{3} (-1) + \frac{1}{3} (0) = 0 < 0.3$$

$$M \left(B, (A, A) \left(\frac{1}{3} \right) (B, B) \left(\frac{1}{3} \right) (C, C) \left(\frac{1}{3} \right) \right) = M \left(B, A \left(\frac{1}{3} \right) B \left(\frac{1}{3} \right) C \left(\frac{1}{3} \right) \right) = \frac{1}{3} (-1) + \frac{1}{3} (1) + \frac{1}{3} (0) = 0 < 0.3$$

$$M \left(C, (A, A) \left(\frac{1}{3} \right) (B, B) \left(\frac{1}{3} \right) (C, C) \left(\frac{1}{3} \right) \right) = M \left(C, A \left(\frac{1}{3} \right) B \left(\frac{1}{3} \right) C \left(\frac{1}{3} \right) \right) = \frac{1}{3} (0) + \frac{1}{3} (0) + \frac{1}{3} (-1.1) < 0 < 0.3$$

This is not a correlated equilibrium since the best response conditioned on the player receiving the signal C is not satisfied:

$$M \left(A, (A, A) \left(\frac{1}{3} \right) (B, B) \left(\frac{1}{3} \right) (C, C) \left(\frac{1}{3} \right) | C \right) = M(A, C) = 0 > -1.1$$

$$M \left(B, (A, A) \left(\frac{1}{3} \right) (B, B) \left(\frac{1}{3} \right) (C, C) \left(\frac{1}{3} \right) | C \right) = M(B, C) = 0 > -1.1$$

- Relation

Dominant \Rightarrow Pure Nash \Rightarrow Mixed Nash \Rightarrow Correlated \Rightarrow Coarse Correlated

1.2 MAB Book

- Coarse correlated equilibrium:

$$\mathbb{E}_{(i,j) \sim \sigma} [M(i, j) - M(i_0, j)] \geq 0 \quad \forall i_0.$$

- Correlated equilibrium:

$$\mathbb{E}_{(i,j) \sim \sigma} [M(i, j) - M(i_0, j) | i] \geq 0 \quad \forall i_0.$$

- What if the signal also specify the other player's action? Convexify the set of pure strategy Nash?

$$\mathbb{E}_{(i,j) \sim \sigma} [M(i, j) - M(i_0, j) | (i, j)] \geq 0 \quad \forall i_0.$$

- Mixed strategy equilibrium?

$$\mathbb{E}_{i \sim \sigma_i, j \sim \sigma_j} [M(i, j) - M(i_0, j) | i] \geq 0 \quad \forall i_0.$$

2 Dominant Strategy Implementability

2.1 PD Games Symmetric Implementation

Suppose the original game is,

Action	Cooperate	Defect
Cooperate	x, x	$0, y$
Defect	$y, 0$	$1, 1$

Suppose the modified game is,

Action	Cooperate	Defect
Cooperate	$x + \delta_1, x + \delta_1$	$0 + \delta_2, y + \delta_3$
Defect	$y + \delta_3, 0 + \delta_2$	$1 + \delta_4, 1 + \delta_4$

$$\min_{\delta} \delta_1^2 + \delta_2^2 + \delta_3^2 + \delta_4^2 \quad (1)$$

such that

$$x + \delta_1 \geq y + \delta_3 + \varepsilon \quad (2)$$

$$0 + \delta_2 \geq y + \delta_4 + \varepsilon \quad (3)$$

Dominant strategy implementation Lagrange multiplier:

$$\mathcal{L} = \delta_1^2 + \delta_2^2 + \delta_3^2 + \delta_4^2 - \lambda (x + \delta_1 - y - \delta_3 - \varepsilon) - \mu (\delta_2 - y - \delta_4 - \varepsilon)$$

$$\frac{\partial \mathcal{L}}{\partial \delta_1} = 2\delta_1 - \lambda = 0$$

$$\frac{\partial \mathcal{L}}{\partial \delta_2} = 2\delta_2 - \mu = 0$$

$$\frac{\partial \mathcal{L}}{\partial \delta_3} = 2\delta_3 + \lambda = 0$$

$$\frac{\partial \mathcal{L}}{\partial \delta_4} = 2\delta_4 + \mu = 0$$

$$\lambda (x + \delta_1 - y - \delta_3 - \varepsilon) = 0$$

$$\mu (\delta_2 - y - \delta_4 - \varepsilon) = 0$$

The interior solution is given by,

$$\lambda = y - x + \varepsilon \geq 0$$

$$\mu = y + \varepsilon \geq 0$$

$$\delta_1 = \frac{1}{2} (y - x + \varepsilon)$$

$$\delta_2 = \frac{1}{2} (y + \varepsilon)$$

$$\delta_3 = -\frac{1}{2} (y - x + \varepsilon)$$

$$\delta_4 = -\frac{1}{2} (y + \varepsilon)$$

The boundary solutions require,

$$\lambda = 0 \Rightarrow x - y - \varepsilon \geq 0$$

$$\mu = 0 \Rightarrow -y - \varepsilon \geq 0$$

These are not feasible because the inequalities are not satisfied for sufficiently small ε .

Bayesian (Nash) implementation Lagrange multiplier:

$$\begin{aligned}
\mathcal{L} &= \delta_1^2 + \delta_2^2 + \delta_3^2 + \delta_4^2 - \lambda (x + \delta_1 - y - \delta_3 - \varepsilon) \\
\frac{\partial \mathcal{L}}{\partial \delta_1} &= 2\delta_1 - \lambda = 0 \\
\frac{\partial \mathcal{L}}{\partial \delta_2} &= 2\delta_2 = 0 \\
\frac{\partial \mathcal{L}}{\partial \delta_3} &= 2\delta_3 + \lambda = 0 \\
\frac{\partial \mathcal{L}}{\partial \delta_4} &= 2\delta_4 = 0 \\
\lambda (x + \delta_1 - y - \delta_3 - \varepsilon) &= 0
\end{aligned}$$

The interior solution is given by,

$$\begin{aligned}
\lambda &= y - x + \varepsilon \geq 0 \\
\delta_1 &= \frac{1}{2} (y - x + \varepsilon) \\
\delta_2 &= 0 \\
\delta_3 &= -\frac{1}{2} (y - x + \varepsilon) \\
\delta_4 &= 0
\end{aligned}$$

The boundary solutions require,

$$\lambda = 0 \Rightarrow x - y - \varepsilon \geq 0$$

This is not feasible because the inequality is not satisfied for sufficiently small ε .

3 Uniqueness of Nash Equilibrium

3.1 Two by Two Games

A game with a unique Nash Equilibrium that is not a dominant strategy equilibrium.

Action	A	B
A	(1, 1)	(0, 0)
B	(0, 1)	(1, 0)

For a general two player two by two game, if it has a unique pure strategy Nash Equilibrium, then at least one player is using a dominant strategy.

Action	A	B
A	(a, b)	(c, d)
B	(e, f)	(g, h)

Proof: without loss of generality, suppose (A, A) is the unique pure strategy Nash Equilibrium. It implies that $a \geq e$ and $b \geq d$. Since (B, B) is not a Nash Equilibrium, either $c > g$ or $f > h$, and in either case A is a dominant strategy for one of the players.

For a symmetric two player two by two game, if it has a unique pure strategy Nash Equilibrium, then it is a dominant strategy equilibrium.

For a symmetric two player two by two game, it has a unique strict pure strategy Nash Equilibrium iff it is a strictly dominant strategy equilibrium.

A game with a dominant strategy equilibrium that is not the unique Nash Equilibrium.

Action	A	B
A	(1, 1)	(1, 0)
B	(1, 1)	(0, 0)

For a general two player two by two game specified below, it has a unique (completely) mixed strategy Nash Equilibrium iff $(b - d)(f - h) < 0$, $(a - e)(c - g) < 0$ and either $(b - d)(a - e) < 0$ or $(f - h)(c - g) < 0$.

Action	β	$1 - \beta$
α	(a, b)	(c, d)
$1 - \alpha$	(e, f)	(g, h)

Proof: suppose the game has a mixed strategy Nash Equilibrium in which the row player mixes action A with probability α and the column player mixes action A with probability β , then it is given by the indifference conditions:

$$\begin{aligned}\alpha b + (1 - \alpha) f &= \alpha d + (1 - \alpha) h, \\ \beta a + (1 - \beta) c &= \beta e + (1 - \beta) g.\end{aligned}$$

Then we have,

$$\begin{aligned}\alpha &= \frac{h - f}{h - f + b - d} \in (0, 1), \\ \beta &= \frac{a - e}{a - e + g - c} \in (0, 1).\end{aligned}$$

If $h > f$, then $\alpha \in (0, 1)$ implies $b > d$, and since (B, B) is not a Nash Equilibrium, $c > g$.

If $f > h$, then $\alpha \in (0, 1)$ implies $d > b$, and since (B, A) is not a Nash Equilibrium, $a > e$.

If $a > e$, then $\beta \in (0, 1)$ implies $g > c$, and since (A, A) is not a Nash Equilibrium, $d > b$.

If $e > a$, then $\beta \in (0, 1)$ implies $c > g$, and since (A, B) is not a Nash Equilibrium, $h > f$.

Therefore, we have $(b - d)(f - h) < 0$, $(a - e)(c - g) < 0$ and $(b - d)(a - e) < 0$.

For a general two player two by two game specified below, it has a unique (completely) mixed strategy Nash Equilibrium if $(b - d)(f - h) \neq 0$, $(a - e)(c - g) \neq 0$ and either $(b - d)(a - e) < 0$ or $(f - h)(c - g) < 0$.

Action	β	$1 - \beta$
α	(a, b)	(c, d)
$1 - \alpha$	(e, f)	(g, h)

Proof: if $(b - d)(f - h) > 0$, the column player has a strictly dominant strategy, $a \neq e$ and $c \neq g$ implies that there is a unique pure strategy Nash.

if $(a - e)(c - g) > 0$, the row player has a strictly dominant strategy, $b \neq d$ and $f \neq h$ implies that there is a unique pure strategy Nash.

if $(b - d)(f - h) < 0$ and $(a - e)(c - g) < 0$, then there is a unique mixed strategy Nash Equilibrium.

3.2 Rosen Characterization

Let $x_1 = \alpha \in [0, 1]$ be the action set S_1 for the row player.

Let $x_2 = \beta \in [0, 1]$ be the action set S_2 for the column player.

Then the payoff functions are,

$$\begin{aligned}\varphi_1(x_1, x_2) &= ax_1x_2 + e(1 - x_1)x_2 + cx_1(1 - x_2) + g(1 - x_1)(1 - x_2), \\ \varphi_2(x_1, x_2) &= bx_1x_2 + f(1 - x_1)x_2 + dx_1(1 - x_2) + h(1 - x_1)(1 - x_2).\end{aligned}$$

The pseudo-gradient function is given by,

$$\begin{aligned}g(x, r) &= \begin{bmatrix} r_1 \nabla_1 \varphi_1(x) \\ r_2 \nabla_2 \varphi_2(x) \end{bmatrix} \\ &= \begin{bmatrix} r_1((a - e)x_2 + (c - g)(1 - x_2)) \\ r_2((b - d)x_1 + (f - h)(1 - x_1)) \end{bmatrix}.\end{aligned}$$

The Jacobian of $g(x, r)$ is given by,

$$G(x, r) = \begin{bmatrix} 0 & r_1((a - e) - (c - g)) \\ r_2((b - d) - (f - h)) & 0 \end{bmatrix}.$$

Then $\sigma(x, r)$ is diagonally strictly concave (Theorem 6) if $G(x, r) + G^T(x, r)$ is negative definite for some $r > 0$, here,

$$(2r_1((a - e) - (c - g)) + 2r_2((b - d) - (f - h)))^2 < 0,$$

which is never true?

3.3 Moulin Characterization

Dominance Solvable Nash Equilibrium is unique. For two by two games, it means at least one player has a dominant strategy, and the other player is not indifferent between the two actions.

3.4 Dominant Strategy Implementability

Suppose the agents only have one type, then the dominant strategy implementation of s^\dagger is,

$$\min \sum_{t=1}^T \|\delta^t\| \quad \text{such that}$$

$$R_i(s_i^\dagger, a_{-i}^t) + \delta_{i, s_i^\dagger, a_{-i}^t}^t \geq R_i(a_i^t, a_{-i}^t) + \delta_{i, a_i^t, a_{-i}^t}^t \quad \forall a_i^t, \forall a_{-i}^t, \forall i, \forall t.$$

Suppose the agents have types given by the belief p of the game, then a direct revelation dominant strategy implementation of $s^\dagger(p)$ is (the agents report p and the designer chooses $s^\dagger(p)$ for the players),

$$\min \sum_{t=1}^T \|\delta^t\| \quad \text{such that}$$

$$R_i(s^\dagger(p_i^t, q_{-i}^t)) + \delta_{i, p_i^t, q_{-i}^t}^t \geq R_i(s^\dagger(q_i^t, q_{-i}^t)) + \delta_{i, q_i^t, q_{-i}^t}^t \quad \forall q_i^t, \forall q_{-i}^t, \forall i, \forall t.$$

Suppose the agents have types given by the belief $p \sim P$ distribution over types is common knowledge, then a direct revelation Bayesian implementation of $s^\dagger(p)$ is,

$$\min \sum_{t=1}^T \|\delta^t\| \quad \text{such that}$$

$$\mathbb{E}_{q_{-i}^t \sim P} R_i(s^\dagger(p_i^t, q_{-i}^t)) + \delta_{i, p_i^t, q_{-i}^t}^t \geq \mathbb{E}_{q_{-i}^t \sim P} R_i(s^\dagger(q_i^t, q_{-i}^t)) + \delta_{i, q_i^t, q_{-i}^t}^t \quad \forall q_i^t, \forall i, \forall t.$$

4 Information Cascade Example

Consider the repeated game with the following stage game,

Actions	Accept	Reject
Accept	θ, θ	$0, 0$
Reject	$0, 0$	$-\theta, -\theta$

Assume a common prior,

$$\theta = \begin{cases} 1 & \text{with probability } p = \frac{1}{2} \\ -1 & \text{with probability } 1 - p = \frac{1}{2} \end{cases}. \quad (4)$$

Now suppose the reward function is given by,

$$R_i^t(\theta = 1) = \begin{cases} 1 & \text{with probability } q = \frac{2}{3} \\ -1 & \text{with probability } 1 - q = \frac{1}{3} \end{cases}, \quad (5)$$

and,

$$R_i^t(\theta = -1) = \begin{cases} 1 & \text{with probability } 1 - q = \frac{1}{3} \\ -1 & \text{with probability } q = \frac{2}{3} \end{cases}. \quad (6)$$

In a world with $\theta = -1$,

In stage 1, with probability $q^2 = \frac{1}{9}$, we have $R_i^1(\theta = -1) = 1$, and the posterior belief that $\theta = 1$ is,

$$\mathbb{P}_i \{ \theta = 1 | R_i^1 = 1 \} = \frac{pq}{pq + (1-p)(1-q)} = \frac{2}{3} > \frac{1}{2}. \quad (7)$$

This implies that both players will select the action Accept.

In stage 2 after the history (Accept, Accept), the posterior belief of both players that $\theta = -1$ is,

$$\mathbb{P}_i \{ \theta = 1 | R_i^2 = -1, h_1 = (\text{Accept}, \text{Accept}) \} = \frac{pq^2(1-q)^1}{pq^2(1-q)^1 + (1-p)(1-q)^2q^1} = \frac{2}{3} > \frac{1}{2}, \quad (8)$$

and,

$$\mathbb{P}_i \{ \theta = 1 | R_i^2 = 1, h_1 = (\text{Accept}, \text{Accept}) \} = \frac{pq^3(1-q)^0}{pq^3(1-q)^0 + (1-p)(1-q)^3q^0} = \frac{8}{9} > \frac{1}{2}. \quad (9)$$

This implies that it doesn't matter what the rewards are in the second stage, both players will select the action Accept.

The same implications hold for $t > 2$. In general, we have,

$$\mathbb{P}_i \{ \theta = 1 | R_i^t = -1, h_1 = \dots = h_{t-1} = (\text{Accept}, \text{Accept}) \} \geq \frac{pq^t(1-q)^{t-1}}{pq^t(1-q)^{t-1} + (1-p)(1-q)^tq^{t-1}} = \frac{2}{3} > \frac{1}{2}, \quad (10)$$

and,

$$\mathbb{P}_i \{ \theta = 1 | R_i^t = 1, h_1 = \dots = h_{t-1} = (\text{Accept}, \text{Accept}) \} \geq \frac{pq^{t+1}(1-q)^{t-2}}{pq^{t+1}(1-q)^{t-2} + (1-p)(1-q)^{t+1}q^{t-2}} = \frac{8}{9} > \frac{1}{2}. \quad (11)$$

Since player i only observes the past actions of player $-i$, but not the actual reward, player i makes the decision based on $t-1$ Accepts, 1 negative reward in stage 1, and at most $t-1$ positive rewards. As a result, the posterior belief that $\theta = 1$ is always at least $\frac{2}{3}$, which implies that it doesn't matter what the rewards are in each stage, both players will select the action Accept.

Note that the players will always choose Accept in a world with $\theta = -1$ in which the unique Nash Equilibrium is (Reject, Reject) in every stage just because the rewards are inconsistent with the state of the world in the first stage. In a more general setting with any p and $q > \frac{1}{2}$, there exists a finite T such that if the reward is different from θ in the first T periods for both players, both players will always choose the non-NE action in all future periods after T .

From the perspective from the adversary, they only have to perturb the reward distribution in the first stage (or first finite number of stages) to ensure that the action pair (Accept, Accept) is chosen, the same action pair will be used in all future stages.