# Beyond Equilibrium: Predicting Human Behaviour in Normal Form Games

by

James Robert Wright

B.Sc., Simon Fraser University, 2000

A THESIS SUBMITTED IN PARTIAL FULFILLMENT OF
THE REQUIREMENTS FOR THE DEGREE OF

MASTER OF SCIENCE

in

The Faculty of Graduate Studies

(Computer Science)

THE UNIVERSITY OF BRITISH COLUMBIA

(Vancouver)

March 2010

# Abstract

It is standard in multiagent settings to assume that agents will adopt Nash equilibrium strategies. However, studies in experimental economics demonstrate that Nash equilibrium is a poor description of human players' actual behaviour. In this study, we consider a wide range of widely-studied models from behavioural game theory. For what we believe is the first time, we evaluate each of these models in a meta-analysis, taking as our data set large-scale and publicly-available experimental data from the literature. We then propose a modified model that we believe is more suitable for practical prediction of human behaviour.

# Table of Contents

## Appendix

# List of Tables

# List of Figures

# Acknowledgements

Many people deserve recognition for their contributions to this work.

First and foremost, I would like to thank my supervisor Kevin Leyton-Brown, for introducing me to behavioural game theory in the first place, and for his continued guidance in every aspect of research.

I am indebted to several other faculty members for their insight and advice. In particular, Yoram Halevy's behavioural economics class was invaluable to me as an intensive survey of the literature, and as an introduction to how economists think. Kevin Murphy's patient advice (and book recommendations) rescued me from more than one bind. I am grateful to both of my second readers, Giuseppe Carenini and Yoram Halevy, for taking the time to review this thesis.

My fellow students have been a constant source of support and ideas. The other members of the GTDT reading group deserve special mention in this regard: Albert Xin Jiang, Baharak Rastegari, Chris Ryan, David Thompson, and Erik Zawadzki. Thanks also to Jennifer Tillett for letting me read an early draft of her thesis when I was first starting out in behavioural game theory.

Finally, thanks to my friends and family for their support and understanding, especially Sarah.

To Sarah

# Quotation

Now, a clever man would put the poison into his own goblet, because he would know that only a great fool would reach for what he was given. I am not a great fool, so I can clearly not choose the wine in front of you. But you must have known I was not a great fool; you would have counted on it, so I can clearly not choose the wine in front of me.

— William Goldman, *The Princess Bride*

# 1 Introduction

Decision making becomes more complicated when it moves beyond the single-agent case, as agents must form beliefs about the strategies (and thus, at least implicitly, the beliefs) of other agents. A standard approach is to assume that agents will adopt Nash equilibrium strategies (Nash, 1950) — that they will jointly behave in a way that ensures that each agent optimally responds to the others. This solution concept has many appealing properties; e.g., in any other strategy profile, one or more agents will regret their strategy choices. However, there are three key reasons why an agent might choose not to adopt such a strategy. First, she may face a computational limitation ("bounded rationality") that prevents her from computing a Nash equilibrium strategy, even if the game has only one. Second, even if she can compute an equilibrium, she may doubt that her opponents can or will do so. Third, when there are multiple equilibria, it is not clear which she should expect the other agents to adopt and hence whether she should play towards one herself — even if all players are perfectly rational.

Problems of bounded rationality and the need for a better normative theory are particularly acute when it comes to the *human* play of games (and, thus, to the design of agents to play against human opponents). Extensive work in experimental economics has established that human subjects often fail to adopt Nash equilibrium strategies even in very simple settings (e.g., see Stahl and Wilson, 1995; Capra et al., 1999; Goeree and Holt, 2001; Costa-Gomes et al., 2001). The relatively new field of *behavioural game theory* (BGT) aims to extend game-theoretic models to account for human behaviour by taking account of human cognitive biases and limitations (Camerer, 2003). Experimental evidence is a cornerstone of BGT, and researchers have developed many models of how humans behave in strategic situations based on experimental data. These models vary quite widely, as each tends to focus on different observed deviations from the standard equilibrium model. Furthermore, most existing work in BGT aims to understand the reasons for behaviour. As a result, there has been an emphasis on *fitting* models to experimental data, rather than attempting to *predict* experimental behaviour.

In Chapter 2, we give the formal game theoretic framework within which this thesis operates. In Chapter 3, we describe four key BGT models — level-$k$ (Costa-Gomes et al., 2001), cognitive hierarchy (Camerer et al., 2004), quantal response equilibrium (McKelvey and Palfrey, 1995), and generalized level-$k$ (Stahl and Wilson, 1994) — along with the behavioural data upon which we based our study (taken from Stahl and Wilson, 1994; Costa-Gomes et al., 2001;

Goeree and Holt, 2001; Rogers et al., 2009). We then present our two key contributions. First, in Chapter 4, we evaluate the quality of the behavioural predictions made by these four models. Second, in Chapter 5, we perform deeper analyses of the elements that make up the models. Overall, we conclude that BGT models, particularly generalized level-$k$, better predict human behaviour than both Nash equilibrium and a feature-based learning model. We also construct a model with roughly equivalent performance to generalized level-$k$ that is conceptually simpler and more parsimonious, in the sense that it assumes more homogeneous agents.

The final two sections of this chapter describe related work in artificial intelligence and economics.

## 1.1   Related work in artificial intelligence

There has been a wide range of work on designing strategies that work well in practice, given both informational and computational restrictions. This line of work is perhaps exemplified by the very influential series of Trading Agent Competitions (Wellman et al., 2007). In the annual Trading Agent Competition, researchers submit agent programs to compete against one another in a multiple market scenario based upon a real-world situation such as travel agents assembling trips for clients, or PC suppliers obtaining components and manufacturing PCs for sale to customers. These scenarios provide a standardized environment for empirically evaluating and comparing techniques for operating under computational bounds (with other, similarly bounded agents), as they are sufficiently complex that finding an optimal strategy is not analytically tractable. Unlike our work, this line of work typically does not specifically target human behaviour, but rather attempts to find algorithms that have empirically good performance in complex multiagent settings, where the other agents are assumed to also be algorithmic.

A great deal of study has also been performed on the theoretical side on alternative solution concepts. This includes set-based solution concepts that attempt to provide a stronger justification for restricting the set of actions considered. In a weak Nash equilibrium at least one of the agents is completely indifferent between the equilibrium strategy and one or more non-equilibrium strategies. Weak Nash equilibria are therefore considered very unstable, since at least one agent has no incentive not to deviate; therefore the fact that a game has a weak Nash equilibrium does not provide a strong justification for excluding non-equilibrium actions when considering what agents might do. Conitzer and Sandholm (2005) proposed an *eliminability criterion* that is intermediate in strength between domination (which may not exist in a given game) and Nash equilibrium (which may exclude more strategies than seems warranted). Similarly, Basu and Weibull (1991) introduced CURB (Closed Under Rational Behaviour) sets, which are sets of strategy profiles such that all the best responses to each strategy profile in the set are also in the set. CURB sets are a more stable solution concept than Nash equilibrium, since no agent can ever fail

to lose utility by deviating to a strategy outside the set. Benisch et al. (2006) study the complexity properties of CURB sets and provide polynomial-time algorithms for finding minimal CURB sets in arbitrary two-player normal form games.

Other theoretical work on solution concepts aims to find formally-justified solution concepts that provide more intuitive or natural solutions to games. In this vein, Halpern and Pass (2009) propose a solution concept called *iterated regret minimization*, where agents act to minimize their regret (rather than maximizing utility) in situations where they have no probabilistic beliefs about the actions of the other agents. The outcomes specified by this solution concept are arguably more intuitive than those specified by Nash equilibrium in many cases (as well as being closer to empirical outcomes in the case of the Traveller's Dilemma (Goeree and Holt, 2001; Becker et al., 2005)).

Another line of work seeks to provide worst-case guarantees on payoffs that do not depend upon other agents' rationality. Tennenholtz (2002) defines $C$-*competitive strategies* as strategies that are guaranteed to obtain at least $1/C$ as much payoff as the agent would obtain in a Nash equilibrium. He shows that $C$-competitive strategies are guaranteed to exist in many settings important to artificial intelligence, including first-price auctions. Hyafil and Boutilier (2006) tackle a similar problem from the perspective of mechanism design. They analyze a class of mechanisms that incrementally elicit type information from agents, until an acceptable bound on efficiency has been attained. Such a mechanism can simultaneously give bounds on the maximum gain that an agent achieves by reporting untruthfully; if this bound is sufficiently small, the cost of computing a profitable deviation may exceed the gain, leading to "approximate incentive compatibility."

The work described above aims to either provide compelling, non-Nash accounts of rationality, or guarantees in the face of arbitrary irrational (non-optimal) behaviour. By contrast, in this thesis we focus on models that aim to accurately model human behaviour specifically, without regard to whether or not it is optimal.

A closely related approach to our work is learning association rules between agents' actions in different games to predict how an agent will play based on its actions in earlier games (Altman et al., 2006). This requires data that identifies agents across games, and cannot make predictions for games that are not in the training dataset. By contrast, the behavioural approaches that we will explore below are able to learn from summarized datasets, and to predict actions in novel games.

## 1.2   Related work in economics

Standard economic models based on expected utility and Nash equilibrium often do not predict actual behaviour. For example, Goeree and Holt (2001) describe several games that experimental subjects play very differently based on modifications that do not affect the predictions of Nash equilibrium for one

or more players. McKelvey and Palfrey (1992) report experimental results in which all game theoretic equilibrium concepts make the same prediction, which the subjects do not follow.

Behavioural models are attempts to extend the standard theory to account for these and many other anomalies. Since the value of a behavioural theory lies in whether it accounts for actual behaviour, evaluating the fit of models against experimental data is very important. A standard technique for evaluating behavioural models is to use a $\chi^2$ statistical test on the fit of *nested* models — models where one is a generalization of the other — to experimental data. Harless and Camerer (1994) evaluate several behavioural generalizations of expected utility in this way. The various models are compared to the expected utility and expected value models only, and not pairwise, since the various behavioural generalizations are not nested. Harless and Camerer (1995) perform a similar evaluation of the performance of the intuitive and sequential refinements to Nash equilibrium.

One technique for comparing non-nested models is to compare the log likelihoods of data given the models. Camerer et al. (2004) compare their cognitive hierarchy model to Nash equilibrium in this way. Rogers et al. (2009) compare the log likelihoods of Nash equilibrium, cognitive hierarchy, and several variants of McKelvey and Palfrey's quantal response equilibrium (McKelvey and Palfrey, 1995). One drawback of comparing dataset log likelihoods directly is that it is difficult to determine how important differences are, unlike $\chi^2$ tests, where one model is either significantly better than another or not. Rogers et al. (2009) address this issue by calculating and reporting best-case and worst-case log likelihoods (the log likelihood of the actual distribution and the log likelihood of a uniform distribution respectively) alongside the model log likelihoods.[1]

A different approach than devising and comparing the fit of various models is to fit a single heterogeneous model in which different "archetypal agents" use different decision-making rules, and then test the hypothesis that various archetypes are present. Stahl and Wilson (1995) propose and test such a model with 5 boundedly rational archetypes, plus a "rational expectations" archetype that best responds to the true empirical distribution of play. Their data rejects the hypothesis that the rational expectations archetype is present, but is consistent with the boundedly rational types. Costa-Gomes et al. (2001) perform a very similar study on 4 non-strategic and 5 strategic archetypes.

Evaluating models based solely upon their fit to a specific dataset leads to a danger that some models may appear stronger than they are due to *overfitting* — in which noisy or atypical aspects of the dataset are fit — rather than due to better modelling the underlying process. One way to guard against this is to divide the experimental data into a *training set*, which is used to estimate the models' parameters, and a *test set*, which is used to test the predictive performance of the models. Stahl and Wilson (1995) manually divide their games into a test and training set. However, since they are not comparing multiple models, they are only able to use this division to show (again via a $\chi^2$

---

[1]For nested models, they also perform $\chi^2$ tests.

test) that the parameters estimated for their model on the training set are not a significantly poor fit for the test set. Similarly, Camerer et al. (2004) perform leave-one-out testing at the game level. For each game in their dataset, they fit their model on the other games and then evaluate the resulting model on the remaining game. Again, because they are not comparing multiple models, this testing is only to check that the performance of their model on out-of-sample data is reasonably robust.

# 2 Framework

In this chapter we give an overview of game theoretic concepts and notation that will be used in later chapters. Chief among these is the normal form game, which is the formalization of multiagent interaction that this study concentrates upon. We also describe three standard solution concepts from game theory that we will refer to later in the thesis.

## 2.1 Normal form games

Game theory abstracts interactive situations to mathematical objects called *games*. The simplest type of game is the *normal form game*.

**Definition 1** (Normal-form game). A *normal form game* is a tuple $(N, A, u)$, where

1. $N$ is a finite set of $n$ agents.

2. $A = A_1 \times \ldots A_n$ is a finite set of *action profiles*, where each $A_i$ is a finite set of actions available to agent $i$ for all $i \in N$. So each $a \in A$ consists of a tuple of actions $a_i$, one for each agent.

3. $u = (u_1, \ldots, u_n)$ is a tuple of utility functions $u_i : A \mapsto \mathbb{R}$, such that $u_i(a)$ is agent $i$'s utility for the outcome where each agent plays their own component of $a \in A$.

A normal form game specifies who the players are, what they can do, and what each agent's utility is for each possible combination of actions by the agents.

Each agent is assumed to play a *strategy* simultaneously. A strategy can be either a *pure strategy*, in which a single action is played deterministically, or a *mixed strategy*, in which an action is chosen stochastically according to a probability distribution over the agent's actions in $A_i$. If only one action is played with positive probability, then it is a pure strategy; if more than one action is played with positive probability, it is a mixed strategy. We use the notation $\Pi(X)$ to represent a probability distribution over the elements of a set $X$. Hence, the set of $i$'s strategies is $\Pi(A_i)$.

The combination of strategies played by each agent is referred to as the *strategy profile*. Each *strategy profile* induces a probability distribution over action profiles. The expected utility to agent $i$ of a given strategy profile is the

|       | $C$     | $D$     |
|-------|---------|---------|
| $C$   | $3,3$   | $-1,5$  |
| $D$   | $5,-1$  | $1,1$   |

Figure 2.1: Prisoner's Dilemma

average of the utilities of each possible action profile, weighted by the probability of the action profile given the strategy profile. Let $s(a)$ be the probability of action profile $a$ induced by strategy profile $s$, and $s_i(a_i)$ be the probability that agent $i$ plays action $a_i$ under strategy $s_i$. Then agent $i$'s expected utility under strategy profile $s$ is defined as

$$u_i(s) = \mathbb{E}_a u_i(a)$$
$$= \sum_{a \in A} s(a) u_i(a)$$
$$= \sum_{a \in A} \left[ \sum_{j \in N} s_j(a_j) \right] u_i(a).$$

For each agent $i \in N$, let $S_i = \Pi(A_i)$ be the set of all distributions over $A_i$. Then $S_i$ is the set of all $i$'s strategies, and $S = S_1 \times \ldots \times S_n$ is the set of all possible strategy profiles. For convenience, we use often use the notation $s = (s_i, s_{-i})$, where $s$ is a strategy profile, $s_i$ is the strategy played by $i$ in $s$, and $s_{-i}$ is the tuple of strategies played by all agents other than $i$ in $s$. Similarly, $S = S_i \times S_{-i}$, where $S_i$ is as above, and $S_{-i}$ is the cross-product of the strategy sets of all agents other than $i$.

Each agent is assumed to be *rational* and *self-interested*. That is, they have preferences that are fully representable by utility functions, and they act to maximize the expected value of their own utility function.[2]

Two-player normal form games are commonly represented by a matrix, with one row for each action of the first player and one column for each action of the second player. The utilities to each agent for each action profile are listed in the corresponding cell of the matrix. See Figure 2.1 for an example. In that game, the utility to the first agent of the action profile where agent 1 plays $D$ and agent 2 plays $C$ is 5. In the same action profile, agent 2 receives utility $-1$.

## 2.2 Solution concepts

A *solution concept* is a criterion for identifying a strategy profile or strategy profiles in a game that are in some way "interesting". A solution concept may

---

[2]Note that an agent's being "self-interested" is not equivalent to the agent's being a sociopath. Agents' utility functions are assumed to contain *all* relevant information about the agents' preferences over action profiles. In particular, this means that if agent $i$ wants agent $j$ to be happy, then this will be reflected in agent $i$'s utility function (by agent $i$ having a higher utility for outcomes in which agent $j$ is happy).

be interpreted as a prescription (i.e., that rational agents ought to play the profile identified by the solution concept), or as a prediction (i.e., that we would expect to observe the given strategy profile in actual plays of the game).

### 2.2.1 Dominant strategies

Consider again the game in Figure 2.1. If agent 2 plays $C$, then agent 1 is better off playing $D$ (for a utility of 5) than $C$ (for a utility of 3). But if agent 2 plays $D$, then agent 1 is also better off playing $D$ (for a utility of 1) than $C$ (for a utility of $-1$). So no matter what agent 2 does (including mixed strategies that include both $C$ and $D$), agent 1 has higher utility for playing $D$ than for playing $C$. In this situation we say that the pure strategy $D$ *dominates* the pure strategy $C$.

**Definition 2** (Dominance). Consider two strategies $s_i, s'_i \in S_i$. Strategy $s_i$ *dominates* $s'_i$ if

1. $\forall s_{-i} \in S_{-i} : u_i(s_i, s_{-i}) \geq u_i(s'_i, s_{-i})$, and

2. $\exists s_{-i} \in S_{-i} : u_i(s_i, s_{-i}) > u_i(s'_i, s_{-i})$.

We say that a strategy $s_i$ is *dominated* if there exists another strategy $s'_i$ that dominates $s_i$. We say that a strategy $s_i$ is *dominant* if it dominates all other strategies $s'_i \in S_i$.

A strategy profile in which each agent plays a dominant strategy is called an *equilibrium in dominant strategies*. The strategy profile $(D, D)$ is an equilibrium in dominant strategies of the Prisoner's Dilemma.

Equilibrium in dominant strategies is a solution concept, since it identifies a specific strategy profile (i.e., one in which each agent plays a dominant strategy). However, not every agent has a dominant strategy in every game (indeed, typically no agent has a dominant strategy), so it is not a solution concept that applies to every game.

### 2.2.2 Nash equilibrium

By far the most commonly used solution concept is the *Nash equilibrium*. The Nash equilibrium is defined in terms of agents' *best-response correspondences*:

**Definition 3** (Best response). A strategy $s_i \in S_i$ is a *best-response* to $s_{-i} \in S_{-i}$ if

$$u_i(s_i, s_{-i}) \geq u_i(s'_i, s_{-i}) \quad \forall s'_i \in S_i.$$

An agent $i$'s *best-response correspondence* is a function $BR_i : S_{-i} \mapsto \mathcal{P}(S_i)$ that maps from each profile of strategies by the agents other than $i$ to the set of $i$'s strategies that are best-responses to that profile. Formally,

$$BR_i(s_{-i}) = \{s^*_i \in S_i \mid \forall s'_i \in S_i : u_i(s^*_i, s_{-i}) \geq u_i(s'_i, s_{-i})\}.$$

|     | $R$      | $P$      | $S$      |
| --- | -------- | -------- | -------- |
| $R$ | $0,0$    | $-1,1$   | $1,-1$   |
| $P$ | $1,-1$   | $0,0$    | $-1,1$   |
| $S$ | $-1,1$   | $1,-1$   | $0,0$    |

Figure 2.2: Rock-Paper-Scissors

|     | $M$   | $F$   |
| --- | ----- | ----- |
| $M$ | $2,1$ | $0,0$ |
| $F$ | $0,0$ | $1,2$ |

Figure 2.3: Battle of the Sexes

**Definition 4** (Nash equilibrium)**.** A *Nash equilibrium* is a strategy profile in which every agent plays a best response to the strategies of the other agents. That is, $s^*$ is a Nash equilibrium if

$$\forall i \in N : s_i^* \in BR_i(s_{-i}^*).$$

**Theorem 1** (Nash 1950)**.** *Every game with a finite number of players and action profiles has at least one Nash equilibrium.*

Note that every equilibrium in dominant strategies is also a Nash equilibrium. So the strategy profile $(D, D)$ is a Nash equilibrium of the Prisoner's Dilemma. We refer to a Nash equilibrium in which all agents play a pure strategy as a *pure strategy Nash equilibrium*. A Nash equilibrium in which one or more players plays a mixed strategy is called a *mixed strategy Nash equilibrium*. For example, in the game of Figure 2.2, the mixed strategy equilibrium $([\frac{1}{3} : R, \frac{1}{3} : P, \frac{1}{3} : S], [\frac{1}{3} : R, \frac{1}{3} : P, \frac{1}{3} : S])$ is the only Nash equilibrium.

It is entirely possible for a game to have both pure strategy and mixed strategy Nash equilibria. In the game of Figure 2.3, there are two pure strategy Nash equilibria, namely $(M, M)$ and $(F, F)$. There is also a mixed strategy Nash equilibrium: $([\frac{2}{3} : M, \frac{1}{3} : F], [\frac{1}{3} : M, \frac{2}{3} : F])$.

A Nash equilibrium is a stable strategy profile, in the sense that no agent has an incentive to *deviate* (play a strategy other than the one prescribed by the equilibrium) given that the other players are playing their equilibrium strategies. How the agents would come to play a Nash equilibrium in the first place — particularly, but not exclusively, in games with multiple or even infinite Nash equilibria — is a separate question.

### 2.2.3 Iterative dominance

Section 2.2.1 defined the concept of *dominated strategies*. Intuitively, it makes sense to assume that no rational agent would play a dominated strategy. Therefore removing a dominated action from a game $G$ results in a smaller game $G'$ that is in some sense strategically equivalent to the original game.

When *common knowledge of rationality* exists, performing these deletions iteratively also yields a strategically equivalent game. An interactive setting has common knowledge of rationality when each agent is rational, and each agent knows that every other agent is rational, and each agent knows that every other agent knows that he is rational, and so on, as an infinite regress.

**Definition 5** (Iterative removal of dominated strategies). Let $G = (N, A, u)$ be a game in normal form. Let $G^0 = G$. Fix some $k > 0$. If $G^k$ contains no dominated actions, then $G^{k+1} = G^k$. Otherwise, $G^{k+1} = (N, A^{k+1}, u)$, where $A^{k+1} = A_1^k \times \ldots \times A_i^k \backslash \{a_i\} \times \ldots \times A_n^k$ and $a_i \in A_i^k$ is a dominated action in $G^k$. Let $\omega$ be the smallest integer such that $G^\omega = G^{\omega+1}$.

Then a pure strategy $a_i'$ *survives iterated removal of dominated strategies* in $G$ if and only if $a_i' \in A_i^\omega$.

Games that have been reduced by iterated removal of dominated strategies will always contain at least one Nash equilibrium of the original game, but may not contain all of them.

In some games, the process of iterated removal of dominated strategies does not terminate until all but one strategy has been eliminated for each player. These games are called *dominance solvable*.

**Example 1** (Traveller's Dilemma). The Traveller's dilemma is a two-player normal-form game. Each agent chooses an integer between 2 and 100. If both agents choose the same number $x$, then they both get a utility of $x$. However, if one agent chooses a smaller number $y$ than the other agent, then the agent that chose $y$ receives utility $y + 2$, and the other agent receives utility $y - 2$.

The Traveller's Dilemma turns out to be solvable by iterated removal of dominated strategies. Consider the pure strategy 100. If the other agent chooses $y \leq 98$, then the strategies 100 and 99 yield exactly the same payoff (viz $y - 2$). If the other agent chooses 99, then it is clearly better to choose 99 than to choose 100. And if the other agent chooses 100, then it is better to choose 99 (for a payoff of $99 + 2 = 101$) than to choose 100 (for a payoff of 100). So 100 is dominated by 99, and can be removed from both agents' action sets.

But in the reduced game, an identical argument results in the removal of 99; and so on until each agent's action set contains only a single strategy, 2. Thus $(2, 2)$ is the only Nash equilibrium of the Traveller's Dilemma.

As we will see in Chapter 4, human behaviour is frequently inconsistent with standard solution concepts. Indeed, in an experiment where expert game theorists played the Traveller's Dilemma for money, the best-performing strategy (97) was markedly different from the unique Nash equilibrium strategy of 2 (Becker et al., 2005).

# 3 Existing BGT Models and Experimental Data

In the first four sections of this chapter, we describe four models which are perhaps the most prominent in behavioural game theory for describing human play of normal form games. For each model, we will first give the model's assumptions and motivation. In most cases the assumptions will describe a family of models, depending upon details of parameters and additional assumptions. We will then give a concrete and formal definition of the specific instantiation of the model that we evaluate. Each model definition will include the definition of a *likelihood function*, which gives the probability of an individual observation given a particular setting of the model's parameters. We fit the model's parameters by maximizing a closely related function (see Section 4.1).

Then, in Section 3.5, we describe the publicly-available experimental data from BGT studies that we used.

## 3.1 Quantal response equilibrium

One prominent behavioural theory asserts that agents become more likely to make errors as those errors become less costly. We refer to this property as *cost-proportional errors*. This can be modeled by saying that agents are expected utility maximizers who noisily estimate each action's expected utility. Formally, agent $i$ maximizes the expected value of $\hat{u}_i(a_i, s_{-i}) = u(a_i, s_{-i}) + \epsilon_{a_i}$, where $\epsilon_{a_i}$ is a zero-mean random variable and $u_i$ is the "true" expected utility.

A *quantal response equilibrium* (QRE) (McKelvey and Palfrey, 1995) is a strategy profile $s^*$ where $\hat{u}_i(s^*) \in \arg\max_{s_i} \hat{u}(s_i, s^*_{-i})$ for all agents $i$. Different distributional assumptions for the $\epsilon_{a_i}$ terms yield different equilibrium concepts. We restrict our attention to the single-parameter *logit equilibrium*: a QRE where error terms are assumed to be independent and identically distributed across agents according to an extreme value distribution. This gives a closed form for the action probabilities:

$$s_i^*(a_i) = \frac{\exp[\lambda \cdot u_i(a_i, s^*_{-i})]}{\sum_{a_i'} \exp[\lambda \cdot u_i(a_i', s^*_{-i})]}, \tag{3.1}$$

where $\lambda$ (the *precision* parameter) indicates how sensitive agents are to utility differences. As $\lambda \to \infty$, logit equilibrium coincides with Nash equilibrium, so Nash equilibrium can be viewed as a special case of logit equilibrium.

**Model Definition 1** (QRE). Let $G = (N, A, u)$ be a normal form game, and $s^*$ be a logit equilibrium of $G$ with precision $\lambda$. Then the likelihood of an observation $(a_i, G)$ is given by

$$L^{QRE}(G, a_i \mid \lambda) = s_i^*(a_i).$$

One criticism of this solution concept is that, although (3.1) is translation-invariant, it is not scale invariant. That is, while adding some constant value to the payoffs of a game will not change its logit equilibria, multiplying payoffs by a positive constant will. This is problematic because utility functions do not themselves have unique scales (Von Neumann and Morgenstern, 1944).

## 3.2 Level-$k$

Another key idea from BGT is that humans can perform only a bounded number of *iterations of strategic reasoning*. The level-$k$ model (Costa-Gomes et al., 2001; Crawford and Iriberri, 2007) captures this idea by associating each agent $i$ with a level $k_i \in \{0, 1, 2, \ldots\}$, corresponding to the number of iterations of reasoning the agent is able to perform. A *level-0 agent* plays randomly, choosing uniformly at random from his possible actions. A *level-k agent*, for $k \geq 1$, best responds to the strategy played by level-$(k-1)$ agents. If a level-$k$ agent has more than one best response, he mixes uniformly over them.

Here we consider a particular level-$k$ model, dubbed Lk, which assumes that all agents belong to levels 0, 1, and 2. [3] Each agent with level $k > 0$ has an associated probability $\epsilon_k$ of making an "error", i.e., of playing an action that is not a best response to their beliefs. However, the agents do not account for these errors when forming their beliefs about how lower-level agents will act.

**Model Definition 2** (Lk). Let $G = (N, A, u)$ be a normal form game and $\pi_{i,k}^{Lk} \in \Pi(A_i)$ denote the distribution over actions that the Lk model predicts for a level-$k$ agent $i \in N$, and $\beta_{i,k}^{Lk} \subseteq A_i$ denote the set of actions that level-$(k+1)$ agents believe a level-$k$ agent $i$ might play. Then

$$\beta_{i,k}^{Lk} = \begin{cases} A_i & \text{if } k = 0 \\ BR_i(\beta_{-i,k-1}^{Lk}) & \text{if } k > 0 \end{cases}$$

$$\pi_{i,k}^{Lk}(a_i) = \begin{cases} |A_i|^{-1} & \text{if } k = 0, \\ (1 - \epsilon_k)/|\beta_{i,k}^{Lk}| & \text{if } k > 0, a_i \in \beta_{i,k}^{Lk}, \\ \epsilon_k/(|A_i| - |\beta_{i,k}^{Lk}|) & \text{otherwise.} \end{cases} \quad (3.2)$$

The likelihood of an observation $(a_i, G)$ is

$$L^{Lk}(G, a_i \mid \alpha_1, \alpha_2, \epsilon_1, \epsilon_2) = \sum_{k=0}^{2} \alpha_k \pi_{i,k}^{Lk}(a_i),$$

---

[3]Costa-Gomes et al. (2001) considered a model containing both level-$k$ agents and also other, non-level-$k$ agents. We study a restricted version of their model which contains level-$k$ agents only.

where $\alpha_0 = 1 - \alpha_1 - \alpha_2$.

## 3.3   Cognitive hierarchy

The cognitive hierarchy model (Camerer et al., 2004), like level-$k$, aims to model agents with heterogeneous bounds on iterated reasoning. It differs from the level-$k$ model in two ways. First, agent types do not have associated error rates; each agent best-responds perfectly to its beliefs. Second, agents best-respond to the full distribution of lower-level types, rather than only the strategy one level below. More formally, every agent again has an associated level $m \in \{0, 1, 2, \ldots\}$. Let $f$ be a probability mass function representing the distribution of the levels in the population. Level-0 agents play (typically uniformly) at random. Level-$m$ agents ($m \geq 1$) best respond to the strategies that would be played in a population described by the conditional distribution $f(j \mid j < m)$.

This family of models has the attractive feature that higher level (i.e., more cognitively capable) agents always have more accurate beliefs than lower-level agents. This is in contrast to level-$k$ models, where the beliefs of higher-than-average level agents get progressively less accurate, since each level thinks the whole population is one level smaller, which becomes progressively less true for higher levels.

Camerer et al. (2004) advocate a single-parameter restriction of the cognitive hierarchy model called *Poisson-CH*, in which the levels of agents in the population are distributed according to a Poisson distribution.

**Model Definition 3** (Poisson-CH). Let $\pi_{i,m}^{PCH} \in \Pi(A_i)$ be the distribution over actions predicted for an agent $i$ with level $m$ by the Poisson-CH model. Let $f(\cdot \mid \tau)$ be the probability mass function for a Poisson distribution with intensity parameter $\tau$. Then

$$BR_{i,m}^{\tau}(\pi^{PCH}) = \underset{a_i' \in A_i}{\arg\max} \sum_{\ell=0}^{m-1} f(\ell \mid \tau) u_i(a_i', \pi_{-i,\ell}^{PCH})$$

is the set of best responses to the truncated Poisson distribution of lower-level agents acting according to $\pi^{PCH}$, and

$$\pi_{i,m}^{PCH}(a_i) = \begin{cases} |A_i|^{-1} & \text{if } m = 0, \\ |BR_{i,m}^{\tau}(\pi^{PCH})|^{-1} & \text{if } m > 0, a_i \in BR_{i,m}^{\tau}(\pi^{PCH}), \\ 0 & \text{otherwise.} \end{cases}$$

For the sake of computability, we add the additional assumption that $m \leq 7$ for all agents.[4] Thus the likelihood for an observation $(a_i, G)$ is

$$L^{PCH}(G, a_i \mid \tau) = \frac{\sum_{m=0}^{7} f(m \mid \tau) \pi_{i,m}^{PCH}(a_i)}{\sum_{m=0}^{7} f(m \mid \tau)}.$$

_____

[4]For all maximum-likelihood estimates of $\tau$, this included 99.9% of the probability mass of the Poisson distribution.

Rogers et al. (2009) noted that cognitive hierarchy predictions often exhibit cost-proportional errors (which they call the "negative frequency-payoff deviation relationship"), even though the cognitive hierarchy model does not explicitly model this effect. This leaves open the question whether cognitive hierarchy (and level-$k$) predict well only to the extent that their predictions happen to exhibit cost-proportional errors, or whether bounded iterated reasoning captures an independent phenomenon.

## 3.4 Generalized level-$k$

Stahl and Wilson (1994) propose a rich model of strategic reasoning that combines elements of the QRE and level-$k$ models; we refer to it as the *generalized level-k model* (GLk). In GLk, agents have one of three levels, as in Lk. Each agent responds to its beliefs *quantally*, playing actions with probability proportional to the exponential of their payoffs, as in QRE. Like Lk, agents believe that the rest of the population has the next-lower type.

**Model Definition 4** (GLk). Let $\pi_{i,k}^{GLk} \in \Pi(A_i)$ denote the distribution over actions that GLk predicts for a level-$k$ agent playing as agent $i$, and $\beta_{-i,1}^{GLk} \in \Pi(A_{-i})$ denote the distribution over actions that a level-2 agent $i$ believes that the level-1 agents will play. Then

$$\pi_{i,0}^{GLk}(a_i) = |A_i|^{-1},$$

$$\pi_{i,1}^{GLk}(a_i) = \frac{\exp[\lambda_1 \cdot u_i(a_i, \pi_{-i,0}^{Glk})]}{\sum_{a_i' \in A_i} \exp[\lambda_1 \cdot u_i(a_i', \pi_{-i,0}^{Glk})]},$$

$$\beta_{i,1}^{Glk}(a_i) = \frac{\exp[\mu \cdot u_i(a_i, \pi_{-i,0}^{Glk})]}{\sum_{a_i' \in A_i} \exp[\mu \cdot u_i(a_i', \pi_{-i,0}^{Glk})]},$$

$$\pi_{i,2}^{GLk}(a_i) = \frac{\exp[\lambda_2 \cdot u_i(a_i, \beta_{-i,1}^{Glk})]}{\sum_{a_i' \in A_i} \exp[\lambda_2 \cdot u_i(a_i, \beta_{-i,1}^{Glk})]}.$$

The likelihood of an observation $(a_i, G)$ is

$$L^{GLk}(G, a_i \mid \alpha_1, \alpha_2, \lambda_1, \lambda_2, \mu) = \sum_{k=0}^{2} \alpha_k \pi_{i,k}^{GLk}(a_i),$$

where $\alpha_0 = 1 - \alpha_1 - \alpha_2$.

The $\{\alpha_1, \alpha_2\}$ parameters represent the proportions of level-1 and level-2 agents in the population. The $\{\lambda_1, \lambda_2\}$ parameters represent the precisions (as in QRE) of level-1 and level-2 agents. Finally, the $\mu$ parameter represents the level-2 agents' belief about the level-1 agents' precision (i.e., it is what the level-2 agents believe $\lambda_1$ to be).

The main difference between GLk and Lk is in the error structure. In Lk, higher-level agents believe that all lower-level agents best-respond perfectly,

although in fact every agent has some probability of making an error. In contrast, in GLk, agents are aware of the quantal nature of the lower-level agents' responses, and have a (possibly-incorrect) belief about the lower-level agents' precision.

With only 3 levels of agent, it is hard to distinguish between beliefs about proportions versus precisions of lower-level agents. The GLk model can approximately represent a level-2 belief that the population contains both level-1 and level-0 agents by a smaller value of $\mu$, since as $\mu \to 0$ the actions that level-2 agents predict become "more random". However, not all combinations of level-1 precision and relative proportions of level-0 and level-1 agents in the population can be represented in this way.

## 3.5 Experimental data

We searched the BGT literature for large-scale sets of publicly-available data from experiments on human play in normal-form games, and identified four relevant studies. Such data is relatively scarce, for two reasons: studies are difficult to conduct because they must follow human-subject protocols, and are expensive because subjects must be paid in order to align their incentives with a game's payoffs.

In Rogers et al. (2009), subjects played 17 normal form games, with payoffs denominated in pennies. In Costa-Gomes et al. (1998) subjects played 18 normal form games, with each point of payoff worth 40 cents. However, subjects were paid based on the outcome of only one randomly-selected game. Goeree and Holt (2001) presented 10 games in which subjects' behaviour was close to that predicted by Nash equilibrium, and 10 other small variations on the same games in which subjects' behaviour was *not* well-predicted by Nash equilibrium. Half of these games were normal form; the payoffs for each game were denominated in pennies. Finally, in Stahl and Wilson (1994) experimental subjects played 10 normal form games, with payoffs denominated in units worth 2.5 cents.

We represent each observation of an action by an experimental subject as a pair $(a_i, G)$, where $a_i$ is the action that the subject took when playing as player $i$ in game $G$. All games were two-player, so each single play of a game generated two observations.

| Name | Size | Contents |
|------|------|----------|
| SW94 | 400 | Data from Stahl and Wilson (1994) |
| CGCB98 | 1296 | Data from Costa-Gomes et al. (1998) |
| GH01 | 500 | Data from Goeree and Holt (2001) |
| RPC09 | 1210 | Data from Rogers et al. (2009) |
| ALL4 | 3406 | Union of above 4 datasets |
| DS | 1638 | All dominance-solvable games from ALL4 |
| GH01-T | 200 | "Treasure treatment" data from Goeree and Holt (2001) |
| GH01-C | 300 | "Contradiction treatment" data from Goeree and Holt (2001) |

Table 3.1: Datasets and their contents.

Table 3.1 lists the datasets that we considered. We built one dataset for each study, named by the source study.[5] We combined the data from all 55 games into a fifth dataset (ALL4). Finally, we also placed the 1638 observations from all the dominance-solvable games into a dataset DS; we speculated that the Lk model would predict especially well on these games, as behaviour in dominance-solvable games was part of the initial motivation behind the level-$k$ model. For the comparisons with Nash equilibrium in Section 4.2, we also separated the data from Goeree and Holt (2001) into the "treasure treatment" data and the "contradiction treatment" data, since these two groups of games were specifically chosen based upon how well or poorly Nash equilibrium predicts for them.

---

[5]Full listings of the games and observations, as well as additional detail on experimental protocols, are available in the source papers. A summary of the games and observations is available at `http://www.cs.ubc.ca/labs/lci/thesis/jrwright/experimental-data.pdf`.

# 4 Analysis of Existing BGT Models

We now describe the first of our contributions: we analyze the effectiveness of each of the BGT models we have just described on each of our datasets. We note that the papers introducing these models all performed some experimental evaluation; however, only one of them made use of data gathered by any other studies (Camerer et al., 2004), none of them made comparisons to any other models other than Nash equilibrium, and none compared generalization performance of different models. We thus believe that our evaluation of BGT models for one-shot games is the broadest ever conducted.

## 4.1 Experimental setup

To evaluate a given model on a given dataset, we performed 10 rounds of 10-fold cross-validation. Specifically, for each round, we randomly divided the dataset into 10 parts of approximately equal size, called "folds". For each of the 10 ways of selecting 9 folds from the 10, we computed the maximum likelihood estimate of the model's parameters based on those 9 folds as described in Section 4.1.1. We then determined the log likelihood of the remaining fold given the prediction. We call the average of this quantity across all 10 folds the *cross-validated log likelihood*.

The 10 rounds of cross-validation constitute a sample of size 10 from the population of possible 10-fold partitions of the dataset. Equivalently, the cross-validated log likelihoods represent a sample of size 10 from the population of possible cross-validated log likelihoods, given the dataset. By the Central Limit Theorem, the average of these cross-validated log likelihoods is approximately normally distributed. As the variance of this distribution is unknown, we use a Student's-$t$ distribution with 9 degrees of freedom to compute the 95% confidence intervals of the average cross-validated log likelihoods (e.g., see Witten and Frank, 2000).

We compared the predictive power of different behavioural models on a given dataset by comparing the average cross-validated log likelihood of the dataset under each model. We say that one model predicted significantly better than another on a given dataset when the 95% confidence intervals for the average cross-validated log likelihoods do not overlap.

We used GAMBIT (McKelvey et al., 2007) to compute QRE and to enumerate

the Nash equilibria of games. We performed computation on the `glacier` cluster of WestGrid (`www.westgrid.ca`), which consists of 840 computing nodes, each with two 3.06GHz Intel Xeon 32-bit processors and either 2GB or 4GB of RAM. In total, our study required approximately 430 CPU days of machine time, primarily for model fitting.

### 4.1.1  Maximum likelihood estimation

Each model $M$ has an associated likelihood function $L^M(G, a_i \mid \vec{\theta})$ that gives the likelihood of a single observation given a particular setting of the model's parameters $\vec{\theta}$. We assume that the observations are independent and identically distributed. This means that the likelihood of a dataset $\mathcal{D}$ is simply the product of the dataset observations:

$$L^M(\mathcal{D} \mid \vec{\theta}) = \prod_{(a_i, G) \in \mathcal{D}} L^M(G, a_i \mid \vec{\theta}).$$

As $n$ grows large, the likelihood of any particular dataset of size $n$ grows small very quickly. To prevent underflow and numerical stability problems, we therefore follow the standard practice of operating on log likelihoods:

$$LL^M(\mathcal{D} \mid \vec{\theta}) = \sum_{(a_i, G) \in \mathcal{D}} \log L^M(G, a_i \mid \vec{\theta}).$$

To compute the maximum likelihood estimate of a model $M$'s parameters $\vec{\theta}$ for a given dataset $\mathcal{D}$, we used the Nelder-Mead simplex algorithm (Nelder and Mead, 1965) to compute $\vec{\theta}^* = \arg\max_{\vec{\theta}'} LL^M(\mathcal{D} \mid \vec{\theta}')$. In general, the likelihood functions of behavioural models are not guaranteed to be convex (in fact, for some models, e.g., Poisson-CH, they can be highly non-convex or even discontinuous). To avoid computing suboptimal local maxima, we started each optimization from 200 quasi-random starting points within the model's parameter space and selected the maximum among the results.[6]

### 4.1.2  Alternatives to cross-validation

We considered two alternative evaluation techniques before settling on cross-validation: randomly dividing the data into a *training set* and a *test set* (e.g., see Bishop, 2006), and *bootstrapping* (e.g., see Efron and Tibshirani, 1993).

In the first technique, the data are randomly divided into a training set, on which model parameters are fit, and a test set, on which the fitted models' performance is evaluated. With sufficiently large datasets, overly-complex models that fit noise in the data (a phenomenon known as *over-fitting*) are likely to perform poorly in such an evaluation, as the test set will not contain the same noise as the training set. However, without distributional assumptions

---

[6]We used 200 starting points as it offered a good balance between computation cost and exhaustiveness.

about the source population, it is not possible to assess the statistical significance of differences in performance from a single training/test set evaluation. Furthermore, a simple training-data/test-data scheme works best when data are plentiful, as both the training sets and test sets need to be reasonably large in order to get reliable evaluations.

In the second technique, $K$ datapoints are sampled with replacement from the original dataset of size $K$. Since it was sampled with replacement, this new, equal-sized *pseudo-sample* may contain some of the original datapoints multiple times, while omitting others. This process is repeated $L$ times, where $L$ is some large number (often on the order of 1000). Each time a pseudo-sample is drawn, the models to be evaluated are trained on the pseudo-sample, and then evaluated on the original datapoints that were not selected. These evaluations are recorded in sorted order, and by removing $Lp/2$ evaluations from the top and bottom of this list, a $1 - p$ confidence interval for the model's performance can be estimated. One variation on this technique first randomly divides the data into training and test sets, draws pseudo-samples from the training set only, and always evaluates on the same test set.

Bootstrapping has the advantage that its statistical justification is straightforward and requires few assumptions. However, it is also extremely expensive in terms of computation (as the training process must be repeated $L$ times per dataset and model). Furthermore, it lacks statistical *power*; that is, the confidence intervals are extremely wide, and therefore models that are actually different will fail to reject the hypothesis that they are the same.

Cross-validation has important advantages over these two techniques in our setting. Average cross-validated log likelihoods admit of statistical analysis, unlike evaluations from a simple training-set/test-set scheme, and the paired $t$-test used to compare average cross-validated log likelihoods is more powerful than bootstrapping. In addition, cross-validation requires much less computing time than bootstrapping (only 100 model evaluations for 10 rounds of 10-fold cross validation, versus $L \approx 1000$ for bootstrapping). Finally, cross-validation uses data more efficiently than a simple training-set/test-set scheme, as more of the data can be used for training in each evaluation.[7]

## 4.2   Comparing to Nash equilibrium

It is straightforward to verify that the unmodified Nash equilibrium solution concept does not effectively predict the behaviour of human subjects. In 82% (45 out of 55) of the games in the ALL4 dataset, *every* Nash equilibrium assigned probability 0 to actions that were actually taken by experimental subjects. This means that treating Nash equilibrium as a prediction resulted in the entire dataset having probability 0, or infinitely negative log likelihood.

Any attempt to use Nash equilibrium for prediction must extend the solution concept to solve two problems: ensuring that no action is assigned probability

---

[7]This efficiency comes at a cost. Since the same data is sometimes used for training and sometimes for testing, cross-validated models can have some bias towards the sample set.
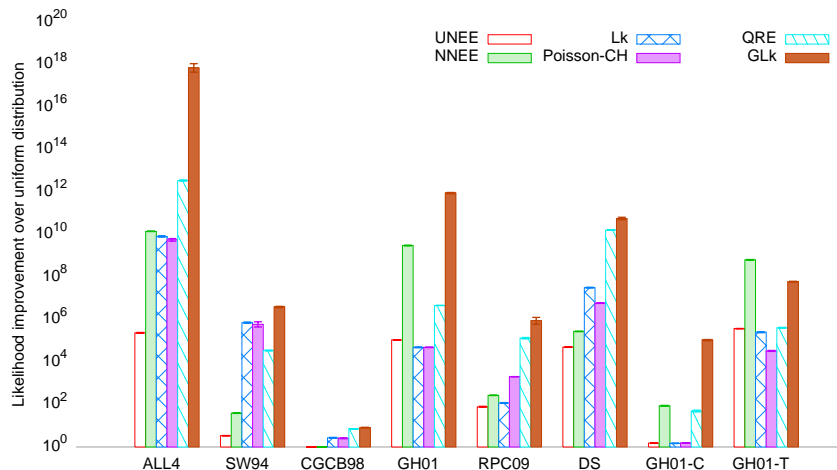
Figure 4.1: Average ratios of initial model likelihoods to random likelihoods, with 95% confidence intervals.

0, and dealing with multiple equilibria. We solved the first problem by adding a parameter ($\epsilon$) representing a proportion of the population that chooses actions at random. We solved the equilibrium selection problem in two ways. The first was to take the average over the predictions of every Nash equilibrium. This is equivalent to having a uniform prior over the equilibria of a game. We called the resulting model *uniform Nash equilibrium with error* (UNEE). Our second solution was to nondeterministically select the Nash equilibrium that was most consistent with the full dataset. We call the resulting model *nondeterministic Nash equilibrium with error* (NNEE). Clearly this model could not be used in practice, as it relies upon "peeking" at the full dataset. It can be understood as a best-case scenario for Nash equilibrium in which the equilibrium-selection problem is solved perfectly.

Figure 4.1 reports the results of the initial comparisons of UNEE, NNEE, and our four BGT models. For each model and each dataset, we give the factor by which the dataset is more likely according to the model's prediction than it would be according to a uniform random prediction. Thus, for example, the ALL4 dataset is approximately $10^{12}$ times more likely according to QRE's prediction than it is according to a uniform random prediction.

On the GH01-T dataset, which was explicitly selected to contain games that were well-described by Nash equilibrium, NNEE had the best prediction performance. However, UNEE predicted significantly worse than both GLk and QRE. As expected, both UNEE and NNEE had worse performance on the GH01-C dataset. More surprisingly, so did the BGT models (although GLk degraded the least, and hence predicted much better than the other models).

UNEE predicted significantly worse than every BGT model on every non-GH01 dataset. Similarly, NNEE predicted significantly worse than GLk on

GH01-C and the combined GH01 dataset, and predicted significantly worse than both GLk and QRE on every non-GH01 dataset. This is unambiguous evidence that BGT models better predict human behaviour than Nash equilibrium.

## 4.3  Comparing BGT models

Referring again to Figure 4.1, we see in most datasets that the model based on cost-proportional errors (QRE) predicted significantly better than the two models based on bounded iterated reasoning. Surprisingly, this was also true in the DS dataset, even though bounded iterated reasoning was largely motivated by observed human behaviour in dominance-solvable games.

In contrast, models based on bounded iterated reasoning (Lk and Poisson-CH) outperformed QRE on SW94. This suggests that bounded iterated reasoning and cost-proportional errors capture distinct underlying phenomena. If that were true, then one would expect that models incorporating both components would predict better than models that incorporate only one or the other of them. This is indeed the case, as GLk incorporates both components and generally outperforms the single-component models. Overall, GLk was the strongest of the BGT models. With the sole exception of NNEE in GH01-T, GLk predicted significantly better than all models in all datasets.

# 5 Deeper Analysis of BGT Models

Based on several questions from our initial evaluation, we then performed a deeper analysis of our four BGT models. To investigate each question, we constructed modified models, and compared their performance to the original models.

## 5.1 Poisson distributions in cognitive hierarchy

Our first question was whether it is reasonable to assume that agent levels have a Poisson distribution in the cognitive hierarchy model. At the best-fitting parameter values for ALL4, this would imply that roughly 75% of agents are level-0, which we consider implausible. We hypothesized that a cognitive hierarchy model assuming some other distribution would better fit the data. To test this hypothesis, we constructed a 4-parameter cognitive hierarchy model (CH4), in which each agent was assumed to have level $m \leq 4$, but where the distributional form was otherwise unrestricted. If CH4 consistently makes better predictions than Poisson-CH, then we can conclude that the restriction to the Poisson distribution was harmful.

**Model Definition 5** (CH4). Analogously to Poisson-CH, we define

$$BR_{i,m}^{\alpha}(\pi^{CH4}) = \arg\max_{a_i' \in A_i} \sum_{\ell=0}^{m-1} \alpha_\ell u_i(a_i', \pi_{-i,\ell}^{CH4})$$

as the best response to the truncated distribution of lower-level agents, with $\alpha_0 = 1 - \alpha_1 - \alpha_2 - \alpha_3 - \alpha_4$. This gives

$$\pi_{i,m}^{CH4}(a_i) = \begin{cases} |A_i|^{-1} & \text{if } m = 0, \\ |BR_{i,m}^{\alpha}(\pi^{CH4})|^{-1} & \text{if } m > 0, a_i \in BR_{i,m}^{\alpha}(\pi^{CH4}), \\ 0 & \text{otherwise.} \end{cases}$$

The likelihood of an individual observation is

$$L^{CH4}(G, a_i \mid \alpha_1, \alpha_2, \alpha_3, \alpha_4) = \sum_{m=0}^{4} \alpha_m \pi_{i,m}^{CH4}(a_i).$$

Figure 5.1: Average likelihood ratios between CH4 and Poisson-CH models, with 95% confidence intervals.



Figure 5.2: Average likelihood ratios between Lk4 and Lk models, with 95% confidence intervals.

Figure 5.3: Average likelihood ratios between normalized QRE and unnormalized QRE models, with 95% confidence intervals.

Figures 5.1–5.3 report the evaluations of the modified models considered in this section, expressed as a ratio between the likelihood of the modified model and the corresponding original model. In Figure 5.1 we can see that the ALL4 dataset is approximately 100 times more likely according to the CH4 model's prediction than it is according to Poisson-CH. CH4 predicted significantly better than the Poisson-CH model on most datasets, and never significantly worse. Overall, we conclude that the assumption of Poisson-distributed agent levels was unhelpful in the cognitive hierarchy model.

Interestingly, although this question was motivated by the high proportion of level-0 agents predicted by the Poisson-CH model, the CH4 model still predicts that 63% of the agents are level-0. This is a substantial reduction compared to Poisson-CH, which predicts 75%, but it is still a surprisingly large proportion.

## 5.2 Are higher-level agents helpful in level-$k$ models?

Both the generalized level-$k$ and level-$k$ models assume that all agents have level $k \leq 2$. Our second question was whether a richer model that allowed for higher-level agents would have better predictive power. To explore this question, we constructed a level-$k$ model with $k \in \{0, 1, 2, 3, 4\}$ (Lk4). We hypothesized that the Lk4 model would have better predictive power than the Lk model.

**Model Definition 6** (Lk4). The Lk4 model is identical to the Lk model, except

24

with more levels. The likelihood function is

$$L^{Lk4}(G, a_i \mid \alpha_1, \alpha_2, \alpha_3, \alpha_4, \epsilon_1, \epsilon_2, \epsilon_3, \epsilon_4) = \sum_{k=0}^{4} \alpha_k \pi_{i,m}^{Lk}(a_i),$$

where $\pi_{i,m}^{Lk}$ is defined in equation (3.2).

As reported in Figure 5.2, the Lk4 model predicted significantly better than the Lk model on all datasets except CGCB98, where there was no significant difference between the two models. However, these differences were small in every case, in spite of the fact that Lk4 has twice as many parameters as Lk.[8] Overfitting does not appear to have influenced these results, as the ratios of test to training log likelihoods were not significantly different between the Lk and Lk4 models. This suggests that there is not a large proportion of higher-level agents that are well-described by the level-$k$ model.

## 5.3   Does payoff scaling matter?

Our third question was whether the payoffs in the different games in the dataset were in appropriate units. The level-$k$ and Poisson-CH models are based only on the best-response structure of the games, and are therefore independent of the units used for payoffs. However, both QRE and generalized level-$k$ are scale dependent. If the precision parameter is held fixed, then QRE will predict differently for two identical games whose payoffs are expressed in different units. When considering a single setting this is not a concern, because the precision parameter can contain a factor to scale a game to appropriately-sized units. However, when data is combined from multiple studies in which payoffs are expressed on different scales, we might worry that the single precision parameter is insufficient to compensate for QRE's scale dependence.

We proposed two hypotheses to explore this question. The first was that subjects were concerned only with relative scales of payoff differences within individual games. To test this hypothesis, we constructed a model (NQRE) that normalizes payoffs within a game to lie in the interval $[0, 1]$ and then predicts based on a QRE of the normalized game. The second was that subjects were concerned with the expected monetary value of their payoffs. To test this hypothesis, we constructed a model (CNQRE) that normalizes payoffs to be denominated in expected cents. If either normalized model consistently predicts significantly better than QRE, then we have evidence for its associated hypothesis.

**Model Definition 7** (NQRE)**.** Let $G = (N, a, u)$ be a normal form game. Define $n(G) = (N, A, v)$, where

$$v_i(a) = \frac{u_i(a) - \min_{a' \in A} u_i(a')}{(\max_{a' \in A} u_i(a')) - (\min_{a' \in A} u_i(a'))}.$$

---

[8]Note that the scales on each figure differ from one another. In particular, Figure 5.2 has a scale that ranges from $10^0$ to $10^1$, compared to Figure 5.1, whose scale ranges from $10^{-1}$ to $10^4$.

Then given a game $G$ and a precision $\lambda$, NQRE predicts according to $s^*$, where $s^*$ is a logit equilibrium of $n(G)$ with precision $\lambda$. The likelihood of an observation $(a_i, G)$ is thus

$$L^{NQRE}(G, a_i \mid \lambda) = L^{QRE}(n(G), a_i \mid \lambda).$$

**Model Definition 8** (CNQRE)**.** Let $G = (N, A, u)$ be a normal-form game. Define $cn(G) = (N, A, v)$, where $v$ is a utility function that returns the expected monetary value of an outcome in cents. Then given a game $G$ and a precision $\lambda$, CNQRE predicts according to $s^*$, where $s^*$ is a logit equilibrium of $cn(G)$ with precision $\lambda$. The likelihood of an observation $(a_i, G)$ is thus

$$L^{CNQRE}(G, a_i \mid \lambda) = L^{QRE}(cn(G), a_i \mid \lambda).$$

Figure 5.3 reports the likelihood ratio between the modified QRE models and QRE. Both NQRE and CNQRE performed worse than the original unnormalized QRE on every dataset except for SW94, where the improvement was very small (although significant). We conclude that subjects responded to the raw payoff numbers, not to the actual values behind those payoff numbers, and not solely to the relative size of the payoff differences (i.e., the subjects' reactions to payoff numbers appear not to have been scale invariant). There are independent reasons to find this plausible, such as the widely-studied "money illusion" effect (Shafir et al., 1997), in which people focus on nominal rather than real monetary values. CNQRE generally made better predictions than NQRE. This may be due to the fact that normalizing to cents distorts the original payoff values less than normalizing to $[0, 1]$.

These results have some very positive implications. They suggest that pooling data from multiple behavioural experiments (as in this study) is likely to yield meaningful results, even if the precise details of the payoff protocols differ.

## 5.4 How useful is explicit cognitive modelling?

All of the models discussed so far make strong assumptions about the process that agents use to choose the action to play. An alternative predictive approach — perhaps more natural to researchers in AI than in behavioural economics — is to predict action frequencies based on some set of features, without reference to the underlying cognitive processes. Our fourth question was whether such a feature-based, process-agnostic model would predict human behaviour as well as BGT models.

To explore this question, we built a feature-based classifier. Action sets differ between games, which complicates classifier construction. In particular, since games in the training set may have completely different action sets than games in the test set, it does not make sense to use actions as the labels for a classifier based on features of the games.

Instead, we built a classifier based on *logit discrete choice models* (Train, 2009), in which features are associated with *actions* rather than with *games*. In our model, DCM, we define $V(a_i)$ to be a linear combination of nine features

Figure 5.4: Average ratios of DCM and QCH likelihoods to best BGT likelihood, with 95% confidence intervals.

of the action: the minimum and maximum utility of playing the action; the maximum regret of playing an action; a measure of strategic domination; four binary features that indicate whether a level-1 (level-2,3,4) agent would play the action in the level-$k$ model;[9] and a constant value of 1. Agents are assumed to choose the available action that maximizes $U(a_i) = V(a_i) + \epsilon_{a_i}$, where $\epsilon_{a_i}$ is a random variable which is independent and identically distributed across agents. For a given set of actions $A_i$, this implies choice probabilities given by

$$\mathbf{P}(a_i) = \frac{\exp[V(a_i)]}{\sum_{a'_i \in A_i} \exp[V(a'_i)]}.$$

**Model Definition 9** (DCM)**.** Let $\mathcal{F}^{DCM}$ be the set of features defined in Table 5.1. DCM is a discrete choice model where

$$V(a_i) = \sum_{f \in \mathcal{F}^{DCM}} w_f f(a_i).$$

The likelihood of an observation $(a_i, G)$ is thus

$$L^{DCM}(G, a_i \mid \{w_f\}_{f \in \mathcal{F}^{DCM}}) = \frac{\exp[V(a_i)]}{\sum_{a'_i \in A_i} \exp[V(a'_i)]}.$$

---

[9]The reader might wonder why we included the behavioural level-$k$ features. The intention was to allow iterated reasoning to play a part in agents' reasoning without dictating its relative importance. We also evaluated a model that used only the five non-behavioural features, which performed worse in every dataset than DCM.

| Feature | Definition |
|---|---|
| MIN-UTILITY | $\max_{a_{-i} \in A_{-i}} u_i(a_i, a_{-i})$ |
| MAX-UTILITY | $\min_{a_{-i} \in A_{-i}} u_i(a_i, a_{-i})$ |
| MAX-REGRET | $\max_{a_{-i} \in A_{-i}} r(a_i, a_{-i})$, where |
| | $r(a_i, a_{-i}) = \left[ \max_{a_i' \in A_i} u_i(a_i', a_{-i}) \right] - u_i(a_i, a_{-i})$ |
| DOMINATED-SUM | $\sum_{a_i^* \in d(a_i)} \sum_{a_{-i} \in A_{-i}} u(a_i^*, a_{-i}) - u(a_i, a_{-i})$ |
| | where $d(a_i) = \{a_i^* \in A_i \mid a_i^* \text{ dominates } a_i\}$. |
| LEVEL-1? | 1 if a level-1 agent (in level-$k$ models) would play this action, or 0 otherwise. |
| LEVEL-2? | 1 if a level-2 agent (in level-$k$ models) would play this action, or 0 otherwise. |
| LEVEL-3? | 1 if a level-3 agent (in level-$k$ models) would play this action, or 0 otherwise. |
| LEVEL-4? | 1 if a level-4 agent (in level-$k$ models) would play this action, or 0 otherwise. |
| CONST | 1 |

Table 5.1: Features for each action $a_i$ of agent $i$ for discrete choice models.

The results of comparing DCM to GLk are shown in Figure 5.4. Although DCM performed somewhat better than GLk in the CGCB98, SW94, and DS datasets, overall it made considerably worse predictions. The likelihood of the combined ALL4 dataset was almost $10^6$ times smaller according to DCM's predictions than according to GLk's. We thus conclude that there is solid empirical support for the practice of explicitly modelling cognitive processes.

## 5.5 Does heterogeneity matter?

The generalized level-$k$ model incorporates heterogeneity of both steps of iterative reasoning and precision of quantal response. Different agent types may have different quantal choice precisions, and higher-level agents' beliefs about the precisions of other levels may differ from both each other and reality. Our final question was whether a more constrained model would predict equally well.

We constructed a model in which non-random agents were constrained to have identical precisions. Further, the agents were constrained to have correct beliefs about the precisions and the relative proportions of lower-level types. This model can also be viewed as an extension of cognitive hierarchy that adds quantal response; hence we called it *quantal cognitive hierarchy*, or QCH. [10] This is similar to the *Truncated QRE* (TQRE) model of Rogers et al. (2009), in that

---

[10]The reader may wonder why we chose to extend cognitive hierarchy rather than a level-$k$ model. In our opinion, the agents' having a true belief about the precision of lower agent types is more consistent with the assumptions of cognitive hierarchy than with those of level-$k$. We did in fact fit a restricted "quantal level-$k$" model as well; its performance was roughly comparable to that of QCH.

all agents have correct beliefs about lower-level types' proportions and strategies. It differs in that agents are homogeneous in precision, while remaining heterogeneous in steps of reasoning, whereas different types must have different precisions in TQRE.

**Model Definition 10** (QCH). For a level-$m$ agent $i$ (with $m > 0$), let $\beta_{-i,m}^{QCH}$ be a vector of distributions $\beta_{j,m}^{QCH}$ representing $i$'s beliefs about the play of the other agents $j$, with

$$\beta_{j,m}^{QCH}(a_j) = \frac{\sum_{\ell=0}^{m-1} \alpha_\ell \pi_{j,\ell}^{QCH}(a_j)}{\sum_{\ell=0}^{m-1} \alpha_\ell}.$$

Level-$m$ agents choose actions with probability

$$\pi_{i,m}^{QCH}(a_i) = \frac{\exp[\lambda \cdot u_i(a_i, \beta_{-i,m}^{QCH})]}{\sum_{a_i' \in A_i} \exp[\lambda \cdot u_i(a_i', \beta_{-i,m}^{QCH})]}.$$

QCH assumes that all agents are level-4 or lower. The likelihood of an observation $(a_i, G)$ is thus

$$L^{QCH}(G, a_i \mid \alpha_1, \alpha_2, \alpha_3, \alpha_4, \lambda) = \sum_{m=0}^{4} \pi_{i,m}^{QCH}(a_i).$$

Figure 5.4 shows the comparison between the prediction performance of QCH and GLk. QCH actually performed somewhat *better* on the combined ALL4 dataset. Overall its performance was similar to GLk's, never performing worse by more than a factor of 10. This suggests that the added flexibility of GLk in terms of heterogeneous beliefs and precisions did not add substantial predictive power.

# 6 Conclusions

## 6.1 Discussion

Looking at the parameter settings that the optimization procedure chooses can give addition insight into the workings of the various models. In this section we consider two general trends in the parameter settings of the behavioural models that we studied.

One feature that appeared in all the fitted models was a high degree of randomness. On the ALL4 dataset, the percentage of the population that the fitted models predicted were behaving in a purely random fashion ranged from a high of 84% (UNEE) to a low of 51% (QCH). Unsurprisingly, models that predicted a lower number of random agents tended to have a better performance. Table 6.1 compares the proportion of random agents predicted (on average) by models that had the concept of a purely-random type with the performance of those same models on the ALL4 dataset.[11] Note that even the best-performing model (QCH) can do no better than to assume that half of the agents are choosing actions completely at random! This indicates that there is still considerable room for improvement in modelling of game behaviour.

Most models had extremely stable estimates for their parameters, with 95% confidence intervals no longer than 0.15. GLk was a striking exception. The parameters specifying relative proportions of different level types in the population ($\alpha_1$ and $\alpha_2$) were roughly as stable as the other models. However, the precision parameters ($\lambda_1$, $\lambda_2$, and $\mu$) had considerably higher variance, with confidence intervals up to 1.96 wide. This is especially interesting in comparison to QCH, whose precision parameter ($\lambda$) had a confidence interval of width 0.02. The higher stability of QCH's parameter estimates is matched by a higher stability of performance. For example, on the ALL4 dataset, QCH had base-10 log likelihood $-156.18 \pm 0.03$, compared to GLk, which had base-10 log likelihood $-157.72 \pm 0.20$.[12] The confidence interval for GLk's performance was nearly 7 times wider than that for QCH. This indicates that QCH is a much more robust model than GLk.

---

[11]See Table A.2 for a full listing of model parameters
[12]See Table A.1 for a complete listing of model log likelihoods

| Model | Random agents | Log likelihood |
|---|---|---|
| UNEE | $0.84 \pm .000$ | $-170.19 \pm 0.02$ |
| NNEE | $0.77 \pm .000$ | $-165.41 \pm 0.01$ |
| Poisson-CH | $0.75 \pm .002$ | $-165.81 \pm 0.06$ |
| Lk | $0.71 \pm .000$ | $-165.64 \pm 0.03$ |
| Lk4 | $0.65 \pm .004$ | $-165.53 \pm 0.03$ |
| CH4 | $0.63 \pm 0.03$ | $-163.82 \pm 0.13$ |
| GLk | $0.59 \pm 0.02$ | $-157.72 \pm 0.20$ |
| QCH | $0.51 \pm .003$ | $-156.18 \pm 0.03$ |

Table 6.1: Average proportion of random agents and (base 10) model log likelihoods on the ALL4 dataset, for models with a fully-random type.

## 6.2 Conclusions and recommendations

To our knowledge, ours is the first meta-study of BGT data and models with a focus on prediction. We explored the properties of four important models from the BGT literature, along with various modifications of these models, by comparing the average cross-validated log likelihoods of the models on BGT data. The division of data into folds was done at the level of individual observations rather than at the level of games, which, to our knowledge, is novel.

Overall, we found that the GLk model had substantially better prediction performance than the other models from the BGT literature that we considered. We would thus recommend the use of GLk by researchers wanting to predict human behaviour in games, especially if maximal accuracy is the main concern. QCH, a conceptually simpler extension of cognitive hierarchy, performed almost as well or better than GLk on all datasets, and had stabler parameter estimates and performance. We recommend the use of QCH in settings for which it is important to be able to interpret the parameters (e.g., in a Bayesian setting where "reasonable" priors need to be determined), when it is important to be able to vary the number of modeled types, and when stable parameter estimates and performance are desirable.

## 6.3 Future work

One of the models that we constructed, DCM, takes a feature-based, non-behavioural approach. With the particular set of features that we chose, DCM substantially underperformed GLk. However, with the right set of features, a discrete choice model might well have better performance. Devising such features is one possible direction for future work.

Another potential direction for future work is to apply these models to prediction problems in multiagent systems, such as bargaining agents or empirical mechanism design. This is distinct from existing work in economics that applies behavioural models to explain observed anomalies (e.g., Crawford and Iriberri

(2007), who use a level-$k$ model that has been extended to Bayesian games to explain higher-than-equilibrium bidding in auctions).

A remaining open problem is to evaluate models that have been extended to account for learning and non-initial play, including repeated-game and extensive-form game settings.

# Bibliography

Altman, A., Bercovici-Boden, A., and Tennenholtz, M. (2006). Learning in one-shot strategic form games. In Fürnkranz, J., Scheffer, T., and Spiliopoulou, M., editors, *ECML*, volume 4212 of *Lecture Notes in Computer Science*, pages 6–17. Springer.

Basu, K. and Weibull, J. (1991). Strategy subsets closed under rational behavior. *Economics Letters*, 36(2):141–146.

Becker, T., Carter, M., and Naeve, J. (2005). Experts playing the traveler's dilemma. Diskussionspapiere aus dem Institut für Volkswirtschaftslehre der Universität Hohenheim 252/2005, Department of Economics, University of Hohenheim, Germany.

Benisch, M., Davis, G., and Sandholm, T. (2006). Algorithms for rationalizability and CURB sets. In *Proceedings of the 21st National Conference on Artificial Intelligence (AAAI-06)*, pages 598–604. AAAI Press.

Bishop, C. (2006). *Pattern recognition and machine learning*. Springer.

Camerer, C., Ho, T., and Chong, J. (2004). A cognitive hierarchy model of games. *Quarterly Journal of Economics*, 119(3):861–898.

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

Capra, M., Goeree, J., Gomez, R., and Holt, C. (1999). Anomalous behavior in a traveler's dilemma? *American Economic Review*, 89(3):678–690.

Conitzer, V. and Sandholm, T. (2005). A generalized strategy eliminability criterion and computational methods for applying it. In *Proceedings of the 20th National Conference on Artificial Intelligence (AAAI-05)*, pages 483–488. AAAI Press.

Costa-Gomes, M., Crawford, V., and Broseta, B. (1998). Cognition and behavior in normal-form games: an experimental study. Discussion paper 98-22, UCSD.

Costa-Gomes, M., Crawford, V., and Broseta, B. (2001). Cognition and behavior in normal-form games: An experimental study. *Econometrica*, 69(5):1193–1235.

Crawford, V. and Iriberri, N. (2007). Level-k auctions: Can a nonequilibrium model of strategic thinking explain the winner's curse and overbidding in private-value auctions? *Econometrica*, 75(6):1721–1770.

Efron, B. and Tibshirani, R. J. (1993). *An introduction to the bootstrap*. Chapman & Hall.

Goeree, J. K. and Holt, C. A. (2001). Ten little treasures of game theory and ten intuitive contradictions. *American Economic Review*, 91(5):1402–1422.

Halpern, J. Y. and Pass, R. (2009). Iterated regret minimization: a new solution concept. In *Proceedings of the 21st International Joint Conference on Artificial Intelligence (IJCAI-09)*, pages 153–158. Morgan Kaufmann.

Harless, D. W. and Camerer, C. F. (1994). The predictive utility of generalized expected utility theories. *Econometrica*, 62(6):1251–1289.

Harless, D. W. and Camerer, C. F. (1995). An error rate analysis of experimental data testing nash refinements. *European Economic Review*, 39(3):649–660.

Hyafil, N. and Boutilier, C. (2006). Regret-based incremental partial revelation mechanisms. In *Proceedings of the 21st National Conference on Artificial Intelligence (AAAI-06)*, pages 672–678. AAAI Press.

McKelvey, R., McLennan, A., and Turocy, T. (2007). Gambit: Software tools for game theory, version 0.2007. 01.30.

McKelvey, R. and Palfrey, T. (1992). An experimental study of the centipede game. *Econometrica*, 60(4):803–836.

McKelvey, R. and Palfrey, T. (1995). Quantal response equilibria for normal form games. *Games and Economic Behavior*, 10(1):6–38.

Nash, J. (1950). Equilibrium points in n-person games. *Proceedings of the National Academy of Sciences of the United States of America*, 36(1):48–49.

Nelder, J. A. and Mead, R. (1965). A simplex method for function minimization. *Computer Journal*, 7(4):308–313.

Rogers, B. W., Palfrey, T. R., and Camerer, C. F. (2009). Heterogeneous quantal response equilibrium and cognitive hierarchies. *Journal of Economic Theory*, 144(4):1440–1467.

Shafir, E., Diamond, P., and Tversky, A. (1997). Money illusion. *Quarterly Journal of Economics*, 112(2):341–374.

Stahl, D. and Wilson, P. (1994). Experimental evidence on players' models of other players. *Journal of Economic Behavior and Organization*, 25(3):309–327.

Stahl, D. and Wilson, P. (1995). On players' models of other players: Theory and experimental evidence. *Games and Economic Behavior*, 10(1):218–254.

Tennenholtz, M. (2002). Competitive safety analysis: Robust decision-making in multi-agent systems. *Journal of Artificial Intelligence Research*, 17:363–378.

Train, K. (2009). *Discrete Choice Methods with Simulation*. Cambridge University Press.

Von Neumann, J. and Morgenstern, O. (1944). *Theory of Games and Economic Behavior*. Princeton University Press.

Wellman, M., Greenwald, A., and Stone, P. (2007). *Autonomous Bidding Agents: Strategies and Lessons from the Trading Agent Competition*. MIT Press.

Witten, I. H. and Frank, E. (2000). *Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations*. Morgan Kaufmann.

# A  Experimental Results

This appendix contains two tables. The first, Table A.1, presents the raw average log likelihoods that the figures in Chapter 4 and Chapter 5 are based upon. There is no entry for CNQRE in the RPC09 and GH01 datasets, since these datasets' games were already denominated in cents, and therefore the QRE and CNQRE predictions coincide.

The second, Table A.2, gives the average maximum likelihood parameter estimates for each model and dataset combination.

| | ALL4 | SW94 | CGCB98 | GH01 | RPC09 | DS |
|---|---|---|---|---|---|---|
| Random | $-404.24$ | $-43.94$ | $-112.34$ | $-115.06$ | $-132.90$ | $-192.72$ |
| UNEE | $-170.19 \pm 0.02$ | $-18.56 \pm 0.01$ | $-48.80 \pm 0.01$ | $-44.93 \pm 0.01$ | $-55.83 \pm 0.02$ | $-79.00 \pm 0.01$ |
| NNEE | $-165.41 \pm 0.01$ | $-17.50 \pm 0.01$ | $-48.80 \pm 0.01$ | $-40.49 \pm 0.02$ | $-55.29 \pm 0.01$ | $-78.26 \pm 0.01$ |
| Lk | $-165.64 \pm 0.03$ | $-13.23 \pm 0.03$ | $-48.36 \pm 0.02$ | $-45.28 \pm 0.02$ | $-55.65 \pm 0.02$ | $-76.20 \pm 0.02$ |
| Poisson-CH | $-165.81 \pm 0.06$ | $-13.32 \pm 0.12$ | $-48.38 \pm 0.02$ | $-45.28 \pm 0.01$ | $-54.40 \pm .001$ | $-76.93 \pm 0.01$ |
| QRE | $-163.02 \pm 0.03$ | $-14.54 \pm 0.02$ | $-47.94 \pm .000$ | $-43.31 \pm .003$ | $-52.59 \pm 0.03$ | $-73.49 \pm 0.01$ |
| GLk | $-157.72 \pm 0.20$ | $-12.50 \pm 0.03$ | $-47.88 \pm 0.02$ | $-38.02 \pm 0.03$ | $-51.78 \pm 0.16$ | $-72.94 \pm 0.06$ |
| CH4 | $-163.82 \pm 0.13$ | $-13.48 \pm 0.16$ | $-48.38 \pm 0.02$ | $-45.13 \pm 0.02$ | $-51.49 \pm 0.07$ | $-75.87 \pm 0.15$ |
| Lk4 | $-165.53 \pm 0.03$ | $-12.80 \pm 0.03$ | $-48.41 \pm 0.04$ | $-44.60 \pm 0.06$ | $-55.50 \pm 0.04$ | $-75.96 \pm 0.03$ |
| NQRE | $-166.62 \pm 0.02$ | $-14.48 \pm 0.02$ | $-48.09 \pm 0.01$ | $-45.10 \pm 0.11$ | $-54.14 \pm 0.02$ | $-79.65 \pm 0.02$ |
| CNQRE | $-166.22 \pm .000$ | $-14.54 \pm 0.02$ | $-47.95 \pm .000$ | - | - | $-77.80 \pm .000$ |
| DCM | $-163.67 \pm 0.06$ | $-12.33 \pm 0.03$ | $-46.91 \pm 0.05$ | $-42.05 \pm 0.07$ | $-52.33 \pm 0.06$ | $-72.18 \pm 0.08$ |
| QCH | $-156.18 \pm 0.03$ | $-13.22 \pm 0.07$ | $-47.87 \pm 0.02$ | $-37.96 \pm 0.03$ | $-51.63 \pm 0.09$ | $-72.52 \pm 0.02$ |

Table A.1: Average (base 10) log likelihoods, with 95% confidence intervals

|  |  | ALL4 | SW94 | CGCB98 | GH01 | RPC09 | DS |
|---|---|---|---|---|---|---|---|
| UNEE | $\epsilon$ | $0.84 \pm .000$ | $0.66 \pm .000$ | $1.00 \pm .000$ | $0.72 \pm .000$ | $0.61 \pm .000$ | $0.84 \pm .000$ |
| NNEE | $\epsilon$ | $0.77 \pm .000$ | $0.55 \pm .000$ | $1.00 \pm .000$ | $0.59 \pm .000$ | $0.69 \pm .000$ | $0.82 \pm .000$ |
| Lk | $\alpha_1$ | $0.21 \pm .000$ | $0.57 \pm 0.01$ | $0.10 \pm .000$ | $0.19 \pm .000$ | $0.21 \pm .000$ | $0.20 \pm .000$ |
|  | $\alpha_2$ | $0.08 \pm .000$ | $0.40 \pm .002$ | $0.04 \pm .000$ | $0.06 \pm .000$ | $0.39 \pm .000$ | $0.29 \pm .000$ |
|  | $\epsilon_1$ | $0.000 \pm .000$ | $0.05 \pm 0.01$ | $0.000 \pm .000$ | $0.000 \pm .000$ | $0.000 \pm .000$ | $0.000 \pm .000$ |
|  | $\epsilon_2$ | $0.000 \pm .000$ | $0.13 \pm .002$ | $0.000 \pm .000$ | $0.000 \pm .000$ | $0.53 \pm .001$ | $0.36 \pm .001$ |
| Poisson-CH | $\tau$ | $0.29 \pm .002$ | $1.40 \pm 0.02$ | $0.12 \pm .000$ | $0.27 \pm .000$ | $0.54 \pm .000$ | $0.29 \pm .001$ |
| QRE | $\lambda$ | $0.03 \pm .000$ | $0.09 \pm .000$ | $0.02 \pm 0.00$ | $0.03 \pm 0.00$ | $0.10 \pm .001$ | $0.03 \pm .000$ |
| GLk | $\alpha_1$ | $0.22 \pm 0.01$ | $0.71 \pm .000$ | $0.89 \pm 0.04$ | $0.32 \pm .000$ | $0.17 \pm 0.07$ | $0.22 \pm .003$ |
|  | $\alpha_2$ | $0.19 \pm 0.01$ | $0.29 \pm .000$ | $0.02 \pm .001$ | $0.33 \pm .000$ | $0.51 \pm 0.05$ | $0.51 \pm 0.03$ |
|  | $\lambda_1$ | $2.60 \pm 0.54$ | $0.20 \pm .000$ | $0.02 \pm .001$ | $1.26 \pm 0.01$ | $3.76 \pm 0.63$ | $5.49 \pm 0.65$ |
|  | $\lambda_2$ | $2.49 \pm 0.43$ | $5.04 \pm 0.61$ | $6.50 \pm 0.96$ | $2.35 \pm .001$ | $1.25 \pm 0.42$ | $0.22 \pm 0.13$ |
|  | $\mu$ | $0.06 \pm .000$ | $4.29 \pm 0.75$ | $6.49 \pm 0.48$ | $1.45 \pm .000$ | $0.06 \pm 0.01$ | $0.20 \pm 0.02$ |
| CH4 | $\alpha_1$ | $0.07 \pm .003$ | $0.38 \pm 0.02$ | $0.02 \pm 0.01$ | $0.02 \pm .000$ | $0.09 \pm .004$ | $0.07 \pm .003$ |
|  | $\alpha_2$ | $0.08 \pm 0.01$ | $0.19 \pm 0.01$ | $0.03 \pm .005$ | $0.01 \pm .000$ | $0.15 \pm .005$ | $0.04 \pm 0.01$ |
|  | $\alpha_3$ | $0.14 \pm 0.01$ | $0.11 \pm 0.01$ | $0.03 \pm .003$ | $0.10 \pm .001$ | $0.12 \pm .003$ | $0.17 \pm 0.01$ |
|  | $\alpha_4$ | $0.07 \pm .004$ | $0.11 \pm 0.01$ | $0.03 \pm .005$ | $0.13 \pm .001$ | $0.14 \pm .004$ | $0.06 \pm 0.01$ |
| Lk4 | $\alpha_1$ | $0.20 \pm .000$ | $0.51 \pm .001$ | $0.12 \pm .001$ | $0.10 \pm .003$ | $0.19 \pm .000$ | $0.18 \pm .000$ |
|  | $\alpha_2$ | $0.06 \pm .000$ | $0.35 \pm .002$ | $0.02 \pm .001$ | $0.10 \pm .001$ | $0.14 \pm .004$ | $0.10 \pm .001$ |
|  | $\alpha_3$ | $0.002 \pm .001$ | $0.13 \pm .002$ | $0.02 \pm .003$ | $0.30 \pm .003$ | $0.13 \pm .004$ | $0.003 \pm .002$ |
|  | $\alpha_4$ | $0.08 \pm .002$ | $0.02 \pm .002$ | $0.02 \pm .004$ | $0.50 \pm .001$ | $0.31 \pm .002$ | $0.34 \pm .003$ |
|  | $\epsilon_1$ | $0.001 \pm .000$ | $0.001 \pm .001$ | $0.003 \pm .000$ | $0.18 \pm 0.01$ | $0.000 \pm .000$ | $0.001 \pm .000$ |
|  | $\epsilon_2$ | $0.004 \pm .001$ | $0.16 \pm .003$ | $0.20 \pm 0.03$ | $0.001 \pm .001$ | $0.46 \pm 0.01$ | $0.01 \pm .001$ |
|  | $\epsilon_3$ | $0.70 \pm 0.05$ | $0.01 \pm 0.01$ | $0.44 \pm 0.07$ | $0.84 \pm .004$ | $0.95 \pm 0.01$ | $0.51 \pm 0.07$ |
|  | $\epsilon_4$ | $0.31 \pm 0.01$ | $0.44 \pm 0.06$ | $0.41 \pm 0.07$ | $1.00 \pm .000$ | $0.47 \pm .002$ | $0.55 \pm .002$ |

Table A.2: Average (base 10) model parameter estimates, with 95% confidence intervals

(Continued from previous page)

| | | ALL4 | SW94 | CGCB98 | GH01 | RPC09 | DS |
|---|---|---|---|---|---|---|---|
| NQRE | $\lambda$ | $2.74 \pm .001$ | $9.11 \pm 0.01$ | $1.23 \pm .003$ | $6.10 \pm 0.02$ | $2.90 \pm .002$ | $2.57 \pm .002$ |
| CNQRE | $\lambda$ | $0.01 \pm 0.00$ | $0.23 \pm .001$ | $0.01 \pm 0.00$ | - | - | $0.01 \pm 0.00$ |
| DCM | $w_{\text{Const}}$ | $-0.30 \pm 0.19$ | $-0.75 \pm 0.22$ | $0.31 \pm 0.08$ | $-0.71 \pm 0.73$ | $-0.21 \pm 0.06$ | $0.29 \pm 0.27$ |
| | $w_{\text{Min-utility}}$ | $-0.000 \pm .000$ | $-0.002 \pm .000$ | $-0.01 \pm .000$ | $-0.004 \pm .000$ | $0.05 \pm .000$ | $0.03 \pm .001$ |
| | $w_{\text{Max-utility}}$ | $-0.003 \pm .000$ | $0.02 \pm .000$ | $0.04 \pm .000$ | $-0.002 \pm .000$ | $0.01 \pm .000$ | $0.02 \pm .001$ |
| | $w_{\text{Max-regret}}$ | $-0.01 \pm .000$ | $-0.02 \pm .000$ | $0.005 \pm .000$ | $-0.02 \pm .000$ | $0.01 \pm .000$ | $-0.01 \pm .000$ |
| | $w_{\text{Dominated-sum}}$ | $0.000 \pm .000$ | $-0.06 \pm .001$ | $-0.01 \pm .000$ | $0.000 \pm .000$ | $-0.01 \pm .000$ | $0.000 \pm .000$ |
| | $w_{\text{Level-1?}}$ | $0.35 \pm .003$ | $0.45 \pm 0.02$ | $-0.51 \pm .000$ | $-0.87 \pm 0.05$ | $0.33 \pm .001$ | $-0.59 \pm 0.03$ |
| | $w_{\text{Level-2?}}$ | $0.32 \pm .002$ | $1.23 \pm 0.01$ | $0.19 \pm .000$ | $1.71 \pm 0.03$ | $0.31 \pm .001$ | $0.14 \pm 0.01$ |
| | $w_{\text{Level-3?}}$ | $0.03 \pm .002$ | $0.33 \pm 0.01$ | $-0.42 \pm 0.10$ | $0.47 \pm 0.02$ | $-0.11 \pm .000$ | $0.38 \pm 0.01$ |
| | $w_{\text{Level-4?}}$ | $0.20 \pm .002$ | $0.45 \pm 0.01$ | $0.02 \pm 0.10$ | $0.25 \pm 0.02$ | $-0.01 \pm .001$ | $-0.43 \pm 0.01$ |
| QCH | $\alpha_1$ | $0.09 \pm .002$ | $0.29 \pm 0.04$ | $0.84 \pm 0.04$ | $0.13 \pm .000$ | $0.10 \pm .002$ | $0.09 \pm .000$ |
| | $\alpha_2$ | $0.16 \pm .001$ | $0.34 \pm 0.04$ | $0.08 \pm 0.04$ | $0.16 \pm .000$ | $0.16 \pm 0.01$ | $0.27 \pm .000$ |
| | $\alpha_3$ | $0.13 \pm .000$ | $0.14 \pm 0.01$ | $0.03 \pm 0.01$ | $0.22 \pm .000$ | $0.13 \pm 0.01$ | $0.03 \pm .000$ |
| | $\alpha_4$ | $0.12 \pm .000$ | $0.11 \pm 0.01$ | $0.04 \pm 0.01$ | $0.22 \pm .000$ | $0.15 \pm .004$ | $0.12 \pm .000$ |
| | $\lambda$ | $3.39 \pm 0.01$ | $0.40 \pm 0.04$ | $0.02 \pm .000$ | $2.11 \pm .002$ | $8.67 \pm 0.46$ | $4.09 \pm .005$ |

Table A.2: Average (base 10) model parameter estimates, with 95% confidence intervals