

Introduction to Psychological Game Theory

Lecture 12, *Experimental Econ. & Psychology*

Pierpaolo Battigalli
Bocconi University

15 October 2020

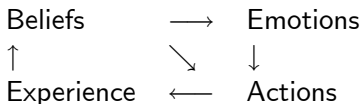
Abstract

Psychological Game Theory (PGT) is a generalization of traditional Game Theory (GT) whereby the utility of outcomes, or—more generally—of whatever actions are taken in the game, may depend on players' endogenous beliefs (i.e., beliefs that depend on the strategic analysis of the game). This generalization allows to incorporate in game theoretic analysis belief-dependent motivations related, for example, to reciprocity concerns, emotions, and image concerns.

Introduction

- Credible promises/threats and reliable communication are essential for cooperation.
- According to standard theory, credibility (incentive compatibility) is related to the value of future interactions.
- But often people cooperate, keep their word, and communicate truthfully even when this is not incentivized by future interactions.
- Emotions like guilt, anger, shame and pride can make people act against their selfish material interests in ways that are often (not always) beneficial to achieve cooperation.
- Many emotions are triggered by beliefs, including beliefs about the beliefs of others (higher-order beliefs).
- Emotions affect behavior in two ways:
 - *direct*: induced action tendencies (e.g., frustration-aggression \Rightarrow carry out threats);
 - *indirect*: anticipated feelings (valence) modify incentives (e.g., keep costly promises to avoid guilt).

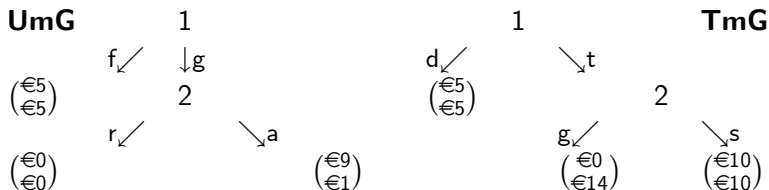
- By letting psychological utility in games depend on endogenous beliefs we can model such phenomena.



- We develop a methodology and illustrate it with some examples/applications.
- We adopt a *subjective* notion of *rationality*: (sequential) best reply to subjective beliefs, with psychological motivations.
- *Caveat*: We do not consider biases, cognitive limitations, and bounded computational abilities, nor do we model how emotions can interfere with cognition.

Stylized dilemmas with implicit threats or promises

- The Ultimatum mini-Game and the Trust mini-Game are very simple game forms representing stylized social dilemmas:



- Ultimatum mini-Game (form):** Fear of rejection may make pl. 1 choose the **fair** allocation. Is the (possibly implicit) *threat* of rejection credible? Yes, if pl. 2 *expected* the **fair** allocation and is sufficiently prone to *anger* (Battigalli, Dufwenberg & Smith, 2019).
- Trust mini-Game (form):** Hope that pl. 2 would **share** may make pl. 1 **trust**. Is the (possibly implicit) *promise* to **share** credible? Yes, if pl. 2 *thinks* pl. 1 *expected* him to **share** and is *guilt averse*.

The following is *inconsistent* with standard social preferences (e.g., inequity or lying aversion), but consistent with our framework and models:

- **Psychology:**

- desire to live up to others' expectations to avoid guilt feelings (Baumeister *et al.*, 1994; Tangney, 1995);
- frustration-aggression hypothesis (Dollard *et al.*, 1939; Frijda, 1993);
- moral behavior to avoid the feeling of shame (Tangney, 1995).

Motivations & Examples (continue)

- **Facts (casual evidence, empirics):**

- Non-returning customers give tips.
- Low offers are often rejected leaving money on the table.
- Unexpected losses by home football/soccer teams are associated with increased domestic violence (Card & Dahl, 2011) or violent crime (Munyo & Rossi 2013).

- **Facts (experimental):**

- **Trust mini-Game:** *correlation between sharing and 2nd-order beliefs of sharing; treatments effects despite no change in the traditional game form representation, which neglects information of inactive players (Charness & Dufwenberg, 2006; Tadelis, 2011; Attanasi et al. 2013).*
- **Ultimatum mini-Game:** *Rejections correlate with (manipulated) initially expected offers (Sanfey, 2009; Xiang et al., 2013, with fMRI; Aina et al., 2020).*
- **Lying/truth-telling** is *not* categorical, i.e., "*all or nothing*" (Fischbacher & Föllmi-Heusi, 2008), it *depends on the payoffs of receivers* (Gneezy, 2005; Battigalli et al., 2013) *and on exposure to passive observers* (Gneezy et al., 2016, Dufwenberg & Duf.jr. 2018).

Formal setting: one-period, sequential game forms

- **Player set:** $I_0 = I \cup \{0\}$, $i \in I$ are **personal** players, 0 is **chance**.
- **Tree of histories:** \bar{H} (*finite*, each prefix of each $h \in \bar{H}$ belongs to \bar{H} as well, including the **empty history** \emptyset).
 - Z , set of **terminal** histories/nodes (game over); H , set of **non-terminal** histories/nodes (including root \emptyset); $\bar{H} = H \cup Z$;
 $Z(h) = \{z \in Z : h \prec z\}$, terminal successors of h .
 - $\iota : H \Rightarrow I_0$ is the **active-players correspondence**;
 $H_i = \{h : i \in \iota(h)\}$, histories where i is active.
 - $A(h) = \times_{i \in \iota(h)} A_i(h)$ is the set of possible **action profiles** given h .
- **Chance probabilities:** $p_0 = (p_0(\cdot|h))_{h:0 \in \iota(h)}$, with $p_0(\cdot|h) \in \Delta(A_0(h))$.
- **Observable actions:** *active players observe earlier choices.*
- **Terminal information:** \mathcal{P}_i is a partition of Z describing what i observes *ex post* ($\mathcal{P}_i(z)$ denotes the cell containing z).
- **Material payoffs:** $\pi_i : Z \rightarrow Y_i$ ($i \in I$), e.g., monetary ($Y_i \subseteq \mathbb{E}\mathbb{R}$).

- **Trait-types:** Θ_i , set of types=personal traits of $i \in I$.
- **First-order beliefs:** set Δ_i^1 of belief systems $\alpha_i = (\alpha_i(\cdot|h))_{h \in H \cup P_i}$ s.t. $\alpha_i(\cdot|h) \in \Delta(\Theta_{-i} \times Z(h))$; given $h \prec h'$ (h prefix of h'), write $\alpha_i(\theta_{-i}, h'|h) = \alpha_i(\{\theta_{-i}\} \times Z(h')|h)$ and $\alpha_i(h'|h) = \alpha_i(\Theta_{-i} \times Z(h')|h)$, with this,
 - *chain rule:* if $(h, a', a'') \in \bar{H}$,
 $\alpha_i((h, a', a'')|h) = \alpha_i((h, a', a'')|(h, a')) \alpha_i((h, a')|h)$,
 - *self vs others indep.:* what i thinks about others' types and simultaneous actions is independent of his action; hence,
 $\alpha_i(\theta_{-i}, (h, a)|h) = \alpha_{i,j}(a_j|h) \times \alpha_{i,-i}(\theta_{-i}, a_{-i}|h)$.
- **Psy-utility:** $u_i : \Theta_i \times Z \times \Delta^1 \rightarrow \mathbb{R}$ with $\Delta^1 = \times_{j \in I} \Delta_j^1$;
 - $u_i(\theta_i, z, \alpha)$, utility of z for type θ_i given $\alpha = (\alpha_j)_{j \in I}$;
 - **note:** i does not know α_{-i} (she consults her 2nd-ord. beliefs to decide);
 - **note:** there are *private values* (standard situation in experiments).

Let $[x]^+ = \max\{x, 0\}$, $\mathbb{E}_{\alpha_i}(\pi_i) = \sum_{z \in Z} \pi_i(z) \alpha_i(z|\emptyset)$ (initially expected payoff), \mathbb{R}_+ = non-negative real n.

- **Guilt aversion**

- $u_i(\theta_i, z, \alpha) = \pi_i(z) - \sum_{j \neq i} \theta_{ij} [\mathbb{E}_{\alpha_j}(\pi_j) - \pi_j(z)]^+$,
 $\theta_i = (\theta_{ij})_{j \neq i} \in \mathbb{R}_+^{\setminus \{i\}}$,
- θ_{ij} = how much i dislikes letting j down,
- u_i does not depend on α_i ; hence, *own-plan independence* (plan = $\alpha_{i,i}$).

- **Disappointment aversion**

- $u_i(\theta_i, z, \alpha) = \pi_i(z) - \theta_i [\mathbb{E}_{\alpha_i}(\pi_i) - \pi_i(z)]^+$, $\theta_i \in \mathbb{R}_+$;
- u_i depends on the whole α_i (including $\alpha_{i,i}$); hence *own-plan dependence*.

Examples: image concerns

Fix function $V : Z \rightarrow \mathbb{R}$, then $\mathbb{E}_{\alpha_i}(V|h) = \sum_{z \in Z(h)} V(z) \alpha_i(z|h)$ denotes the conditional expectation of V given h .

- **Image concerns: good/bad behavior**

- Z_i^G (resp. Z_i^B), paths where i took **good** (resp. **bad**) actions, $\mathbf{I}_i^G : Z \rightarrow \{0, 1\}$ indicator fun. of Z_i^G (\mathbf{I}_i^B similar),
- $u_i(\theta_i, z, \alpha) = \pi_i(z) + \sum_{j \neq i} \theta_{ij} [\mathbb{E}_{\alpha_j}(\mathbf{I}_i^G | \mathcal{P}_j(z)) - \mathbb{E}_{\alpha_j}(\mathbf{I}_i^B | \mathcal{P}_j(z))]$.
- $\theta_{ij} \geq 0$, how much i cares about the opinion of j .

- **Image concerns: good/bad traits**





- $\theta_i = (\theta_i^I, \theta_i^R)$, $\theta_i^I \geq 0$: **intrinsic**-motivation trait,
- $\theta_i^R = (\theta_{ij}^R)_{j \neq i} \in \mathbb{R}_+^{\setminus \{i\}}$: **reputational**-motivation trait,
- $u_i(\theta_i, z, \alpha_j) = \pi_i(z) + \theta_i^I [\mathbf{I}_i^G(z) - \mathbf{I}_i^B(z)] + \sum_{j \neq i} \theta_{ij}^R \mathbb{E}_{\alpha_j} [\tilde{\theta}_i^I | \mathcal{P}_j(z)]$,
- where $\tilde{\theta}_i^I$ denotes a trait of i unknown to (uncertain for) j .
- u_i does not depend on α_i ; hence, *own-plan indep.* (plan= α_i).

- **Second-order beliefs:** Δ_i^2 set of 2^{nd} -order belief systems
 $\beta_i = (\beta_i(\cdot|h))_{h \in H}$ s.t.
 - $\beta_i(\cdot|h) \in \Delta(\Theta_{-i} \times Z(h) \times \Delta^1)$, the *chain rule* and *self vs others independence* hold;
 - derive 1st-order beliefs α_i by "marginalization".
- **Expected utility of actions:** For $h \in H_i$, $a_i \in A_i(h)$,
 $\bar{u}_{i,h}(a_i; \beta_i) = \mathbb{E}_{\beta_i}(u_i|h, a_i)$.
- **Local best replies:** $a_i^* \in \arg \max_{a_i \in A_i(h)} \bar{u}_{i,h}(a_i; \beta_i)$.
- **Rational planning:** Given $\alpha_{i,j}$ derived from β_i , for every $h \in H_i$,
 $\alpha_{i,i}(a_i^*|h) > 0 \Rightarrow a_i^* \in \arg \max_{a_i \in A_i(h)} \bar{u}_{i,h}(a_i; \beta_i)$ (intrapersonal equilibrium).

Consider the **Trust mini-Game** with **perfect terminal information** ($\mathcal{P}_i(z) = \{z\}$ for every $i \in I$ and $z \in Z$).






- **Exercise:**



- Let $Z_2^G = \{(t, s)\}$, $Z_2^B = \{(t, g)\}$ (sharing is good, grabbing is bad).
- Consider **image concerns** of pl. 2 **for traits**, with $\Theta_2 = \{0, \bar{\theta}_2^I\} \times \{0, \bar{\theta}_2^R\}$, $\bar{\theta}_2^I, \bar{\theta}_2^R > 0$.
- $\beta_2(\cdot|t)$ assigns probability $\frac{1}{2}$ to α_1' and α_1'' , which are such that $\mathbb{E}_{\alpha_1'}(\tilde{\theta}_2^I | (t, g)) = \mathbb{E}_{\alpha_1''}(\tilde{\theta}_2^I | (t, g)) = 0$, $\mathbb{E}_{\alpha_1'}(\tilde{\theta}_2^I | (t, s)) = \frac{1}{2}\bar{\theta}_2^I$, and $\mathbb{E}_{\alpha_1''}(\tilde{\theta}_2^I | (t, s)) = \bar{\theta}_2^I$ [α_1' deems 0 and $\bar{\theta}_2^I$ equally likely given (t, s) , α_1'' is certain of $\bar{\theta}_2^I$ given (t, s)].
- Find values of $\bar{\theta}_2^I$ and $\bar{\theta}_2^R$ such that pl. 2's best reply is to share, and values of $\bar{\theta}_2^I$ and $\bar{\theta}_2^R$ such that 2's best reply is to grab.




-  BATTIGALLI, P. (2020): *Game Theory: Analysis of Strategic Thinking*. Typescript, Bocconi University. [Downloadable from webpage, optional.]
-  BATTIGALLI, P. (2020): *Mathematical Language and Game Theory*. Typescript, Bocconi University. [Downloadable from webpage, optional.]
-  BATTIGALLI, P., C. CORRAO, AND M. M. DUFWENBERG (2019): “Incorporating Belief-Dependent Motivation in Games,” *Journal of Economic Behavior & Organization*, **167**, 185-218. [Downloadable from webpage, optional.]
-  BATTIGALLI, P., AND M. DUFWENBERG (2020): “Belief-Dependent Motivations and Psychological Game Theory,” *Journal of Economic Literature*, forthcoming.

Additional references

-  AINA, C., P. BATTIGALLI, AND A. GAMBA (2020): “Frustration and Anger in the Ultimatum Game: An Experiment,” *Games and Economic Behavior*, **122**, 150-167.
-  ATTANASI G., P. BATTIGALLI AND R. NAGEL (2013): “Disclosure of Belief-Dependent Preferences in the Trust Game,” IGER Working Paper 506, Bocconi University.
-  BATTIGALLI P., G. CHARNESS AND M. DUFWENBERG (2013): “Deception: The Role of Guilt,” *Journal of Economic Behavior & Organization*, 93, 227-232.
-  BATTIGALLI P., M. DUFWENBERG AND A. SMITH (2019): “Frustration, Aggression, and Anger in Leader-Follower Games,” *Games and Economic Behavior*, **117**, 15-39.

-  CARD, D. AND G. B. DAHL (2011): “Family Violence and Football: The Effect of Unexpected Emotional Cues on Violent Behavior,” *The Quarterly Journal of Economics*, 126, 103–143.
-  CHARNESS G. AND M. DUFWENBERG (2006): “Promises and Partnership,” *Econometrica*, 74, 1579-1601.
-  DOLLARD, J., L. DOOB, N. MILLER, O. MOWRER, AND R. SEARS (1939): *Frustration and Aggression*. Yale University Press, New Haven, CT.
-  DUFWENBERG, M., AND M. DUFWENBERG JR., (2018): “Lies in disguise—A theoretical analysis of cheating,” *Journal of Economic Theory*, 175, 248-264
-  FISCHBACHER, U., AND F. FOLLMI-HEUSI (2008): “Lies in disguise: An experimental study on cheating,” *Journal of the European Economic Association*, 11, 525–547.

-  FRIJDA, N. H. (1993): “The Place of Appraisal in Emotion,” *Cognition and Emotion*, 7, 357–387.
-  GNEEZY, U. (2005): “Deception: The role of consequences,” *American Economic Review*, 95, 384–394.
-  GNEEZY, U., A. KAJACKAITE, AND J. SOBEL (2018): “Lie Aversion and the Size of a Lie,” *American Economic Review*, 108, 419–453.
-  MUNYO, I. AND M. ROSSI (2013): “Frustration, euphoria, and violent crime,” *Journal of Economic Behavior & Organization*, 89, 136–142.

-  TADELIS S. (2011): “The Power of Shame and the Rationality of Trust,” typescript, UC Berkeley.
-  SANFEY, A. (2009): “Expectations and Social Decision-Making: Biasing Effects of Prior Knowledge on Ultimatum Responses,” *Mind and Society* 8, 93–107.
-  XIANG, T., T. LOHRENZ, AND R. MONTAGUE (2013): “Computational Substrates of Norms and Their Violations during Social Exchange,” *Journal of Neuroscience*, 33, 1099–1108.