

# Introduction to Psychological Game Theory

## Lecture 1, *Psychological GT and Experiments*

Pierpaolo Battigalli  
Bocconi University and IGER

Summer School Soleto 2025

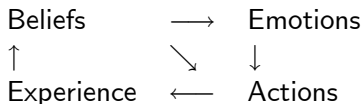
## Abstract

**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 (Battigalli & Dufwenberg 2009, 2022; Battigalli, Corrao & Dufwenberg 2019).

# Introduction

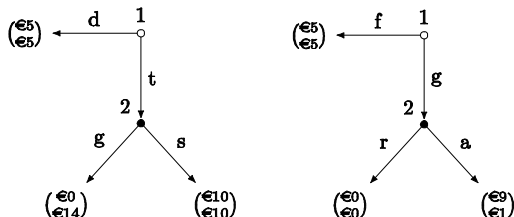
- Credible promises/threats and reliable communication are essential for cooperation.
- According to standard theory, credibility (incentive compatibility) is related to the value of future interactions.
- But often people cooperate, keep their word, and communicate truthfully even when this is not incentivized 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 **f**air allocation. Is the (possibly implicit) *threat* of rejection credible? Yes, if pl. 2 *expected* the **f**air allocation and is sufficiently prone to *anger* (Battigalli, Dufwenberg & Smith 2019).
  - **Trust mini-Game (form):** Hope that pl. 2 would **s**hare may make pl. 1 **t**rust. Is the (possibly implicit) *promise* to share credible? Yes, if pl. 2 *thinks* pl. 1 *expected* him to **s**hare and is *guilt averse*.

The following is *inconsistent* with standard social preferences (e.g., inequity or lying aversion), but consistent with the PGT 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)

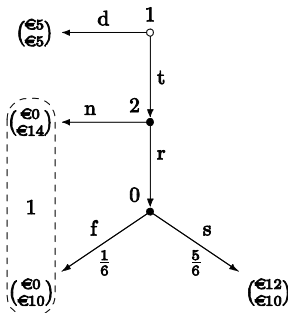
- **Facts (casual evidence, empirics,** see also survey by Charness & Fehr 2025):
  - Non-returning customers give tips.
  - Low offers are often rejected leaving money on the table.
  - Unexpected losses by home football/soccer teams are associated with increased domestic violence (Card & Dahl, 2011) or violent crime (Munyo & Rossi 2013).
- **Facts (experimental):**
  - **Trust mini-Game:** *correlation between sharing and 2<sup>nd</sup>-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. 2025).*
  - **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 2013), it *depends on the payoffs of receivers* (Gneezy, 2005; Battigalli et al., 2013) *and on exposure to*

# 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  $p_0(\cdot|h) \in \Delta(A_0(h))$ .
- **Observed 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 ( $Y_i \subseteq \mathbb{R}$ ).



## Example: TmG with random outcome



- Omitting unnecessary parentheses:

- $H_1 = \{\emptyset\}$ ,  $A_1(\emptyset) = \{d, t\}$ ,  $H_2 = \{t\}$ ,  $A_2(t) = \{n, r\}$ ,  
 $H_0 = \{(t, r)\}$ ,  $A_0(t, r) = \{f, s\}$ ,  $Z = \{d, (t, n), (t, r, f), (t, r, s)\}$ ;
- $p_0(s | (t, r)) = \frac{5}{6}$ ,  $\pi_i(d) = €5$  ( $i = 1, 2$ ),  $\pi_1(t, n) = €0 = \pi_1(t, r, f)$ ,  
 $\pi_1(t, r, s) = €12$ ,  $\pi_2(t, n) = €14$ ,  $\pi_2(t, r, f) = €10 = \pi_2(t, r, s)$ ;
- Players observe their monetary payoffs and have perfect recall:  
 $\mathcal{P}_1 = \{\{d\}, \{(t, n), (t, r, f)\}, \{(t, r, s)\}\}$ ,  $\mathcal{P}_2$  finest part. of  $Z$ .

# Formal setting: beliefs & psychological utility

- **Trait-types:**  $\Theta_i$ , set of types=personal traits of  $i \in I$ .
- **First-order beliefs:** set  $\Delta_i^1$  of belief systems  $\alpha_i = (\alpha_i(\cdot|h))_{h \in H \cup \mathcal{P}_i}$  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,  
 $\alpha_i(\theta_{-i}, (h, a)|h) = \alpha_{i,i}(a|h) \times \alpha_{i,-i}(\theta_{-i}, a_{-i}|h)$ .
- **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  $\theta_i$  given  $\alpha = (\alpha_j)_{j \in I}$ ;
  - **note:**  $i$  does not know  $\alpha_{-i}$  (she consults her 2<sup>nd</sup>-ord. beliefs to decide);
  - **note:** there are *private values* (standard situation in experiments).

# Examples: guilt and disappointment

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

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






- **Second-order beliefs:**  $\Delta_i^2$  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 1<sup>st</sup>-order beliefs  $\alpha_i$  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  $\beta_i$ , for all  $h \in H_i$  and  $a_i^* \in A_i(h)$ ,  $\alpha_{i,j}(a_i^*|h) > 0 \Rightarrow a_i^* \in \arg \max_{a_i \in A_i(h)} \bar{u}_{i,h}(a_i; \beta_i)$  (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  $\alpha'_1$  and  $\alpha''_1$ , 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  $\mathbb{E}_{\alpha''_1}(\tilde{\theta}_2^I|(t, s)) = \bar{\theta}_2^I$  [ $\alpha'_1$  deems 0 and  $\bar{\theta}_2^I$  equally likely given  $(t, s)$ ,  $\alpha''_1$  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**.






# References






-  BATTIGALLI, P. (2025): *Mathematical Language and Game Theory*. Typescript, Bocconi University. [Downloadable from webpage]
-  BATTIGALLI, P., E. CATONINI, AND N. DE VITO (2025): *Game Theory: Analysis of Strategic Thinking*. Typescript, Bocconi University. [Downloadable from webpage]
-  BATTIGALLI, P., R. CORRAO, AND M. M. DUFWENBERG (2019): "Incorporating Belief-Dependent Motivation in Games," *J. Econ. Behav. Organ.*, 167, 185-218.
-  BATTIGALLI, P., AND M. DUFWENBERG (2009): "Dynamic Psychological Games," *J. Econ. Theory*, 144, 1-35.
-  BATTIGALLI, P., AND M. DUFWENBERG (2022): "Belief-Dependent Motivations and Psychological Game Theory," *J. Econ. Lit.*, 60, 833-882.




# Additional references

-  AINA, C., P. BATTIGALLI, AND A. GAMBA (2020): "Frustration and Anger in the Ultimatum Game: An Experiment," *Games Econ. Behav.*, 122, 150-167.
-  ATTANASI G., P. BATTIGALLI, E. MANZONI, AND R. NAGEL (2025): "Disclosure of Belief-Dependent Preferences in the Trust Game," *Econ. Theory*, 10.1007/s00199-025-01645-5.
-  BATTIGALLI P., G. CHARNESS AND M. DUFWENBERG (2013): "Deception: The Role of Guilt," *J. Econ. Behav. Organ.*, 93, 227-232.
-  BATTIGALLI P., M. DUFWENBERG AND A. SMITH (2019): "Frustration, Aggression, and Anger in Leader-Follower Games," *Games Econ. Behav.*, 117, 15-39.
-  BAUMEISTER, R., A. STILLWELL, AND T. HEATHERTON (1994): "Guilt: An Interpersonal Approach". *Psychol. Bull.*, 115, 243-267.



-  CARD, D. AND G. B. DAHL (2011): “Family Violence and Football: The Effect of Unexpected Emotional Cues on Violent Behavior,” *Q. J. Econ.*, 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,” *J. Econ. Theory*, 175, 248-264
-  FISCHBACHER, U., AND F. FOLLMER-HEUSI (2013): “Lies in disguise: An experimental study on cheating,” *J. Europ. Econ. Ass.*, 11, 525–547.

-  CHARNES, G., AND E. FEHR (2025): "Social Preferences: Fundamental Characteristics and Economic Consequences," *J. Econ. Lit.*, 63, 440–514.
-  FRIJDA, N. H. (1993): "The Place of Appraisal in Emotion," *Cogn. Emot.*, 7, 357–387.
-  GNEEZY, U. (2005): "Deception: The role of consequences," *Am. Econ. Rev.*, 95, 384–394.
-  GNEEZY, U., A. KAJACKAITE, AND J. SOBEL (2018): "Lie Aversion and the Size of a Lie," *Am. Econ. Rev.*, 108, 419–453.
-  MUNYO, I. AND M. ROSSI (2013): "Frustration, euphoria, and violent crime," *J. Econ. Behav. Organ.*, 89, 136–142.

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