# Mind Games

Peter Dayan MPI for Biological Cybernetics

Andreas Hula Read Montague Nitay Alon Joe Barnby Lion Schulz Jeff Rosenchein

#### Plan

- theory of mind
- the p-beauty game
- recursive theory of mind in shopping
- recursive theory of mind in the ultimatum game
  - dynamic updating
  - defection

#### ToM as inverse reinforcement learning



#### $P(Rewards|Actions,Costs) \propto P(Action|Costs,Rewards)p(Rewards)$

Jara-Ettinger et al. (2016), Baker et al. (2017), and many more

## Recursive theory of mind: the p-Beauty Game

- Keynes; Nagel
  - choose number between 1-100
  - person closest to  $^{2}/_{3}$ rds of the mean wins prize





## The problem setting - intuition



## A three stage disinformation game



#### A three stage disinformation game



## The theory of mind depths

**DoM(-1) buyer:** The naïve buyer



**DoM(0) seller:** The inverse reinforcement learning seller (a.k.a. the Naïve Utility Calculus)



**DoM(1) buyer:** Hacking the inverse reinforcement learner

**Observer** awareness



**DoM(2) seller:** Defending against the hack *A skeptical utility calculus* 



**DoM(3) buyer:** Planning through the skeptical utility calculus



## Recap: Theory of mind levels



Gmytrasiewicz & Doshi, 2005

## Behaviors of interest



## Recap: Theory of mind levels

DoM(-1) Buyer





Reinforcement Learning (RL)



## Recap: Theory of mind levels



#### DoM(0): "Naïve Utility Calculus" Seller



## Recap: Theory of mind levels



#### DoM(1): Hacking the "Naïve Utility Calculus"



#### DoM(1): Hacking the "Naïve Utility Calculus"



## Recap: Theory of mind levels



## DoM(2): The "Skeptical Utility Calculus"





## Recap: Theory of mind levels



#### DoM(3): The "Cornered Hacker"



#### Intermediate summary



#### Intermediate summary

#### • theory of mind and pure reward maximization give rise to

- observed agents (partially) hiding their preferences
- observing agents selectively ignoring or reinterpreting the signals sent by the observed agents
- belief updating in the face of opponent
- information theoretic view



## Plan

- theory of mind
- the p-beauty game
- recursive theory of mind in shopping
- recursive theory of mind in the ultimatum game
  - dynamic updating
  - defection

## The Ultimatum Game

GDP pc GINI index NMean offer Mean reject IDV PDI AUTH TRUST COMP (1)(2)(3) (4) (5) (6) (7) (8) (9) (10)Country Austria 1 39.21 16.10 55 11 -0.050.32 6.78 12955 23.1 Bolivia 37.00 0.001721 42.0 1 0.23 Chile 34.00 6.70 23 63 1.10 5.94 4890 56.5 1 78 Ecuador 2 34.50 7.50 8 2830 46.6 France 3 40.24 30.78 71 68 -0.150.23 5.97 13918 32.7 Germany 36.70 9.52 67 35 -1.300.38 6.75 11666 30.0 1 Honduras 45.70 23.05 1385 53.7 1 Indonesia 46.63 14.63 14 78 2102 36.5 4 9843 Israel 5 41.71 17.73 54 13 35.5 Japan 3 44.73 19.27 54 -1.580.42 5.52 15105 24.9 46 -0.650.30 Yugoslavia 1 44.33 26.67 27 76 7.07 4548 31.9 44.00 27 57.5 Kenya 4.0064 914 1 1842 33.2 Mongolia 2 35.50 5.00Netherlands 38 -0.550.56 5.60 13281 31.5 2 42.25 9.24 80 Papua New-2 40.50 33.50 1606 50.9 Guinea Paraguay 51.00 0.002178 59.1 1 1.75 0.05 2092 Peru 26.00 4.8016 64 6.54 46.2 1 Romania 2 36.95 23.50 0.16 7.32 2043 28.2Slovakia -0.553 43.17 12.67 0.23 6.97 4095 19.5 Spain 26.66 29.17 51 57 0.600.34 5.70 9802 38.5 1 Sweden 35.23 18.18 71 31 -1.350.66 6.78 13986 25.01 38.2 Tanzania 37.50 19.25 27 64 534 4 UK 2 34.33 23.38 89 35 0.10 0.44 6.19 12724 32.6 US East 17945 22 40.54 17.15 91 40 1.11 0.50 6.70 40.1US West 9.41 91 17945 42.64 40 1.110.50 6.70 40.16 Zimbabwe 2 43.00 8.50 1162 56.8

Oosterbeek et al, 2004

#### The Ultimatum Game



$$u_{S}^{t}(\eta_{S}, a_{S}^{t}, a_{R}^{t}) = (1 - a_{S}^{t} - \eta_{S}) * a_{R}^{t}$$
$$u_{R}^{t}(\eta_{R}, a_{S}^{t}, a_{R}^{t}) = (a_{S}^{t} - \eta_{R}) * a_{R}^{t}$$

## Sender

- DoM(-1):
  - random: just acts randomly
  - threshold: myopic use of lower and upper bounds
    - lower: increase from rejection
    - upper: decrease from acceptance
    - softmax policy within bounds
- DoM(1):
  - pretends to be random

DoM(-1) vs DoM(0)



## DoM(1) vs DoM(0)



## DoM(2) vs DoM(-1) : 'Paranoia'

A Receiver DoM(2) Beliefs



- DoM(1) masquerades as random
  - so true randomness is hard to infer implying cost for the receiver
  - DoM(-1) threshold behaviour is atypical for DoM(1) random-pretenders
    - so infer that it is actually random
    - and so lose out!

#### Self-Protection

how can DoM(0) protect against DoM(1)?



#### Credible Threats





## Summary

- from inverse RL to deception, protection, threats
  - dynamic belief updating; non-stationarity; mixed motives
- depth of mentalizing
- critical dependence on DoM(-1) strategy
- difficulty of being too smart
- difficulty of being insufficiently smart
  - need to detect manipulation (non-likelihood based)
  - need an א-policy
- personality disorders: irritation