

The Relationship of Economic Theory to Experiments¹

February 3, 2009

David K. Levine²

Abstract: The link between economic theory and experimental data is much tighter than is commonly supposed. Many presumed paradoxes arise because the theory is incorrectly applied. I go through several examples, emphasizing the theory as seen by a theorist. The main problem with the theory is that in some instances it lacks predictive power – I highlight where this is the case and current theoretical work designed to remedy the problem.

¹ I am grateful to NSF grant SES-03-14713 for financial support, to Drew Fudenberg and Tom Palfrey for many conversations on this topic and to Guillaume Freche for encouraging me to do this.

² John H. Biggs Distinguished Professor, Department of Economics, Washington University in St. Louis.
Email: david@dklevine.com

1. Introduction

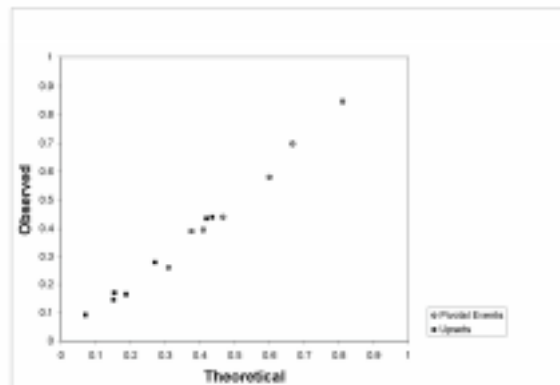
The relationship between economic theory and experimental evidence is controversial. One could easily get the impression from reading the experimental literature that economic theory has little or no significance for explaining experimental results. The point of this essay is that this is a tremendously misleading impression. Economic theory makes strong predictions about many situations, and is generally quite accurate in predicting behavior in the laboratory. Most familiar situations where the theory is thought to fail, the failure is to properly apply the theory, and not in the theory failing to explain the evidence.

That said, economic theory still needs to be strengthened to deal with experimental data: the problem is that in too many applications the theory is correct only in the sense that it has little to say about what will happen. Rather than speaking of whether the theory is correct or incorrect, the relevant question turns out to be whether it is useful or not useful. In many instances it is not useful. It does not predict how players will play in unfamiliar situations. It buries too much in individual preferences without attempting to understand how individual preferences are related to particular environments. This latter failing is especially true when it comes to preferences involving risk and time, and in preferences involving interpersonal comparisons – altruism, spite and fairness.

By way of contrast, in many circumstances equilibrium is robust to modest departures from assumptions about selfish and rational behavior. In these circumstances, the simplest form of the theory – Nash equilibrium with selfish preferences – explains the data quite well. As we shall explain – in this case predictions about aggregate behavior are quite accurate. Predictions about individual behavior are better explained by a perturbed form of Nash equilibrium – now widely known as Quantal Response equilibrium.

2. Equilibrium Theory That Works

The central theory of equilibrium in economics is that of Nash equilibrium. Let us see how that theory works in a reasonably complex voting situation. The model is adapted from Palfrey and Rosenthal [1985]. There are N voters divided into two groups, supporters of candidate A and candidate B. The number of voters is odd and divisible by three and can take on the values $\{3,9,27,51\}$. Unlike Palfrey and Rosenthal the two groups are not equal in size, rather group B is larger than group A. In the landslide treatment, there are twice as many members of B as of A. In the tossup treatment there is one more voter in group B than in group A. The voters may either vote for their preferred candidate or abstain, and the rule is simple majority. The members of the winning group receive a common prize of 105, while those in the losing group receive 5. In case of a tie, both groups receive 55. Voting is costly: the costs are private information and drawn independently and randomly on the interval $[1,55]$. Players are told the rules in a common setting, and they get to play 50 times.

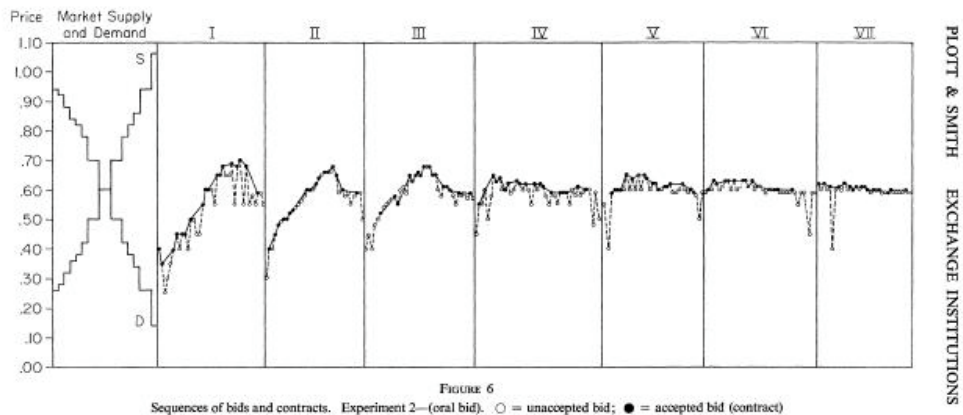


Computing the Nash equilibrium of this game is sufficiently difficult that it cannot be done by hand, nor is it possible to prove that there is a unique equilibrium. However, the equilibrium can be computed numerically, and grid searches show that there is only one equilibrium. The key to equilibrium is the probability of pivotal events: the benefit of casting a vote depends on the probability of being pivotal in an election. A good test of Nash equilibrium then is to compare the theoretical probability of a voter being pivotal – that is, of a close election, versus the empirical frequency observed in the

laboratory. The graph above from Levine and Palfrey [2007] plots the theoretical probability on the horizontal axis and the empirical frequency on the vertical axis. If the theory worked perfectly, the points should align on the forty-five degree line. They do. Despite the fact that both theoretically and from observing fifty data points it is no easy matter to infer the probability of being pivotal – the theory works nearly perfectly.

It deserves emphasis that the when we speak of “theory” here we are speaking entirely of a theoretical computation. In finding the Nash equilibrium probabilities of being pivotal *no* parameters are fit to the data: no estimation is done whatever. A pure computation is compared to live data, and the fit is nearly perfect.

The other central theory in economics besides Nash equilibrium is the competitive equilibrium of a market. In modern theory, this can be viewed as the Nash equilibrium of a mechanism in which traders reveal preferences to a market that then determines the equilibrium – with the exact details of the market clearing mechanism of no importance. Experiments on competitive equilibrium – generally in which the market clearing mechanism is a double-oral auction in real time – have been conducted many times, dating back at least to the work of Smith [1962]. The results are highly robust:



competitive equilibrium predicts the outcome of competitive market experiments with a high degree of accuracy, with experimental markets converging quickly to the competitive price. One typical picture is the history of bids in an experiment by Plott and Smith [1978] showing the convergence to the competitive equilibrium at a price of 60. Again notice that the competitive price of 60 is computed from purely theoretical considerations – no parameters are fit to the data.

This picture of data that nearly perfectly fits purely theoretical computations is true for a wide variety of experiments and is very much at odds with the viewpoint that experimental results somehow prove the theory wrong. Indeed the theory fits much better than models that must be estimated in order to fit noisy field data.

3. Equilibrium Theory that Does Not Fail

Moving past theory that predicts accurately and well, there are a set of experiments in which equilibrium – especially the refinement of subgame perfection – apparently fails badly. One such example is the ultimatum bargaining game. Here one player proposes a division of \$10 in nickels, and the second player may either accept or reject the proposal. If she accepts then the money is divided as agreed upon. If she rejects the game ends and neither player receives any money. Subgame perfection predicts that the second player should accept any positive amount, and so the first mover should get at least \$9.95. The data below from Roth et al [1991] shows that this is scarcely the case.

x	Offers	Rejection Probability
\$2.00	1	100%
\$3.25	2	50%
\$4.00	7	14%
\$4.25	1	0%
\$4.50	2	100%
\$4.75	1	0%
\$5.00	13	0%
	27	

US \$10.00 stake games, round 10

Nobody offers less than \$2.00 and most offers are for \$5.00, which is the usual amount that the first player earns. Superficially, it would be hard to imagine a greater rejection of

a theory than this. Moreover, like competitive market games, these results have been replicated many times under many conditions.

Despite appearances, theory is consistent with these results – it is the misapplication of the theory that leads to the apparent anomaly. First, the computation of the subgame equilibrium is based on the assumption that players are selfish – that they care only about their own money income. This assumption – which has nothing to do with equilibrium theory, but is merely an assertion about the nature of players’ utility functions – is clearly rejected by the data. A selfish player would not reject a positive offer – this fact is the basis for calculating the subgame perfect equilibrium. However, the data clearly shows that five out of twenty-seven positive offers are rejected. The data – not to speak of common sense – shows that many players find low offers offensive in the sense that they prefer nothing at all to a small share of the pie. A “theory” based on the assumption of selfish preferences will naturally fail to explain the data. However, there is nothing in the logic of rationality, Nash equilibrium, or subgame perfection that requires players to have selfish preferences.

It is true in the mainstream theory of competitive markets economists typically assume that people are selfish. This is not because economists believe that people are selfish – I doubt you could find a single economist who would assert that – but rather because in competitive markets it does not matter whether or not people are selfish because they have no opportunity to engage in spiteful or altruistic behavior. Consequently it is convenient for computational purposes to model people in those environments as being selfish. That should not be taken to mean that this useful modeling tool should be ported to other inappropriate environments, such as bargaining situations.

Surprisingly, even the theory of selfish preferences does not do so badly as a cursory inspection of the data might indicate. Nash equilibrium – as opposed to subgame perfection – allows any offer to be an equilibrium: it is always possible that any lower offer than the one the first player makes might be rejected with probability one, while the current offer is accepted. Nash equilibrium rules out two less obvious features of the data. It rules out a heterogeneity of offers, and it rules out offers being rejected in equilibrium (if players are truly selfish). It is a mistaken view of the theory that leads to the conclusion that this is a large discrepancy. Any theory is an idealization. Players’ exact preferences, beliefs, and so forth are never going to be known exactly to the modeler. As

a result, the only meaningful theory of Nash equilibrium is Radner's [1980] notion of epsilon equilibrium. This requires only that no player lose more than epsilon compared to the true optimum – which in practice can never be known by the players. The correct test of the goodness of fit of Nash equilibrium in experimental data is not whether the results look like a Nash equilibrium, but rather whether players losses (epsilon) is small relative to what they might have had.

The correct calculation of the departure of the facts from the theory, in other words, is to determine how much money a player who had available the experimental data could have earned, and compare it to how much that player actually earned. To the extent this is a large amount of money, we conclude the theory fits poorly. To the extent it is a small amount of money we conclude the theory fits well. This is regardless of whether the data “appear like” a Nash equilibrium or not. The key point is that allowing a small epsilon in certain games can result in a large change in equilibrium behavior. That large change does not contradict the theory of equilibrium – it is predicted by the theory of equilibrium.

For the ultimatum game, Fudenberg and Levine [1997] calculated the losses player suffered from playing less than optimal strategies given the true strategies of their opponents. Out of the \$10 on the table, players only lose on average about \$1.00 per game.

This is not the end of the story however. Nash equilibrium, as least as it is currently viewed, is supposed to be the equilibrium in which players understand their environment, including how their opponents play. It is supposed to be the outcome of a dynamic process of learning – indeed, it may accurately be described as a situation where no further learning is possible. This is important in the games in which the theory worked: in the voting experiment players played 50 times and so had a great deal of experience. Similarly, in the double oral auctions players got to participate in many auctions and equilibrium occurs only after they acquire experience. In the ultimatum game, players got to play only ten times. More important, in an extensive form game where players are informed only of the outcomes and not their opponents strategies, players would have to engage in expensive active learning to achieve a Nash equilibrium, and without a great deal of repetition and patience, they have no incentive to do so. In ultimatum bargaining in particular, the first mover can only conjecture what might

happen if she demanded more – in ten plays there is relatively little incentive or opportunity to systematically experiment with different offers to see which will be rejected or accepted. If the game were played 100 times, for example, then it would make sense to try demanding a lot to see if perhaps the opponent would be willing to accept bad offers. In 10 repetitions such a learning strategy does not make sense.

A weaker theory than Nash equilibrium – but one more suitable to the ultimatum bargaining environment – is that of self-confirming equilibrium introduced in Fudenberg and Levine [1993]. This asserts that players optimize given correct beliefs about the equilibrium path, but does not require that they know correctly what happens off the equilibrium path, as they do not necessarily observe that. This makes a difference when computing the amount of money players “lose” relative to the true optimum. In ultimatum bargaining as we observed the first movers cannot know what will happen if they demanded more. So setting making a demand that is too low is not a “knowing” error, in the sense that the player has no way to know whether it is an error or not. This leads us to compute not just the losses made by a player relative to the true optimum, but to compute how much of those losses are “knowing losses” meaning the player might reasonably know that they are making a loss. Self-confirming equilibrium is a theory that predicts that knowing losses should be low – but makes no prediction about unknowing losses.

For the ultimatum game, Fudenberg and Levine [1997] also calculated the knowing losses. On average players lose only \$0.33 per game, and this is due entirely to second players turning down positive offers – which as we noted has nothing to do with equilibrium theory at all. It is interesting to compare the impact of preferences (the spiteful play of the second players) versus that of learning (the mistaken offers of the first players). Player on average lose \$0.33 due to having preferences that are not selfish, and they lose on average \$0.67 due to the fact that they lack adequate opportunity to learn about their opponents strategies. The losses due to the deviation of preferences from the assumption of selfish behavior are considerably less than the losses due to incomplete learning.

The message here is not that theory does well with ultimatum bargaining. Rather the message is that theory is weak with respect to ultimatum bargaining – very little data in this game could be inconsistent with the theory. Rather by applying the theory

inappropriately, the conclusion was reached that the theory is wrong, while the correct conclusion is that the theory is not useful. Modern efforts in theory are quite rightly directed towards strengthening the theory – primarily by better modeling the endogenous attitudes of players towards one another as in Levine [1998], Fehr and Schmidt [1999], Bolton and Ockenfels [2000], or Gul and Pesendorfer [2004].

Another important effort is to try to capture the insight of epsilon equilibrium – that when some players deviate a little from equilibrium play, this may greatly change the incentives of other players – without losing the predictive power of Nash equilibrium. The most important effort in that direction is what has become known from the work of McKelvey and Palfrey [1995] as quantal response equilibrium. This allows for the explicit possibility that player make random errors. Specifically, if we denote by the utility that a player receives from her own pure strategy s_i and opponents mixed strategy σ_{-i} by $u_i(s_i, \sigma_{-i})$, and let $\lambda_i > 0$ be a behavioral parameter, we define the propensity with which different strategies are played by

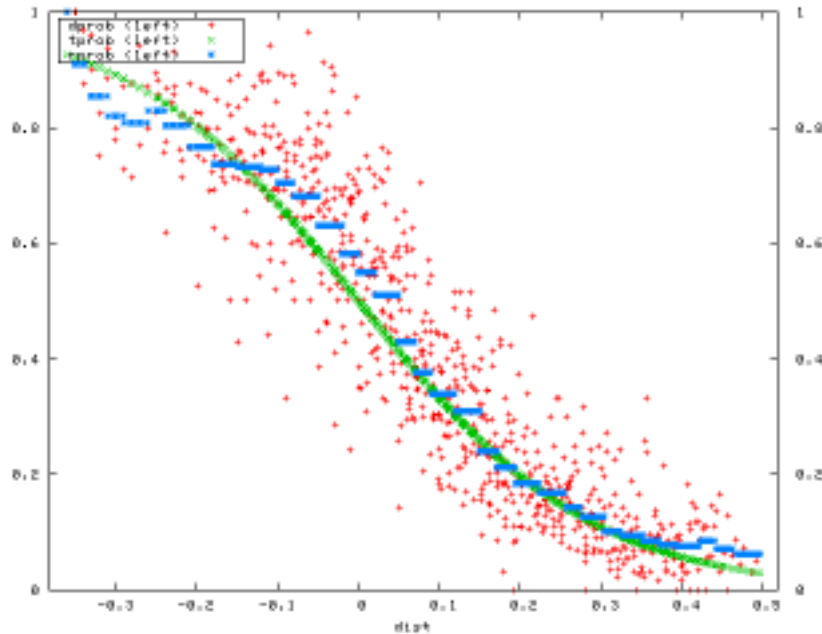
$$p_i(s_i) = \exp(\lambda_i u_i(s_i, \sigma_{-i})).$$

Quantal response theory then predicts that the mixed strategies that will be employed are given by normalizing the propensities to add to one

$$\sigma_i(s_i) = p_i(s_i) / \sum_{s_i'} p_i(s_i').$$

This theory, like Nash equilibrium, makes strong predictions. As $\lambda_i \rightarrow \infty$, these predictions in fact converge to those of Nash equilibrium. One important strength of this theory is that it allows for substantial heterogeneity at the individual level. This is important, because experimental data is quite noisy, and individual behavior generally heterogeneous.

A good example of this is in the Levine and Palfrey [2007] voting experiment described in the first section. The aggregate fit of the theory was very good, but at the

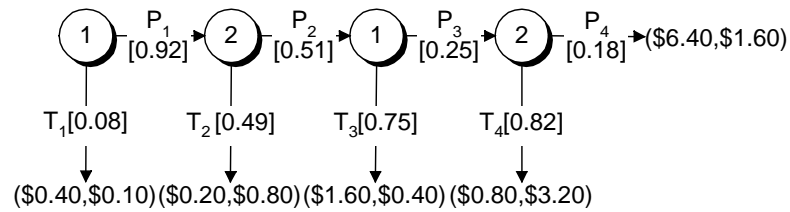


individual level, the theory fits poorly. The figure below taken from that paper shows the empirical probability with which a voter participates as a function of the loss from participating. If the loss is positive, Nash equilibrium predicts the probability of participation should be zero; if it is negative, the probability of participation should be one, and the data should align itself accordingly. The individual data in the form of red crosses and the aggregated data in the form of blue lines show that this is by no means true. When losses and gains are small, the probability of participation is relatively random – near 50%. As the loss from participating increases, the probability of participating decreases – but it hardly jumps from 1 to 0 as the threshold of indifference is crossed. However, the gradual decline seen in the data is exactly what is predicted by quantal response equilibrium. Quantal response predicts that when players are near indifferent they effectively randomize. As incentives become stronger they play more optimally. The green line shows the best fit quantal response function where λ_i is estimated from the data. As can be seen, it fits the individual level data quite well.

A key idea here is that in the aggregate quantal response equilibrium may or may not be sensitive to values of λ_i that are only moderately large. In some games, such as

the voting game, it makes little difference to aggregate behavior what λ_i is, since some voters over-voting makes it optimal for other voters to under-vote. Similarly in the market games, individual errors do not matter much at the aggregate level. The important thing is that we can always compute the quantal response equilibrium and determine how sensitive the equilibrium is to changes in λ_i .

We can tell a similar tale of poorly applied subgame perfection in the other



famous “rejection” of theory, the centipede game of McKelvey and Palfrey [1992]. The extensive form of the game is shown above. There are two players, and each may take 80% of the pot or pass, with the pot doubling at each round. Backwards induction says to drop out immediately. In fact, as the empirical frequencies in the diagram show, only 8% of players actually do that. As in ultimatum bargaining, the evidence seems to fly in the face of the theory. Again, a closer examination shows that this is not the case.

In a sense, this centipede game is the opposite of ultimatum. In ultimatum the apparent discrepancy with theory was driven by the fact that second movers are spiteful in the sense of being willing to take a small loss to punish an ungenerous opponent. In centipede the discrepancy is driven by altruism – by the willingness of a few players to suffer a small loss to provide a substantial reward to a generous opponent. The crucial empirical fact is that 18% of players will make a gift to their opponent in the final round. Notice that it costs them only \$1.60 to give a gift worth \$5.60. These gifts change the strategic nature of the game completely. With the presence of gift-givers, the true optimal strategy for each player is to stay in as long as possible. If you are the first mover stay in and hope you get lucky in the final round. If you are the second mover and make it to the final round, go ahead and grab then.

Most of the losses in centipede are actually suffered by players (foolishly misapplying subgame perfection?) who do not realize that they should stay in as long as possible, and so drop out too soon. Overall losses were computed by Fudenberg and

Levine [1997] to be about \$0.15 per player per game. However, if you drop out too soon, you never discover that there were players giving money away at the end of the game, so those losses are not knowing losses. The only knowing losses are the gifts by players in the final round. These amount to only \$0.02 per player per game. Notice that as in ultimatum, failed learning is responsible for substantially greater losses than deviation in preferences from the benchmark case of selfishness.

4. What Experiments Have Taught Us

Experimental economics has certainly taught us where the theory needs strengthening – as well as settling some long-standing methodological issues. For example, the issue of “why should we expect Nash equilibrium” has always had two answers. One answer is that players introspectively imagine that they are in the shoes of the other player, and reason their way to Nash equilibrium. This theory has conceptual problems, especially when there are multiple equilibria. It also has computational issues – for example there is a great deal of evidence that the game in which commuters choose routes to work during rush hour is in equilibrium although individual commuters certainly do not compute solutions to the game. Never-the-less in principle, players might, at least in simpler games, employ a procedure such as the Harsanyi and Selten [1988] tracing procedure. Experimental evidence, however, decisively rejects the hypothesis that the first time players are exposed to a game they manage to play a Nash equilibrium. As a result the current view – for example in Fudenberg and Levine [1998] – is that if equilibrium is reached, it is through learning. For example, the rush hour traffic game is known from the work of Monderer and Shapley [1996] to be a potential game, and such games have been shown, for example by Sandholm [2001], to be stable under a wide variety of learning procedures.

As Nash equilibrium cannot predict the outcome of one-off games, one area of theoretical research is to investigate models that can. The most promising models are the type models of Stahl and Wilson [1995]: here players are viewed as having different levels of strategic sophistication. At the bottom level, players play randomly; more sophisticated player optimize against random opponents; even more sophisticated players optimize against opponents who optimize against random opponents, and so forth. Experimental research, for example by Costa-Gomes et al [2001], shows that these

models can explain a great deal of first-time play, as well as the details of how players reason. The greatest lacuna in this literature, is that it has not yet been well tied in to a theory of learning: we have a reasonable theory of first-time play, and a reasonable theory of long-term play, but the in-between has not been solidly modeled.

The second area we highlighted above is the area of interpersonal preferences: altruism and spite. As mentioned, there are a variety of models including Levine [1998], Fehr and Schmidt [1999], Bolton and Ockenfels [2000], or Gul and Pesendorfer [2004], that attack this problem, but there is not as yet a settled theory.

There is one “emperor has no clothes” aspect of experimental research. This involves attitudes towards risk. The standard model of game theory supposes that players’ preferences can be represented by a cardinal utility function. The deficiency in this theory was highlighted by Rabin’s [2000] paradox

“Suppose we knew a risk-averse person turns down 50-50 lose \$100/gain \$105 bets for any lifetime wealth level less than \$350,000, but knew nothing about the degree of her risk aversion for wealth levels above \$350,000. Then we know that from an initial wealth level of \$340,000 the person will turn down a 50-50 bet of losing \$4,000 and gaining \$635,670.”

The point here is that in the laboratory players routinely turn down 50-50 lose \$100/gain \$105 gambles, and even more favorable gambles. Yet this is not only inconsistent with behavior in the large – it is off by (three!!) orders of magnitude. Roughly, the stakes in the laboratory are so small, that any reasonable degree of risk aversion implies risk neutrality for laboratory stakes – something strongly contradicted by the available data.

There are various possible theoretical fixes, ranging from the prospect theory of Tversky and Kahneman [1974] to the dual self approach of Fudenberg and Levine [2006], but it is fair to say that there is no settled theory, and that this is an ongoing important area of research.

5. Conclusion

The idea that experimental economics has somehow overturned years of theoretical research is ludicrous. A good way to wrap up, perhaps, is with the famous

prisoner's dilemma game. No game has been so much studied either theoretically or in the laboratory. One might summarize the widespread view as: people cooperate in the laboratory when the theory says they should not. *Caveat emptor*. The proper antidote to that view can be found in the careful experiments of Dal Bo [2005]. The proper summary of that paper is: standard Nash equilibrium theory of selfish players works quite well in predicting the laboratory behavior of players in prisoner's dilemma games.

What experimental economics has done very effectively is to highlight where the theory is weak, and there has been an important feedback loop between improving the theory – quantal response equilibrium being an outstanding example – and improving the explanation of experimental facts.

References

- Bolton, GE and A Ockenfels [2000]: “ERC: A Theory of Equity, Reciprocity, and Competition,” *American Economic Review*
- Costa-Gomes, M, VP Crawford, B Broseta [2001]: “Cognition and Behavior in Normal-Form Games: An Experimental Study,” *Econometrica*
- Dal Bo, P [2005]: “Cooperation under the Shadow of the Future: Experimental Evidence from Infinitely Repeated Games,” *American Economic Review*
- Fehr, E and KM Schmidt [1999]: “A Theory Of Fairness, Competition, and Cooperation,” *Quarterly Journal of Economics*
- Fudenberg, D and DK Levine [1997]: “Measuring Players' Losses in Experimental Games,” *Quarterly Journal of Economics*
- Fudenberg, D and DK Levine [1998] *The Theory of Learning in Games* , MIT Press
- Fudenberg, D and DK Levine [1993]: “Self-Confirming Equilibrium,” *Econometrica*
- Fudenberg, D and DK Levine [2006]: “A Dual Self Model of Impulse Control,” (with D. Fudenberg), *American Economic Review*, 96: 1449-1476
- Gul, F and W Pesendorfer [2004]: “The Canonical Type Space for Interdependent Preferences”
- Harsanyi, JC and R Selten [1988]: *A general theory of equilibrium selection in games*
- Levine, DK [1998]: “Modeling Altruism and Spitefulness in Experiments,” *Review of Economic Dynamics*
- Levine, DK and TR Palfrey [2007]: “The Paradox of Voter Participation: A Laboratory Study,” *American Political Science Review*, 101: 143-158
- McKelvey, RD and TR Palfrey [1992]: “An Experimental Study of the Centipede Game,” *Econometrica*
- McKelvey, RD and TR Palfrey [1995]: “Quantal Response Equilibria for Normal Form Games,” *Games and Economic Behavior*
- Monderer D and LS Shapley [1996]: “Potential Games,” *Games and Economic Behavior* 14:124-143
- Plott, CR and VL Smith [1978]: “An Experimental Examination of Two Exchange Institutions,” *Review of Economic Studies*

- Rabin, M [2000]: “Risk Aversion and Expected-utility Theory: A Calibration Theorem,” *Econometrica*
- Radner, R [1990]: “Collusive Behavior in Noncooperative Epsilon-Equilibria of Oligopolies with Long but Finite Lives,” *Journal of Economic Theory*
- Roth, AE, V Prasnikar, M Okuno-Fujiwara, and S Zamir [1991]: “Bargaining and Market Behavior in Jerusalem, Ljubljana, Pittsburgh, and Tokyo: An Experimental Study,” *American Economic Review*
- Sandholm, WH [2001]: “Potential Games with Continuous Player Sets,” *Journal of Economic Theory*
- Stahl, DO and PW Wilson [1995]: “On Players’ Models of Other Players: Theory and Experimental Evidence,” *Games and Economic Behavior*
- Tversky, A and D Kahneman [1974]: “Judgment under Uncertainty: Heuristics and Biases,” *Science*