# Hysteresis effects of changing parameters of noncooperative games

David H. Wolpert,[1] Michael Harré,[2] Eckehard Olbrich,[3] Nils Bertschinger,[3] and Juergen Jost[3, *]

[1]*NASA Ames Research Center, MailStop 269-1, Moffett Field, CA 94035-1000*, `david.h.wolpert@nasa.gov`
[2]*Centre for the Mind, The University of Sydney, Australia*
[3]*Max Planck Institute for Mathematics in the Sciences, Inselstrasse 22, D-04103 Leipzig, Germany*
(Dated: March 2, 2011)

We adapt the method used by Jaynes to derive the equilibria of statistical physics to instead derive equilibria of bounded rational game theory. We analyze the dependence of these equilibria on the parameters of the underlying game, focusing on hysteresis effects. In particular, we show that by gradually imposing individual-specific tax rates on the players of the game, and then gradually removing those taxes, the players move from a poor equilibrium to one that is better for all of them.

## INTRODUCTION

The Maximum Entropy (Maxent) principle is an information-theoretic formalization of Occam's razor. It says that if we are given the expectation values of some functions of a system's state, then we should predict that the associated distribution is the one with minimal information (i.e., maximal entropy) consistent with those expectations [1, 2]. Maxent provides a succinct way to derive much of statistical physics [3, 4], e.g., the canonical ensemble.

Noncooperative game theory [5–8] is the foundation of conventional economics. It uses provided utility functions of a set of human "players" to predict how the players will model one another. It then uses this to predict the players' joint behavior.

Many recent applications of statistical physics to economics analyze it at a coarse-grained level, bypassing its game-theoretic foundation. Here we build on [9] and apply Maxent to game theory, thereby introducing statistical physics techniques into the foundation of economics.

In this application of Maxent, there is a separate expectation value for each player. In contrast, when applying Maxent to derive the canonical ensemble, there is a single expectation value (of the system's energy). Accordingly, rather than the canonical ensemble's single Boltzmann distribution, involving a single Hamiltonian and a single "temperature", we derive a separate Boltzmann distribution for each player, involving only that player's utility function, and a "temperature" unique to that player. The players' Boltzmann distributions are coupled, and the joint solution provides a bounded rationality version of the Nash equilibrium (NE) of game theory, where each player's inverse temperature quantifies their rationality.

We analyze the dependence of this modified NE on the parameters of the underlying game, focusing on bifurcation behavior and hysteresis effects. In particular, we show how by gradually imposing taxes on the players, and then gradually removing them, the joint behavior of the players can be moved from a poor equilibrium to a Pareto-superior one. (This can even be done if we require that the players agree to each infinitesimal change in tax rates, since each such change increases every player's expected utility.) This is particularly interesting given estimates that non-OECD countries could increase their wealth by one third by moving from their current equilibrium to a different one. Next we introduce three toy models of how a society can modify tax rates: via "socialism", a "market", or "anarchy". We then compare these three models in terms of the associated discounted sum of total utilities along the path of tax rates.

## BACKGROUND

Many different axiomatic arguments establish that the amount of syntactic information in a distribution $P(y)$ increases as the Shannon entropy of that distribution, $S(P) \equiv -\sum_y P(y)ln[\frac{P(y)}{\mu(y)}]$ [1, 3, 4], decreases. (Here $\mu$ is a measure over $y$, often interpreted as a prior probability distribution. In this paper for simplicity we assume it is uniform.) This provides a way to formalize "Occam's razor": given limited prior data concerning $P(y)$, predict $P(y)$ is the distribution with minimal information (maximum entropy) consistent with that data. This variational formalization of Occam's razor is called the maximum entropy principle (**Maxent**). When the data concerning $P(y)$ is expectation values of functions under $P$, Maxent has proven extremely accurate in domains ranging from signal processing to supervised learning [2]. Jaynes used it to derive statistical physics [3], e.g., having the data be the expected energy of a system or its expected number of particles of various types.

A finite, strategic form noncooperative game consists of a set of $N$ **players**, where each player $i$ has her own set of allowed **pure strategies** $X_i$ of size $|X_i| < \infty$. A **mixed strategy** is a distribution $q_i(x_i)$ over $X_i$. The joint distribution over $X \equiv \bigtimes_i X_i$ is $q(x) = \prod_i q_i(x_i)$, and is called a **strategy profile**.

Each player $i$ has a **utility function** $u_i : X \to \mathbb{R}$. So given strategy profile $q$, the expected utility of player $i$ is $\mathbb{E}(u_i) = \sum_x \prod_j q_j(x_j)u_i(x)$ where $q_{-i}(x_{-i}) \equiv \prod_{j \neq i} q_j(x_j)$. The **Nash equilibrium (NE)** is the strategy profile defined by having every player $i$ set $q_i$ to maximizes $i$'s expected utility, i.e., $\forall i$, $q_i = \text{argmax}_{q_i'}\left[\sum_x q_i'(x_i)q_{-i}(x_{-i})u_i(x)\right]$. In general, this set of coupled equations has multiple solutions.

A well-recognized problem of using the NE to predict real-world behavior is its assumption that every player chooses their optimal mixed strategy given the strategies of the other

players, which is called **(common knowledge) full rationality**. This assumption is violated (often badly) in many experimental settings [10, 11]. Our modified NE derived using Maxent accommodates such **bounded rationality**.

## MAXENT AND QUANTAL RESPONSE EQUILIBRIA

To predict what $q$ the players in a given $N$-player game $\Gamma$ will adopt, first pick one of the players, $i$. Consider a counter-factual situation, where $i$ has the same move space and utility function as in $\Gamma$, but rather than have a set of $N-1$ other humans set the distribution over $X_{-i}$, an inanimate stochastic system sets that distribution, to some $q_{-i}(x_{-i})$. In general, due to her limited knowledge of $q_{-i}$, limited computational power, etc., $i$ will choose a suboptimal $q_i$, i.e., $q_i \notin \text{argmax}_{p_i}[\mathbb{E}_{p_i q_{-i}}(u_i)]$. To quantify this bounded rationality, in analogy to Jaynes' derivation of the canonical ensemble, constrain $q_i$ so that $\mathbb{E}_{q_i,q_{-i}}(u_i)$ has some (nonmaximal) $K_i$ value for the given $q_{-i}$. Then Maxent says

$$q_i = argmax_{q_i'}\left[S(q_i) + \beta_i[\mathbb{E}_{q_{-i}}(u_i) - K_i]\right], \quad (1)$$

$$q_i(x_i) \propto \exp[\beta_i \mathbb{E}_{q_{-i}}(u_i \mid x_i)]. \quad (2)$$

where $\beta_i$ is the Lagrange parameter enforcing the constraint $\mathbb{E}_{q_i}(u_i) = K_i$. Note that as $\beta_i \to \infty$, $i$ becomes increasingly rational, whereas as $\beta_i \to 0$, she becomes increasingly irrational.

Next, recall that by the axioms of utility theory [12], *all* that player $i$ is concerned with in choosing her mixed strategy is the resultant expected utility. Accordingly, we presume that if the best $i$ can do is choose a particular $q_i$ when $q_{-i}$ is set by an inanimate system, she would also choose $q_i$ if she faces that same distribution $q_{-i}$ when it is set by other humans.

Generalizing, Maxent says that Eq. 2 should hold simultaneously for all $N$ players $i$, with player-specific Lagrange parameters. This gives a set of $N$ coupled non-linear equations for $q$. Brouwer's fixed point theorem [13] guarantees that set always has a solution, and in general it has more than one.[1]

This prediction for $q$ is not based on a model of bounded rational human behavior derived from experimental data. It is based on desiderata concerning the prediction process, not on a model of the system being predicted. Nonetheless, it is intriguing to note that maximizing Shannon entropy has a natural interpretation in terms of a common model of human bounded rationality, involving the cost of computation. To see this, recall that $-S(q_i)$ measures the amount of information in the distribution $q_i$. Say we equate the cost to $i$ of computing $q_i$ with this amount of information.[2] Then under the Maxent

solution, player $i$ minimizes the cost of computing her mixed strategy, subject to a constraint for the value of her expected utility. (This constraint acts as an "aspiration level".) Under this interpretation, $\beta_i$ quantifies $i$'s cost of computing $q_i$, in units of expected utility.

Alternatively, we can derive Eq. 2 with a different behavioral model that does not use Maxent. We do this by assuming that each player $i$ is purely rational, i.e., maximizes their expected utility, but does so subject to a constraint that their computational cost $-S(q_i)$ is less than some particular upper limit.

Future work involves incorporating experimental data concerning human behavior as additional constraints in the Maxent. The resultant Maxent solution could be viewed as a refined version of our two behavioral models.

Solutions for $q$ to the $N$ coupled equations given by Eq. 2 are typically called "logit Quantal Response Equilibria" (QRE) in game theory [19–22]. They have been independently suggested several times as a way to model human players [14, 23–27]. In all this earlier work the logit distribution is not derived from first principles.[3] Nor is it related to information theory, or the cost of computation. Rather typically the logit QRE has been used as an *ad hoc*, few-degree of freedom model of bounded rational play. As such it has been widely and successfully used to fit experimental data concerning human behavior.[4]

## THE SHAPE OF THE QRE SURFACE

To analyze the QRE surface of Eq. 2, we express it as a set of functional relationships, $q_i = f_i(q_{-i}, \beta_i)$, $q_{-i} = f_{-i}(q_i, \beta_{-i})$. A bifurcation may occur if for some $i$

$$\frac{\partial f_i}{\partial q_{-i}} \frac{\partial f_{-i}}{\partial q_i} \frac{\partial q_i}{\partial \beta_i} + \frac{\partial f_i}{\partial \beta_i} - \frac{\partial q_i}{\partial \beta_i} = 0 \quad (3)$$

cannot be solved for $\frac{\partial q_i}{\partial \beta_i}$, i.e., if $\det(\frac{\partial f_i}{\partial q_{-i}} \frac{\partial f_{-i}}{\partial q_i} - \text{Id}) = 0$. To illustrate this and related phenomena, we consider games between a Row and Column player, each with two pure strategies. The first is the famous "battle of the sexes" coordination game [5], where the utility functions are

$$\begin{array}{cc} 2|1 & 0|0 \\ 0|0 & 1|2 \end{array} \quad (4)$$

where the first (second) entry in each cell is the Row (Column) player's utility for the associated pure strategy profile.

---

[1] An alternative Maxent approach would use it to set the entire joint distribution $q(x) = \prod_i q_i(x_i)$ at once, rather than use it to set each $q_i$ separately and then impose self-consistency. However there are difficulties in choosing what constraints to use under this approach. See [9].

[2] Other models of the cost of computation can be found in [14–18].

---

[3] The QRE literature justifies the logit distribution by appealing to choice theory [28], where it arises if double-exponential noise is added to player utility values. However that double-exponential noise assumption is never axiomatically justified in choice theory; it is assumed for the calculational convenience that it results in the logit distribution.

[4] The logit distribution in Eq. 2 also arises in Reinforcement Learning [29–32], as a way to design artificial agents that learn from experience.
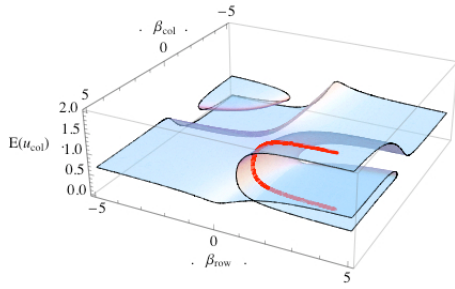
predicts that the the joint expected utility of the bargain reached is $\operatorname{argmax}_{\vec{u} \in T} [\prod_{i=1}^{N} u_i]$.

We can use the Nash bargaining concept to predict what change to $\vec{\beta}$ the players would agree to under a "market" where they bargain with one another to determine that change. To do this we fix the set of all allowed bargains to the set of all pairs $\vec{\beta}$ such that $\|\vec{\beta} - \vec{\beta}(t)\|^2 \leq 2\Delta^2$, where $\vec{\beta}(t)$ is the current joint $\beta$. We also choose $\vec{d}$ to be the joint expected utility at $\vec{\beta}(t)$. So under Nash bargaining, at each iteration $t$, the players choose the change in joint $\beta$, $\delta\vec{\beta}$, that maximizes the product

$$\left[ \mathbb{E}(u_{Row} \mid \vec{\beta}(t) + \delta\vec{\beta}) - \mathbb{E}(u_{Row} \mid \vec{\beta}(t)) \right] \times$$

$$\left[ \mathbb{E}(u_{Col} \mid \vec{\beta}(t) + \delta\vec{\beta}) - \mathbb{E}(u_{Col} \mid \vec{\beta}(t)) \right]$$

subject to $\|\delta\vec{\beta}\| \leq 2\Delta^2$. As in the other two procedures, we use first order approximations in this one, to evaluate the two differences in expected utilities.

In all three procedures the total change in $\vec{\beta}$ in any step never exceeds $\sqrt{2}\Delta$. This adiabaticity reduces the computational burden on the players, by not changing the game too much from one timestep to the next. (Similar assumptions are called comparitive statics in economics [34].)

As in standard economics, we can quantify how good a full path produced by a procedure is for society as a whole by calculating the discounted sum of future utilities along the path,

$$Q \equiv \sum_{t' > 0} (1 + \gamma)^{t - t'} \sum_{i=1}^{N} \mathbb{E}(u_i(t')) \qquad (5)$$

So we can compare the three procedures by calculating the $Q$'s for the paths they generate starting from some shared $\vec{\beta}$ at time $t = 0$. We did this for several representative initial $\vec{\beta}$'s for the surface in Fig. 3. Anarchy always did worse than the other two procedures. Those others are compared to each other in Fig. 4. When the discounting factor $\gamma$ is large (i.e., we are more concerned with near-term than long-term utility) the market procedure does better, otherwise socialism does.

All three procedures are local, looking only a single step into the future. A procedure that also considers the QRE surface's global geometry will produce better paths in general. In particular, such global information allows us to consider paths where a player loses expected utility for certain periods, but in the end all players are better off. Fig. 1 highlights such a path, along which player Column always benefits but player Row loses initially, before ultimately benefitting. (A cross-section of the expected utility of Row along the path is shown in Fig. 2. ) Note that player Row might demand compensation to agree to follow such a path where they temporarily lose expected utility, e.g. in terms of a subsidy paid for with a bond that is repaid by all players at the end of the path.

Particularly interesting issues arise when setting full paths under the market model, if the players use discounted sum
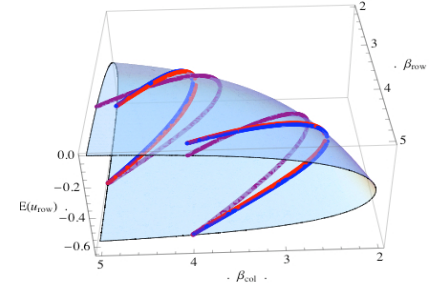


FIG. 3. A QRE surface with paths shown for the anarchy (red), socialism (blue) and market (purple) procedures. As in Fig. 1, the x and y axes are player rationalities, $\beta_{row}$ and $\beta_{col}$, and the z axis is expected utility (this time of player Row).
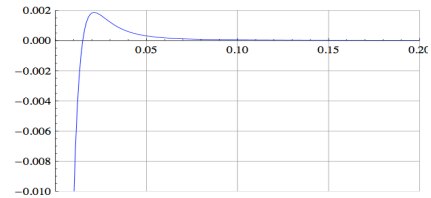


FIG. 4. The difference between the discounted sums of future expected utilities of the two players under the "socialism" and "market" procedures, plotted against the discounting factor $\gamma$.

of future utilities to value full paths. For example, say that at $t = 0$ society starts to follow a path $\vec{\beta}_0(t)$ that is a Nash bargaining solution then. Then in general, for $t' > 0$, the path $\vec{\beta}_{t'}(t)$ that is a Nash bargaining solution for full paths starting from $\vec{\beta}_0(t')$ is not a truncation of $\vec{\beta}_0(t)$ to $t > t'$. There is an inconsistency across time. This raises many interesting issues concerning binding commitments, what it means for a path chosen by bargaining to be renegotiation-proof, etc.

Multiple folds will exist for the QRE surfaces involving many kinds of game parameters, not just tax rates. Often such parameters will be set externally, perhaps in a noisy process. When this is the case, the QRE surface tells us how stable player behavior is against that external noise. For example, say the players are on the top fold of the surface in Fig. 1, with $\vec{\beta} = (2, 4)$, so the joint behavior is near an edge of the QRE surface. In this situation, small external noise may lead the players to "fall off the edge", and undergo a discontinuous jump to the lower surface. Moreover, even if the players managed to (adiababitically slowly) restore their original rationalities after such a jump, they would end up on the middle fold of the region where $\beta_{row}$ is near 2, not on the good fold

they started in. Due to this, when an economic situation exhibits such qualitative features, it may behoove society to stay away from such edges in the QRE surface, even if that lowers total expected utility.

———————

* Santa Fe Institute, 1399 Hyde Park Road, Santa Fe, NM 87501, USA

[1] T. Cover and J. Thomas, *Elements of Information Theory* (Wiley-Interscience, New York, 1991).

[2] D. Mackay, *Information Theory, Inference, and Learning Algorithms* (Cambridge University Press, 2003).

[3] E. T. Jaynes, Physical Review **106**, 620 (1957).

[4] E. T. Jaynes and G. L. Bretthorst, *Probability Theory : The Logic of Science* (Cambridge University Press, 2003).

[5] D. Fudenberg and J. Tirole, *Game Theory* (MIT Press, Cambridge, MA, 1991).

[6] R. B. Myerson, *Game theory: Analysis of Conflict* (Harvard University Press, 1991).

[7] M. Osborne and A. Rubenstein, *A Course in Game Theory* (MIT Press, Cambridge, MA, 1994).

[8] R. Aumann and S. Hart, *Handbook of Game Theory with Economic Applications* (North-Holland Press, 1992).

[9] D. H. Wolpert, in *Complex Engineered Systems: Science meets technology*, edited by D. Braha, A. Minai, and Y. Bar-Yam (Springer, 2004) pp. 262–290.

[10] C. Camerer, *Behavioral Game Theory: Experiments in Strategic Interaction* (Princeton University Press, 2003).

[11] C. Starmer, Journal of Economic Literature **38**, 332 (2000).

[12] J. von Neuman and O. Morgenstern, *Theory of Games and Economics Behavior* (Princeton university Press, 1944).

[13] C. Aliprantis and K. C. Border, *Infinite Dimensional Analysis* (Springer Verlag, 2006).

[14] D. Fudenberg and D. K. Levine, *The Theory of Learning in Games* (MIT Press, Cambridge, MA, 1998).

[15] S. Hart, Econometrica **73**, 1401 (2005).

[16] A. Rubinstein, *Modeling Bounded Rationality* (MIT press, 1998).

[17] S. Russell and D. Subramanian, Journal of AI Research , 575 (1995).

[18] M. Georgeff, B. Pell, M. Pollack, M. Tambe, and M. Wooldridge, in *Intelligent Agents V*, Lecture Notes in Computer Science, Vol. 1555 (Springer, Berlin / Heidelberg, 1999) pp. 1–10.

[19] R. D. McKelvey and T. R. Palfrey, Games and Economic Behavior **10**, 6 (1995).

[20] R. D. McKelvey and T. R. Palfrey, Japanese Economic Review **47**, 186 (1996).

[21] R. D. McKelvey and T. R. Palfrey, in *Handbook of Experimental Economics Results*, Vol. 1 (North Holland, 2008) pp. 541–548.

[22] S. P. Anderson, J. Goeree, and C. A. Holt, Southern Economic Journal **69**, 21 (2002).

[23] J. Shamma and G. Arslan, IEEE Trans. on Automatic Control **50**, 312 (2004).

[24] D. Fudenberg and D. Kreps, Games and Economic Behavior **5**, 320 (1993).

[25] J. R. Meginniss, Proceedings of the American Statistical Association, Business and Economics Statistics Section , 471 (1976).

[26] S. P. Anderson, J. Goeree, and C. A. Holt, Scandanavian journal of economics , 21 (2004), preprint title was "Stochastic game th eory: adjustment to equilibrum under noisy directional learning.

[27] S. Durlauf, Proceedings Natl. Acad. Sci. USA **96**, 10582 (1999).

[28] K. E. Train, *Discrete Choice Methods with Simulation* (Cambridge University Press, 2003).

[29] R. H. Crites and A. G. Barto, in *Advances in Neural Information Processing Systems - 8*, edited by D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo (MIT Press, 1996) pp. 1017–1023.

[30] J. Hu and M. P. Wellman, in *Proceedings of the Fifteenth International Conference on Machine Learning* (1998) pp. 242–250.

[31] D. H. Wolpert and K. Tumer, Journal of Artificial Intelligence Research **16**, 359 (2002).

[32] D. H. Wolpert, K. Tumer, and E. Bandari, Physical Review E **69** (2004).

[33] D. Wolpert and N. Kulkarni, in *Proceedings of the 2008 NASA/ESA Conference on Adaptive Hardware and Systems*, edited by A. Erdogan (IEEE Press, 2008).

[34] T. J. Kehoe, in *The New Palgrave: A Dictionary of Economics*, edited by J. Eatwell, M. Milgate, and P. Newman (Palgrave Macmillan, 1987) pp. 517–520.