

# Lecture 5: The Impact of Beliefs on Payoffs - Psychological Game Theory

## 1. Psychological game theory (Geanakoplos et al 1989)

First approach to internalize impact of beliefs on payoffs in games.

$n$ -player, finite extensive game form  $F$  (no payoff function yet)

$a_i$ : pure strategy of player  $i$  - associates with each information set  $h$  controlled by  $i$  an action available at  $h$ ;  $a_i \in A_i$

$\sigma_i$ : behavior strategy of player  $i$  - associates with each information set  $h$  controlled by  $i$  a probability distribution over the actions available at  $h$ .

$\Sigma_i$ : set of behavior strategies of  $i$

$\sigma = (\sigma_1, \sigma_2 \dots \sigma_n)$ : strategy profile

$\Sigma = \prod_i \Sigma_i$ : set of strategy profiles

$\sigma_{-i} = (\sigma_1, \dots, \sigma_{i-1}, \sigma_{i+1}, \dots, \sigma_n)$ : profile of strategies of all players but  $i$

$\Sigma_{-i} = \prod_{j \neq i} \Sigma_j$

$M_i^{init,1}$ : First order initial beliefs - specifies, the (uncertain) initial belief of player  $i$  about which (possibly mixed) strategies all other players  $j$  play

$\implies$

$M_i^{init,1}$ : probability distribution over  $\Sigma_{-i}$

$M_i^{init,2}$  : Second order initial belief - specifies, the (uncertain) belief of player  $i$  about the first order beliefs all other players  $j$

$M_i^{init} = \prod_{k=1}^{\infty} M_i^{init,k}$ : hierarchy of initial beliefs of player  $i$

$\Theta_i^{init}$ : set of all possible initial belief hierarchies of player  $i$

$M^{init} = (M_1^{init}, M_2^{init}, \dots, M_n^{init})$ : profile of initial belief hierarchies

Note: only initial beliefs (of all orders) are described by  $M_i^{init}$

## Payoffs

$$u_i : \Sigma \times \Theta_i^{init} \rightarrow \mathbb{R}$$

Psychological game  $\Gamma = (F, (u_i)_n)$

$\Gamma(M^{init})$ : Standard extensive form game with payoffs derived for fixed belief hierarchy  $M^{init}$ .

Un-contradictory belief hierarchy:  $M^{init}$  does not contradict strategy profile  $\sigma$ , if every player  $i$  believes with probability 1 that the other players  $-i$  play  $\sigma_{-i}$ , if every player  $i$  believes with probability 1 that every other player  $k$  believes with probability 1 that all other players  $-k$  play  $\sigma_{-k}$ , etc.

$M^{init}(\sigma)$ : The unique  $M^{init}$  that does not contradict  $\sigma$ .

Definition:  $(M^{init*}, \sigma^*)$  is a psychological Nash equilibrium of  $\Gamma$ , if:

i)  $M^{init*} = M^{init}(\sigma^*)$

ii) for all  $i$  and all  $\sigma_i \in \Sigma_i$  it holds:

$$u_i(M_i^{init*}, \sigma^*) \geq u_i(M_i^{init*}, \sigma_i, \sigma_{-i}^*)$$

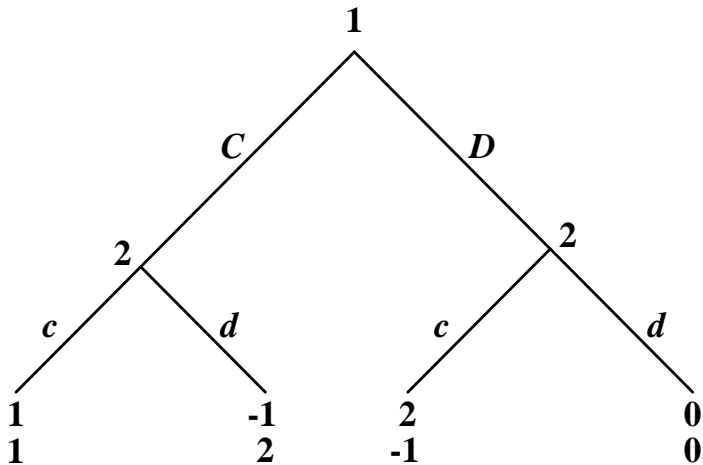
Definition:  $(M^{init*}, \sigma^*)$  is a subgame perfect psychological equilibrium of  $\Gamma$ , if:

i)  $(M^{init*}, \sigma^*)$  is a psychological Nash equilibrium of  $\Gamma$ .

ii)  $\sigma^*$  is a subgame perfect equilibrium of  $\Gamma(M^{init*})$ .

Theorem: Take a psychological game  $\Gamma$  with payoff-functions that are continuous (for a suitable topology). Then  $\Gamma$  has a subgame perfect psychological equilibrium.

Main problem with this approach:



Strategy combination  $(C, cc)$  plus consistent initial beliefs  $\implies$  1 is perceived as kind by 2 .

But: When initial beliefs are not updated, 1 is perceived as kind even after (unexpected) action *D*.

Reason: Payoffs depend only on initial beliefs not contradicting strategy combination. Off the equilibrium path, the beliefs which influence the payoffs do not get updated - "off equilibrium path, motivation remains as if deviation from equilibrium path has never occurred".

⇒ sequential reciprocity model is not a psychological game in sense of Geanakoplos et al.

To overcome these problems:

## 2. Dynamic psychological games (Battigalli and Dufwenberg 2009)

Introduces new concept of psychological games, which allows for updating of beliefs and payoffs that depend on beliefs of others.

Informal description:

Conditional probability system  $M$ : profile of belief hierarchies for any information set  $h$ .

Belief updating:

$M$  has to be such that at any information set  $h$ , all players believe (and believe that all other players believe, etc.) that all players have chosen all the actions leading to  $h$  with probability 1.

For all "parts of the strategy" not "on the way" to  $h$  : beliefs remain unchanged



$M$  is consistent with  $\sigma$ , if:

at the root, every players believes (and believes that the others believe, etc) that  $\sigma$  is played.

belief updating follows the rule above.

payoffs depend on  $M$  and on strategy combination

Definition:  $(M^*, \sigma^*)$  is a psychological sequential equilibrium, if:

i)  $M^*$  is consistent with  $\sigma^*$

ii) At any information set  $h$ , the player controlling  $h$  chooses an action which maximizes his expected payoff, given the beliefs he holds at  $h$ .

Theorem: Take a psychological game  $\Gamma$  with payoff-functions that are continuous (for a suitable topology). Then  $\Gamma$  has a psychological sequential equilibrium.

For payoffs which are constant in  $M$ , the psychological sequential equilibrium is equivalent to the sequential equilibrium (in the sense of Kreps and Wilson 1982).

Approach also allows for a general definition of rationalizability for this type of games.

### 3. Applications of Psychological Game Theory

Sequential reciprocity: Rabin 1993, Dufwenberg and Kirchsteiger 2005

Anxiety: Caplin and Leahy 2001 and 2004

Social Respect: Bernheim 1994

#### 3.1 Guilt in Games (Battigalli and Dufwenberg 2007)

Basic Idea: For each strategy of player  $j$  there is a maximum material payoff the other players can give  $j$ . If  $j$ 's actual material payoff is below this maximum, the other players are guilty vis-a-vis player  $j$ .

finite extensive form game with complete information (for simplicity)

$s_i$ : pure strategy of  $i$

$a_i$ : mixed strategy of player  $i$

$a_i(s_i)$ : probability, that pure strategy  $s_i$  is realized when mixed strategy  $a_i$  is played.

$\pi_i(s)$ : material payoff of player  $i$  for pure strategy profile  $s$ .

$b_i^h$ : first order (point) beliefs of player  $i$  about the strategies of players  $-i$  at information set  $h$

$b_i$ : first order (point) beliefs of player  $i$  about the strategies of players  $-i$  at all information sets  $h$

$b_i^0$ : first order (point) belief of player  $i$  about the strategies of the other players  $-i$  at the root of the game 0 - initial first order belief

$\pi_i(s_i, b_i^0)$ : Expected material payoff of player  $i$  if he chooses pure strategy  $s_i$  and the other players behave according to  $b_i^0$ .

$D_j(s_j, s_{-j}, b_j^0)$ : measure of how much damage is done to player  $j$  by the other players, if they choose  $s_{-j}$  instead of the expected  $b_j^0$ :

$$D_j(s_j, s_{-j}, b_j^0) = \max [0, \pi_j(s_j, b_j^0) - \pi_j(s_j, s_{-j})]$$

Player  $i$  not alone responsible for damage  $\implies$

$G_{ij}(s_j, s_i, s_{-j,i}, b_j^0)$ : measure of how much damage is done to player  $j$  by player  $i$ , if player  $i$  chooses  $s_i$  and the other players  $-j, i$  choose  $s_{-j,i}$  instead of  $b_j^0$

$$G_{ij}(s_j, s_i, s_{-j,i}, b_j^0) = D_j(s_j, s_i, s_{-j,i}, b_j^0) - \min_{\bar{s}_i \in S_i} D_j(s_j, \bar{s}_i, s_{-j,i}, b_j^0)$$

$G_{ij}(s_j, s_i, s_{-j,i}, b_j^0)$  measures  $i$ 's guilt vis-a-vis  $j$

## Preferences with simple guilt

$$u_i^{sg}(s, b_{-i}^0) = \pi_i(s) - \sum_{j \neq i} Y_{ij} G_{ij}(s_j, s_i, s_{-j,i}, b_j^0)$$

with  $Y_{ij}$  measuring  $i$ 's sensitivity of his guilt to  $j$ .

Preferences with guilt from blame (basic idea): After the game has been played,  $j$  forms an opinion on the amount of guilt  $i$  expected to have vis-a-vis  $j$  by choosing  $a_i$ , and  $j$  blames  $i$  according to this opinion. Requires second and third order beliefs.

Solution Concept: Psychological sequential equilibrium (PSE) as defined above

Impact of guilt depends of course on particular strategic situation

## General results:

For any two-player simultaneous move-game form, all PSE of the game with  $Y_{ij} = 0$  for all  $i, j$  are also PSE for all the games with  $Y_{ij} \neq 0$  with simple guilt as well as with guilt from blame.

"Equilibria of the game with selfish players are also equilibria of the games with players sensitive to guilt".

Opposite does not hold

Take any game form, and any pure strategy combination  $s$  with an outcome of material payoffs that is not strongly paretodominated by the material payoff outcome of another pure strategy combination. Then for sufficiently high sensitivities  $Y_{ij}$  there exists a PSE with  $s$  as equilibrium strategy combination.

"Any undominated outcome is an equilibrium outcome for high enough guilt sensitivities"