

# COUNTERFACTUAL PREDICTION IN COMPLETE INFORMATION GAMES: POINT PREDICTION UNDER PARTIAL IDENTIFICATION\*

SUNG JAE JUN<sup>†</sup> AND JORIS PINKSE<sup>‡</sup>

Center for the Study of Auctions, Procurements and Competition Policy  
Department of Economics  
The Pennsylvania State University

this version: February 2019

We study the problem of counterfactual prediction in discrete decision games with complete information, pure strategies, and Nash equilibria: the presence of multiple equilibria poses unique challenges. We introduce multiple types of counterfactuals to establish sharp identified bounds for their prediction probabilities. We propose and compare various point prediction methods, namely midpoint prediction, an approach using a Dirichlet-based prior, a maximum entropy method, and minmax with an entropy constraint. On balance, we conclude that the maximum-entropy approach is the least of several evils. Our results have implications for counterfactual prediction in other models with partial identification.

**Key words:** complete information games, counterfactual prediction, partial identification, maximum entropy, Dirichlet process, minmax decisions.

**JEL Codes:** C01; C10; C57.

---

\*We thank Victor Aguirregabiria for a great discussion at the 2016 Canadian Econometrics Study Group meeting at London, Ontario. We thank Tim Armstrong, Paul Grieco, Ken Hendricks, Rohit Lamba, Andrés Aradillas-López, Richard Blundell, Andrew Chesher, Joachim Freyberger, Ron Gallant, Bruce Hansen, Marc Henry, Nail Kashaev, Michael Leung, Francesca Molinari, Roger Moon, Ulrich Müller, Hashem Pesaran, Jack Porter, Mar Reguant, Bruno Salcedo, Xiaoxia Shi, Elie Tamer, Guofu Tan, conference participants at the Penn State-Cornell econometrics conference at Ithaca, a conference at the Institute for Economics at Academia Sinica, and seminar participants at Cambridge University, the Centre for Microdata and Practice and the Institute for Fiscal Studies at University College London, Erasmus University, the Institute for New Economic Thinking at the University of Southern California, the London School of Economics, the Research Institute of Economics and Management, Simon Fraser University, Syracuse University, Texas A&M, the Tinbergen Institute, the University of British Columbia, the Wang Yanan Institute for Studies in Economics, the University of Wisconsin at Madison, and Western University for helpful comments and suggestions. Finally, we thank the coeditor, the associate editor, and three anonymous referees for their helpful comments.

<sup>†</sup>619 Kern Bldg, University Park, PA 16802, suj14@psu.edu

<sup>‡</sup>616 Kern Bldg, University Park, PA 16802, joris@psu.edu

# 1. Introduction

**1.1 Problem:** We are interested in prediction problems in which the prediction of interest is only partially identified, particularly ones in which a point prediction is *also* desired. One prominent class of examples of models featuring partial identification consists of discrete decision games with complete information and Nash equilibria in pure strategies, in which *counterfactual prediction probabilities* are desired. There are many applications in which such games arise: examples can be found in Bresnahan and Reiss, 1991b; Kooreman, 1994; Soetevent and Kooreman, 2007; Jia, 2008; Ciliberto and Tamer, 2009; Bajari, Hong, and Ryan, 2010; Grieco, 2014, among others. The models associated with such games are often *incomplete* in the sense that the values of the payoff variables do not always specify a unique outcome. For instance, payoffs can be such that there is room for exactly one player to ‘enter,’ but it could be either one. The problems associated with identification and estimation of payoffs in such games, which are mostly a consequence of the incompleteness of the model, have been studied extensively and it is well-known how to address them. For instance, Tamer (2003) relies on the assumption that covariates have sufficiently large support and Bajari, Hahn, et al. (2011) use finite mixtures to deal with multiple equilibria. See also Kline and Tamer (2012) and Kline (2015). However, the incompleteness of the model poses unique challenges for the problem of counterfactual prediction even if the payoff structure is known in its entirety.

To facilitate the analysis plus the exposition and to maximize the ratio of intuition to technical detail, we focus on the simplest possible (and somewhat hackneyed) case that has all relevant features: a single-shot game with two players, and binary decisions, albeit that we allow for both strategic substitutability and strategic complementarity (Bulow, Geanakoplos, and Klemperer, 1985). The problems discussed here arise a fortiori in more general scenarios, including ones with more than two players, nonbinary decisions (Aradillas-Lopez, 2011), mixed strategies, dynamics (Aguirregabiria and Mira, 2007; Bajari, Benkard, and Levin, 2007; Pakes, Ostrovsky, and Berry, 2007), incomplete information (Seim, 2006; Liu, Vuong, and Xu, 2013; Xu, 2014), and more general solution concepts (Aradillas-López and Tamer, 2008; Kashaev, 2015; Kashaev and Salcedo, 2015; Magnolfi and Roncoroni, 2016). Although our general approach can be used to address many such problems, the results established here only provide intuition for similar models with more general discrete action spaces.

Player payoffs are functions of unobservables  $e$  and observables  $x$ , both of which are known to the players.<sup>1</sup> The observables  $x$  are exogenous in the sense that they are independent of all unobservables in the model. For some combinations of  $e, x$  there exists only a single Nash equilibrium in pure strategies and hence a single value  $y$  of the outcome variables  $y$ : the set of such values of  $e$  given  $x = x$  is denoted by  $S_y(x)$ . For other combinations of  $e, x$  there can exist multiple (in our case two) pure strategy Nash equilibria:  $(1, 0)$  and  $(0, 1)$  for the case of strategic substitutability and  $(0, 0)$ ,  $(1, 1)$  in the case of strategic complementarity. The model is hence ‘incomplete’ (Tamer, 2003) in the sense that the same values of the payoff variables  $e, x$  can lead to different outcomes. For given  $x$ ,  $S_m(x)$  denotes the region of values  $e$  in which payoffs do not produce a unique

---

<sup>1</sup>We use bold typeface to denote random variables.

outcome. We therefore call  $S_m(x)$  the *multiplicity region*. Most of the formal results of this paper are established for the strategic substitutability case, where the extension to strategic complementarity is natural. The following discussion will assume strategic substitutability unless otherwise specified. For the purpose of intuition, it is adequate to think of the game considered in this paper as an entry game, but the scope is broader than that.

We do not assume that we know why a particular outcome arises in the multiplicity region. Identification of the payoff structure does not require knowledge of such reasons (e.g. Tamer, 2003; Kline, 2015), but such knowledge *is* required for the counterfactual prediction probabilities to be identified. If such knowledge were available then it should be used to obtain point identification and our analysis becomes redundant. Thus, we will seek to *characterize* instead of to *specify* players' behavior in the multiplicity region. For an unknown function  $p$ , we *characterize* the behavior of players by a probability

$$p = p(e, v, x), \quad (1)$$

for which

$$\mathbb{P}\{y = (1, 0) \mid e = e, p = p, x = x\} = p, \quad e \in S_m(x), \quad (2)$$

where  $v$  is unobserved and  $p$  belongs to the set  $\mathcal{Q}(e, x) = \{p : \exists 0 < v \leq 1 : p(e, v, x) = p\}$ . Thus,

$$p(e, v, x) = \mathbb{P}\{y = (1, 0) \mid e = e, v = v, x = x\}, \quad e \in S_m(x).$$

The variable  $v \sim U(0, 1)$  represents potential unobserved 'market' heterogeneity (separate from the payoff shifters  $e, x$ ) that can affect the probability of the outcome  $(1, 0)$  occurring if  $e \in S_m(x)$ . We shall assume the function  $p$  to be weakly increasing in  $v$  such that a higher value of  $v$  makes the outcome  $(1, 0)$  no less likely. We can derive results analogous to ours absent this assumption, but bounds obtained without the monotonicity assumption would be wider than the ones that we have and having  $p$  weakly increasing in  $v$  is attractive from an intuition perspective. Alternatively, one can think of  $p(e, v, x)$  as a reduced-form representation of  $p$  based on the conditional quantile function of  $p$  given  $e, x$ , where we allow for the possibility that  $p$  can be 'random' even after controlling for the payoff shifters  $e$  and  $x$ ;<sup>2</sup> hence the characterization  $p = p(e, v, x)$  cannot be tested or refuted. Indeed, there always exists an alternative characterization in which  $v$  is absent. However, by using the characterization that includes  $v$  the class of counterfactuals we can discuss becomes richer: in section 2, we provide several examples in which payoff shifters alone do not determine the equilibrium selection probability.

Since  $p$  is the *probability* that  $(1, 0)$  is realized in the multiplicity region, we need another (uniformly distributed) unobservable  $u$  to complete the description of which outcome is realized:  $e, u, v, x$  are mutually independent, which is not restrictive for reasons discussed in section 2. Thus, if  $e \in S_m(x)$  then  $y = (1, 0)$  arises if and only if  $u \leq p$ . Indeed, we can write  $y = y(e, u, v, x)$ , where  $y$  is an unknown function to be defined formally in section 2. Recall that for  $e \notin S_m(x)$ , the payoff variable values specify a unique outcome and hence

<sup>2</sup>For instance, suppose that the conditional distribution function  $F_{p|e, x}$  of  $p$  given  $e, x$  is strictly increasing. Then we can define  $v = F_{p|e, x}(p|e, x)$ , which is uniformly distributed and independent of  $e, x$  by construction. We can now set  $p(e, v, x) = F_{p|e, x}^{-1}(v|e, x)$ . If  $p$  is discrete given  $e, x$ , then the arguments can be extended by using the usual generalized-inverse adjustments for CADLAG functions.

the values of  $\mathbf{u}$ ,  $\mathbf{v}$  are then irrelevant. For  $e \in S_m(\mathbf{x})$ , however, all four variables  $(e, \mathbf{u}, \mathbf{v}, \mathbf{x})$  are germane.

Recall from (1) that we use the characterization  $\mathbf{p} = \rho(e, \mathbf{v}, \mathbf{x})$  and that in the multiplicity region  $\mathbf{u} \leq \mathbf{p} \Leftrightarrow \mathbf{y} = (1, 0)$ . There are examples in the literature (e.g. Bjorn and Vuong, 1984; Jia, 2008; Grieco, 2014) that are nested in our characterization: section 2 contains a discussion. We know of no other papers in the econometrics or applied microeconomics literature that allow for heterogeneity like  $\mathbf{v}$  in the present context.

Note that we have two random variables,  $\mathbf{u}$  and  $\mathbf{v}$ , that together represent unobserved market heterogeneity not affecting payoffs:  $\mathbf{v}$  affects the *probability*  $\mathbf{p}$  that a particular equilibrium is reached in the multiplicity region and a combination of  $\mathbf{u}$  and  $\mathbf{p}$  produces the actual outcome. The reason for having both  $\mathbf{u}$  and  $\mathbf{v}$ , then, is that we want to allow for counterfactuals in which some, but not necessarily all, of the market heterogeneity is fixed. Further, a model with only  $\mathbf{v}$  is more general than a model with only  $\mathbf{u}$  since  $\rho$  need not be invertible in  $\mathbf{v}$ . Finally, there are examples in the economic theory literature in which  $\mathbf{u}$  and  $\mathbf{v}$  have a natural structural interpretation that is lost once either of the two variables is removed from the model; section 2 contains several such examples. Losing the structural interpretation implies losing the ability to conduct a number of interesting thought experiments.

We now turn our attention to the issue of counterfactual prediction. We denote the counterfactual outcomes of  $\mathbf{x}$ ,  $\mathbf{y}$  by  $\mathbf{x}^*$ ,  $\mathbf{y}^*$ , respectively, and are interested in predicting  $\mathbf{y}^*$ . Since  $\mathbf{y}^*$  is binary, we focus on the conditional probabilities of the counterfactual outcome given the observables, i.e.

$$\mathbb{P}(\mathbf{y}^* = y^* \mid \mathbf{x}^* = x^*, \mathbf{x} = x, \mathbf{y} = y), \quad (3)$$

where  $y, y^* \in \mathcal{S} = \{(0, 0), (0, 1), (1, 0), (1, 1)\}$ . We also condition on  $\mathbf{x}$ ,  $\mathbf{y}$  in (3) since their values contain information about the values of  $e, \mathbf{u}, \mathbf{v}$ .

The objective in this paper is to construct predictions of counterfactual outcomes at the population level if parts of the environment are held fixed: we take the distribution of observables as known. In all counterfactuals,  $\mathbf{x} = x$  is replaced by a new value  $\mathbf{x}^* = x^*$ , where  $x^* = x$  is allowed. The definition of the counterfactual outcome of interest  $\mathbf{y}^*$  will further depend on how the unobservables  $e, \mathbf{u}, \mathbf{v}$  are treated. In the simplest (and least interesting) case,  $e, \mathbf{u}, \mathbf{v}$  are replaced with independent copies  $e^*, \mathbf{u}^*, \mathbf{v}^*$  (i.e.  $e, \mathbf{u}, \mathbf{v}$  are redrawn), in which case  $\mathbf{y}^* = y(e^*, \mathbf{u}^*, \mathbf{v}^*, x^*)$ . Then, the *regression prediction*

$$q(y^* \mid x^*, x, y) = \tilde{\pi}_{y^*}(x^*) = \mathbb{P}(\mathbf{y} = y^* \mid \mathbf{x} = x^*) \quad (4)$$

is the same as (3), which is hence conveniently point-identified. However, if at least one of the unobservables  $e, \mathbf{u}, \mathbf{v}$  is not redrawn (i.e. is the same in the counterfactual), then the regression prediction (4) differs from the counterfactual prediction probability (3).

We consider various scenarios where different combinations of  $e, \mathbf{u}, \mathbf{v}$  are not redrawn. For example, one might be interested in the counterfactual where  $e$  is fixed: fixing only  $e$  is a different thought experiment than fixing both  $e$  and  $\mathbf{v}$  or fixing  $e, \mathbf{u}, \mathbf{v}$ . We study all eight scenarios but in the main text we focus on two representative cases, i.e. fixing  $\mathbf{v}$  and fixing both  $e$  and  $\mathbf{v}$ . We provide results for the leading cases in sections 3.1 and 3.2 and for

the other scenarios in section 3.3 and appendix B. We denote the counterfactual prediction probabilities in (3) by respectively

$$q_v(y^* | x^*, x, y) \quad \text{and} \quad q_{ev}(y^* | x^*, x, y),$$

when  $v$  is fixed and when both  $e$  and  $v$  are fixed: the subscripts denote which variables are fixed. We define similar symbols analogously. Unlike the regression prediction  $q$ , neither  $q_v$  nor  $q_{ev}$  is generally point-identified because they depend on the function  $p$ , which is *not* (point-) identified.

In addition to counterfactual prediction probabilities of individual outcomes, we consider the case in which it is the number of competitors instead of their identity that is of interest, as in e.g. Bresnahan and Reiss (1991a). For instance, a government may wish to minimize the occurrence of ‘food deserts’ and a regulatory agency would typically prefer duopoly to monopoly, but may be less concerned about the identity of the monopolist. We show in section 3 that when fixing  $v$ , the counterfactual probability of having a monopoly collapses to the corresponding regression prediction and is hence point-identified, although  $q_v$  is not point-identified.<sup>3</sup> However if both  $e$  and  $v$  are fixed then even the probability of monopoly is only partially identified.<sup>4</sup>

There are several solutions to our partial identification problem. One conceptually straightforward approach is to treat the counterfactual prediction problem itself as a partial identification problem and use bounds on the counterfactual probabilities, i.e. to work out the (sharp) *identified set* of counterfactual probability values.<sup>5</sup> Inference under partial identification is now a well-studied topic, also: for instance, a robust Bayes approach (e.g. Kitagawa, 2012) is a possibility. We determine the sharp identified sets for both  $q_v(y^* | x^*, x, y)$  and  $q_{ev}(y^* | x^*, x, y)$  in sections 3.1 and 3.2, respectively: the bounds for all other cases are provided in appendix B.

Reguant (2016) considers a related problem; she proposes a general methodology for *computing* valid bounds (via mixed integer programming) to equilibrium counterfactual outcomes, which are functions of equilibrium strategies. Her results, while interesting, are of limited relevance here. We do provide bounds to our counterfactual prediction probabilities which we show (in our proofs) can be characterized as constrained optimization problems, but they are sharp explicit bounds. Further, the constrained optimization problems in our case feature inequality constraints that are different for upper and lower bounds and obtaining bounds is not the main focus of our paper.

While bounds are useful and using point predictions *alone* can be suboptimal (Manski, 2015), point predictions are nevertheless valuable. First, for policy purposes, a ‘best guess’ may be desired since the optimal policy may depend on the value of the predicted quantity.<sup>6</sup> Second, bounds can be too wide to be practically useful. Third, if bound predictions are used as an input to something else, e.g. a function or an optimization problem, then the identified set of the ultimate object can be wider still or can have the consistency of Swiss

---

<sup>3</sup>Point identification of the number of competitors does not extend to the case with more than two choices and/or players.

<sup>4</sup>We thank Ken Hendricks for suggesting this case.

<sup>5</sup>An early work on bounds can be found in Fréchet (1935). We thank an anonymous referee for this reference.

<sup>6</sup>See Jun and Pinkse (2017) for an analogous example in the estimation of auction models.

cheese instead of being convex, as is desirable for inference with set-identified models.<sup>7</sup> Fourth, by their very definition, bounds correspond to extreme cases: it is ex ante unclear how much attention one should pay to the most extreme eventualities. Fifth and specific to our prediction problem, the bounds for some counterfactuals are discontinuous in  $x$  and the jump in the bound can be as large as 50 percentage points, i.e. half the parameter space: it is unlikely to inspire policy makers with confidence if counterfactual predictions are highly unstable. The reasons for this discontinuity will be explored, e.g. in section 7. In sum, it is useful to have point predictions *in addition to* bounds.

One intuitive point prediction choice is the midpoint  $\theta_m$  of the identified set. Indeed, Song (2014) has shown that (an efficient estimator of) the midpoint of an identified set can be justified on decision-theoretic grounds by using a ‘local asymptotic minimax regret’ criterion. But Song’s environment is substantially different from ours and he is concerned with estimation, not prediction. Here, because the bounds can be discontinuous in  $x$ , so can the midpoints. Indeed, midpoint predictions can have discontinuities in  $x$  that are as large as 25 percentage points. Further, midpoint predictions are by definition an average of two extremes. Third, as we discuss in section 3.3 that  $q_e, q_{ev}, q_{euv}$  are ranked<sup>8</sup> but their bounds coincide. In other words, even though it is known that e.g. (omitting arguments)  $q_e \leq q_{ev} \leq q_{euv}$ , the bounds and hence the midpoint predictions are the same. The maximum entropy predictions discussed below do not have this problem. Finally, there is an additional ‘inconsistency’ problem: as we show in section 7.1, there does not generally exist a single function  $p$  that can generate midpoint predictions for different counterfactuals.

We develop several alternatives to simply using bounds, i.e. to taking the midpoint. Our preferred method is to use maximum entropy from information theory. We further develop two additional decision-theoretic methods. Descriptions of and tradeoffs between these three approaches are discussed in section 1.2.

**1.2 Alternatives to bounds:** Decision-theoretic alternatives to using bounds entail defining a loss function  $\ell$  which measures the distance between the infeasible prediction (which depends on  $p$ ) and the ‘decision’ (which does *not* depend on  $p$ ). We focus on decisions *at the population level* that use the complete knowledge of the model, not a particular set of observations. For example, if  $\theta_0$  is the prediction of interest then complete knowledge of the model only provides the bounds  $\mathcal{B}_L \leq \mathcal{B}_U$  for  $\theta_0$  and these bounds depend only on population objects and they are hence point-identified. Then, the (unconstrained) minmax decision with a quadratic loss is formulated by  $\min_{\theta \in [\mathcal{B}_L, \mathcal{B}_U]} \max_{\theta_0 \in [\mathcal{B}_L, \mathcal{B}_U]} (\theta - \theta_0)^2$ , which yields the solution  $\theta = (\mathcal{B}_L + \mathcal{B}_U)/2$ , i.e. the midpoint decision. There is no randomness here and therefore loss and risk are the same thing: so  $\max_{\theta_0 \in [\mathcal{B}_L, \mathcal{B}_U]} (\theta - \theta_0)^2$  is the maximum risk.

As mentioned earlier, the unconstrained minmax approach leads trivially to the midpoint prediction, which has downsides that have already been discussed. We consider two alternatives: adding a constraint based on entropy to the minmax formulation and

<sup>7</sup>The nonconvexity problem arises e.g. in Haile and Tamer (2003), which they address by imposing additional technical assumptions on the shape of a ‘pseudoprofit function.’

<sup>8</sup>The order can be ascending or descending, depending on the values of  $y, y^*$ .

minimizing average risk that uses a ‘prior’ on  $p$ .

The first alternative, which is described in section 5, assumes a symmetric convex loss function, e.g. a quadratic, and minimizes maximum loss/risk, where the maximum is now taken over all functions  $p$  whose entropy is not below some specified bound  $\mathcal{E}^*$ . This is a hybrid method in that it uses both information theory and decision theory. Indeed, if  $\mathcal{E}^*$  is chosen to be the maximum attainable value then one obtains the maximum entropy solution. At the other extreme, if one chooses  $\mathcal{E}^* = -\infty$  then one obtains the midpoint solution. A different way of looking at this is to note that the entropy restriction reduces the size of the identified set and that the hybrid method then selects the midpoint of this smaller identified set, which also happens to be an interval.

The second alternative is investigated in section 6. There, using Dirichlet processes we specify a new class of probability distributions over  $p$ -functions that is consistent with the distribution of observables.

However, whether one uses the approach of section 6 or something else, the decision-theoretic approaches do not solve the fundamental problem at hand. For instance, there exists no *natural* probability measure on the parameter space, i.e. the class of functions  $p$ : depending on what measure one assigns to the parameter space any prediction in the identified set can be generated (Aumann, 1961). In other words, the minimum-average-risk approach relocates the problem of choosing a prediction from the identified set to the choice of a measure on the parameter space that is consistent with the distribution of observables.<sup>9</sup> Further, the Dirichlet approach only uses discrete probability distributions, whereas  $p$  can — or indeed is more likely to be — continuously distributed. Nevertheless, for the Dirichlet-based probability measures considered in section 6, the decision-theoretic approach generates predictions that are similar to our preferred method, *maximum entropy*, which is discussed in section 4 in detail.

With the maximum entropy method of the information theory literature (Jaynes, 1957a; Jaynes, 1957b; Golan, Judge, and Miller, 1996), one selects the probability distribution that best represents the current state of knowledge as measured by the entropy. The maximum entropy solution complements the information contained in the data with a criterion that is consistent with the (at least) seven hundred year old principle of *Occam’s razor*: if there are multiple explanations for the same phenomenon then one should choose the simplest one. In our context, this translates into choosing the function  $p$  such that out of all candidates that are consistent with the distribution of observables  $(x, y)$  and the model assumptions, the random variable  $p$  has the distribution that is closest to a uniform (given  $e, x$ ).<sup>10</sup> There are numerous philosophical justifications for the application of Occam’s razor; we refer the reader to Baker (2013).

There are justifications for the use of maximum entropy other than Occam’s razor. Indeed, Golan (2017) lists sets of axioms for the problem of choosing a probability distribution absent identification. Theorem 2.1 in Golan (2017) establishes that the maximum entropy solution is the only one that satisfies five stated axioms.<sup>11</sup>

Finally, although information theory and decision theory both have their respective

---

<sup>9</sup>In the other possibility there is an issue of how to choose  $\mathcal{E}^*$ . Setting  $\mathcal{E}^*$  to the two extreme possibilities leads to the maximum entropy solution or the midpoint prediction, which we have already discussed above.

<sup>10</sup>With maximum entropy, the closest choice is unique up to trivial deviations.

<sup>11</sup>Golan (2017) attribute their theorem 2.1 to Shore and Johnson (1980).

merits, one does not need to make a choice of one above another. Indeed, as Topsøe (1979), Harremoës and Topsøe (2001), and Grünwald and Dawid (2004) show, maximum entropy itself has a decision–theoretic interpretation as well. For instance, Grünwald and Dawid (2004) show that maximizing entropy and minimizing worst–case expected loss are each other’s dual in a statistical game in which nature reveals values from an unknown distribution, a decision maker specifies a distribution, and loss is measured by a log–loss function.<sup>12</sup>

Maximum entropy is sometimes confused with the Bayesian approach. Although one of Jaynes’s intentions behind maximum entropy was to provide Bayesians with a sensible prior, the method itself can equally be used in a classical context. Further, since maximum entropy entails an optimization problem whose constraints correspond to the information available in the data, the maximum entropy solution coincides with the standard classical solution in the case of point identification: the constraints then provide a unique solution. In other cases, maximum entropy provides the minimal amount of additional information needed to provide uniqueness. In other words, maximum entropy can be thought of as providing ‘second class’ information.<sup>13</sup>

Here, the function  $p$  selected by maximum entropy is flat in  $e$  in the multiplicity region, which is consistent with equilibrium–selection mechanisms used in the literature (e.g. Bjorn and Vuong, 1984; Jia, 2008; Bajari, Hong, and Ryan, 2010). Note, however, that in our case this is an *outcome* instead of an *assumption*.

One can cook up measures other than the one provided by maximum entropy, so using maximum entropy is not altogether free of arbitrariness. For instance, using a distance criterion other than the one used in maximum entropy (i.e. Kullback–Leibler divergence) yields a different solution. Using maximum entropy on a monotonic transformation of  $p$  will generally also lead to a different solution. Indeed, the exercise of making a single guess about the value of a parameter that is only partially identified is inherently arbitrary. But the use of maximum entropy has been justified extensively in the information theory literature: recall the five axioms listed in Golan and see e.g. Cover and Thomas (2012). Second, since the random variable at the center of the problem is the selection probability  $p$ , it is more natural to deal with the distribution of  $p$  than that of e.g.  $p^2$ . Further, with maximum entropy it is straightforward to introduce additional information (restrictions) to the problem whereas with the decision–theory based Dirichlet–like method doing so is both complicated and cumbersome. As noted above, for the choices made in section 6, the Dirichlet and maximum entropy approaches appear to generate similar predictions while the maximum entropy approach is considerably easier to implement. Finally, as we show in section 8, it is straightforward to address inferential issues for maximum–entropy predictions. Compared to Bayesian alternatives such as pragmatic Bayesian approxima-

<sup>12</sup>If nature draws a value  $z$  of  $\mathbf{z}$  from a density  $f_0$  and a statistician chooses  $f$  to specify the distribution, then the log–loss is defined by  $-\log f(z)$ : i.e. the smaller value  $f$  takes at  $z$ , the larger the loss is. Then, the expected log–loss is given by  $-\mathbb{E} \log f(z)$ . Grünwald and Dawid (2004) note that maximizing entropy with respect to  $f_0$  can be written as  $\max_{f_0} \min_f -\mathbb{E} \log f(z)$  and they discuss switching the order of max and min to reformulate it as a minmax decision problem.

<sup>13</sup>There are other contexts in which one wishes to select a single function from a set, e.g. theorem A.1 in Chen and Pouzo (2012), but the context plus the considerations and implications of choosing one function over another are entirely different.



tions described in e.g. Murray and Snelson (2006), our approach turns out to be a simple application of the classical Delta method once the payoff parameters are estimated.

**1.3 Contribution:** In sum, we see the main contributions of this paper as follows. To our knowledge, we are the first to study the problem of counterfactual point–prediction in games of complete information featuring a partially identified infinite–dimensional parameter. We derive formulas for the infeasible counterfactual prediction probabilities  $q_v(y^* | x^*, x, y)$  and  $q_{ev}(y^* | x^*, x, y)$  as functions of the unknown, partially identified, function  $\rho$  and construct corresponding identified sets in the form of sharp bounds; see section 3. We further derive bounds for the case in which the object of interest is the number instead of the identity of ‘entrants.’ We propose, compare, and contrast various point prediction methods. We develop a new decision theory–based point prediction method in the spirit of the Dirichlet–process literature in section 6 and a new point prediction method based on the maximum entropy concept from the information theory literature in section 4.<sup>14</sup> The hybrid method described in section 1.2 is developed in section 5. We demonstrate the virtues of the maximum entropy approach (and to a lesser extent the Dirichlet approach) in a number of examples in section 7. Section 8 demonstrates how inference can be conducted on our counterfactual prediction probabilities. Finally, analogs of our maximum entropy approach can be used to do counterfactual analysis in other models with partial identification, albeit that in our case the partially identified parameter is a function, which is more complicated than if it were finite–dimensional.

## 2. Setup

**2.1 Binary decision games:** We consider standard two player binary decision games with complete information and pure strategies. In the normal form in figure 1, the functions  $\tau_1, \tau_2$  and payoff variables  $e_1, e_2, \mathbf{x}$  are known to both players: their respective choices are denoted  $y_1, y_2$  and the function  $\tau_i$  depends on  $x$  and  $y_{3-i}$  for  $i = 1, 2$ . We assume that  $\mathbf{x}, y_1, y_2$  are observable to the econometrician, but  $e_1, e_2$  are not. Further,  $\mathbf{e} = (e_1, e_2)$  is assumed to be independent of  $\mathbf{x}$ , which can be vector–valued. We assume that both the  $\tau_i$ –functions and the distribution of  $\mathbf{e}$  are known or, failing that, are identified, e.g. via a result like that in Tamer (2003): we do this because identification and estimation of payoffs have already been studied extensively.<sup>15</sup>

For given values of  $\mathbf{e}$  and  $x$ , the normal form in figure 1 allows for either strategic complementarity ( $\tau_i(x, 0) < \tau_i(x, 1)$  for  $i = 1, 2$ ; you choosing option 1 increases the payoff for me to choose option 1) or strategic substitutability ( $\tau_i(x, 0) > \tau_i(x, 1)$  for  $i = 1, 2$ ; you choosing option 1 reduces the payoff for me to choose option 1). In many cases there is

---

<sup>14</sup>Others have used the notion of entropy in the context of partial identification, albeit with a different purpose. Indeed, Schennach (2014) uses entropy as a way of *constructing* the identified set in an estimation problem, whereas we use it to *select* a point from the identified set in a prediction problem. However, her problem and method are much closer substitutes to Galichon and Henry (2011) than to ours.

<sup>15</sup>Our counterfactual analysis does not allow the payoff structure or the distribution of payoff variables to change. However, such changes could in principle be accommodated as long as the new payoff structure and distribution of payoff variables are given or can be identified using other methods.

|   |                         |  |   |
|---|-------------------------|--|---|
|   | 2                       | 0  | 1 |
| 1 | 0, 0                    | $0, \tau_2(x, 0) + e_2$                  |   |
| 0 | $\tau_1(x, 0) + e_1, 0$ | $\tau_1(x, 1) + e_1, \tau_2(x, 1) + e_2$ |   |

Figure 1: Normal form

only a single Nash equilibrium in pure strategies, but there are also cases in which there are two pure Nash equilibria.<sup>16</sup>

Indeed, consider a given value  $x$  (of  $\mathbf{x}$ ) and let  $\bar{\tau}_i(x) = \max\{\tau_i(x, 0), \tau_i(x, 1)\}$  and  $\underline{\tau}_i(x) = \min\{\tau_i(x, 0), \tau_i(x, 1)\}$ . Then there are five distinct regions: four regions in which there is a unique Nash equilibrium in pure strategies ( $S_{00}(x), S_{01}(x), S_{10}(x), S_{11}(x)$ ) and one *multiplicity region*,  $S_m(x)$ ; see figure 2 for examples. Thus,  $S_{10}(x)$  is the set of  $e = (e_1, e_2)$  values for which (1, 0) is the unique Nash equilibrium if  $\mathbf{x} = x$ . There are two pure strategy Nash equilibria whenever  $e = (e_1, e_2)$  falls in the multiplicity region  $S_m(x) = [-\bar{\tau}_1(x), -\underline{\tau}_1(x)] \times [-\bar{\tau}_2(x), -\underline{\tau}_2(x)]$ .

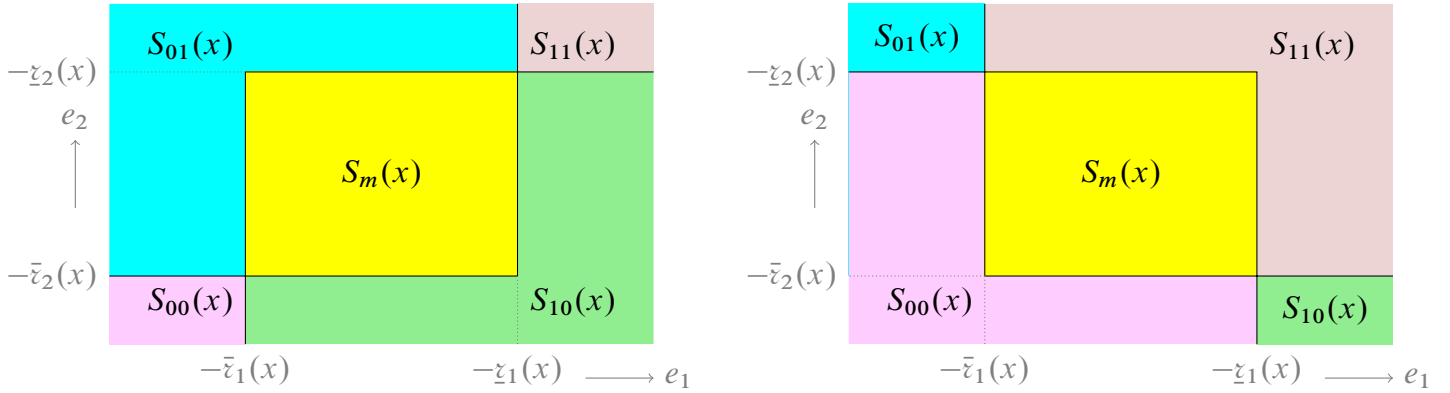


Figure 2: Regions of  $e$ -values corresponding to either a specific equilibrium or multiple ones for the cases of strategic substitutability (i.e.  $\bar{\tau}_i(x) = \tau_i(x, 0)$  for  $i = 1, 2$ , on the left) and strategic complementarity (i.e.  $\bar{\tau}_i(x) = \tau_i(x, 1)$  for  $i = 1, 2$ , on the right).

If  $e$  belongs to  $S_m(x)$  and  $\tau_i(x, 0) < \tau_i(x, 1)$  for  $i = 1, 2$  then one obtains a coordination game with (0, 0) and (1, 1) as the two Nash equilibria in pure strategies. If instead  $e$  belongs to  $S_m(x)$  and  $\tau_i(x, 0) > \tau_i(x, 1)$  for  $i = 1, 2$  then we have an entry game with (0, 1) and (1, 0) as the two Nash equilibria in pure strategies. We discuss examples of coordination games with two (pure strategy Nash) equilibria that are consistent with the normal form in figure 1 in section 2.3.1 and entry games with two equilibria in section 2.3.2. These examples illustrate the need for the additional random variable  $v$  used in the sketch of the counterfactual prediction problem in section 2.2. A discussion of the counterfactuals themselves is postponed until section 3.

<sup>16</sup>Since  $e$  is continuously distributed, the probability of more than two pure strategy Nash equilibria is zero. If  $e$  were discrete then it is possible that the probability that  $e$  were e.g. at the intersection of  $S_m, S_{11}, S_{01}, S_{10}$  would be nonzero.

**2.2 Prediction:** One can relate the outcome  $(y_1, y_2)$  to payoff variables  $(e, x)$  by

$$\begin{cases} y_1 = \mathbb{1}\{\tau_1(x, y_2) + e_1 \geq 0\}, \\ y_2 = \mathbb{1}\{\tau_2(x, y_1) + e_2 \geq 0\}, \end{cases} \quad (5)$$

where  $x$  is a vector of exogenous covariates,  $e = (e_1, e_2)$  are errors that are independent of  $x$ . The model in (5) is *incomplete* in that  $x$  and  $e$  do not necessarily determine a unique outcome  $y = (y_1, y_2)$  due to the possible presence of multiple Nash equilibria, as noted above (see e.g. Bresnahan and Reiss, 1991a; Tamer, 2003). Incompleteness, here, means that there can be multiple pairs  $(y_1, y_2)$  that satisfy (5).

With the exception of section 2.3.1, we assume that there is strategic substitutability, i.e.  $\tau_i(x, 1) \leq \tau_i(x, 0)$  for  $i = 1, 2$  and all  $x$ . We do this merely to avoid duplication of notation and arguments.

The properties of the model in (5) have been studied extensively and conditions under which  $\tau_1, \tau_2$  and the distribution of  $(e, x)$  are identified are well-understood (e.g. Tamer, 2003). From hereon, we therefore take the functions  $\tau_1, \tau_2$  and the distribution of  $(e, x)$  as given.

Since we are taking the functions  $\tau_i$  and the distribution of  $(e, x)$  as given, the  $S$ -regions are given, also, and so are their probabilities of occurring, i.e.

$$\pi_y(x) = \mathbb{P}\{e \in S_y(x)\}, \quad \pi_m(x) = \mathbb{P}\{e \in S_m(x)\}.$$

We further denote

$$\tilde{\pi}_y(x) = \mathbb{P}(y = y \mid x = x).$$

As mentioned in the introduction, we are interested in predicting the counterfactual outcome  $y^*$  if the same game is played again under various circumstances. To do this, we have to provide a (partial) description of the equilibrium determination process in the multiplicity region.

Recall from (2) that for some random variable  $p = p(e, v, x)$  and all  $x$ , if  $e \in S_m(x)$  then

$$\forall p \in \mathcal{Q}(e, x) : \mathbb{P}\{y = (1, 0) \mid e = e, p = p, x = x\} = p, \quad (6)$$

where  $v$  has a standard uniform distribution and is independent of  $e, x$ . The function  $p$  is unknown and (weakly) monotonic in  $v$  for every  $e, x$ .

So from the perspective of the researcher, there is a ‘probability’  $p$  that a particular equilibrium arises in the multiplicity region, where  $p$  can depend on (both observable and unobservable) market heterogeneity. This looks similar to the equilibrium selection mechanism used in Grieco (2014), where  $p(e, v, x)$  depends only on the payoff shifters  $e, x$ , and also nests the structures in Bjorn and Vuong (1984), where  $p(e, v, x) = \mu_m(x)$ , and Jia (2008), where  $p(e, v, x) = 1$ . Our specification is a *characterization* instead of an *assumption* and it is consistent with *any* equilibrium refinement.<sup>17</sup> The same can be said

<sup>17</sup>A. Kalai and E. Kalai (2012) (and references therein) provide refinements for two by two games, albeit that in A. Kalai and E. Kalai (2012) utility is transferable, which means that the outcome would not have to correspond to a Nash equilibrium in pure strategies in the corresponding game without transfers. Other references on equilibrium selection include Harsanyi and Selten (1988) and Kandori, Mailath, and Rob (1993), which we will discuss in more detail in section 2.3.

for Grieco (2014) — and hence also for Bjorn and Vuong (1984) — but, as noted in the introduction, we want to allow for the possibility that  $v$  is fixed and  $u$  changes, i.e. that the unobserved payoff–irrelevant heterogeneity that affects the *equilibrium–selection probability* is fixed but that the equilibrium selection probability itself can be different from zero and one. Indeed, in other contexts it is common to have vector–valued unobservables (e.g. Athey and Imbens, 2007; Briesch, Chintagunta, and Matzkin, 2012): the problem at hand can often be described by a scalar unobservable, but doing so changes the nature of the problem as is argued in Kasy (2011).

Furthermore, as will become apparent in sections 2.3.1 and 2.3.2, there are instances in which  $u, v$  do have a structural interpretation that is lost once either of them is dropped and hence having both  $u$  and  $v$  allows for more interesting counterfactuals than would a single such unobservable. The flipside of this argument is that a characterization in terms of two unobservables  $u, v$  does limit the nature of the counterfactual experiments discussed in section 3 if the underlying structural model of players’ behavior in the multiplicity region features multiple unobservables that feed into  $v$  and that an empiricist is interested in the effects of holding fixed one or more of the structural model unobservables instead of  $v$ : our paper does not speak to that possibility. However, the presence of two payoff–irrelevant unobservables ( $u, v$ ) instead of one, as in Grieco (2014), increases the likelihood that any explicit model of behavior in the multiplicity region can be mapped to our specification.

Equation (6) is silent about the behavior of  $p$  when  $e \notin S_m(x)$  since the value of  $p$  is then immaterial for the selection of  $y$ . Even if  $e \in S_m(x)$ , however, (6) does not restrict the dependence of  $e$  and  $p$ .

In view of (6) we can represent the outcome  $y$  by

$$y = y(e, u, v, x) = \begin{cases} y, & e \in S_y(x), \\ (1, 0), & e \in S_m(x) \text{ and } u \leq p(e, v, x), \\ (0, 1), & e \in S_m(x) \text{ and } u > p(e, v, x), \end{cases} \quad (7)$$

where  $u$  is uniform and independent of  $e, v, x$ .

In view of (7), (6) *completes* the model in that it enables us to determine the outcome as a function of observables and unobservables. However, since  $p$  is unknown — and indeed not identified — the structure that (6) imposes is by itself of limited help for the prediction problem studied in this paper.

**2.3 Equilibrium selection:** As promised, we now discuss examples to illustrate the importance of allowing for an additional unobservable  $v$ ; the fact that  $p$  can depend on the payoff shifters  $e$  and  $x$  is intuitive. If  $e$  belongs to  $S_m(x)$  then there are two pure strategy Nash equilibria. Recall that we *characterize* equilibrium selection by a random variable  $p = p(e, v, x)$ , where  $p$  is the *probability* that a given equilibrium, e.g.  $(1, 0)$ , is selected. The examples in this subsection are divided into two parts, depending on the nature of whether there is strategic complementarity or strategic substitutability.

**2.3.1 Coordination games:** As noted, if  $v_i(x, 0) < v_i(x, 1)$  for  $i = 1, 2$  then one obtains a coordination game in  $S_m(x)$  in which  $(0, 0)$  and  $(1, 1)$  are (the only) pure strategy Nash

equilibria. Since there are two pure strategy Nash equilibria, the payoffs provide insufficient information to produce a unique outcome. We now discuss equilibrium selection for this case.

As an example, consider the payoff matrix in figure 3.

|   |      |      |   |
|---|------|------|---|
| 1 | 2    | 0    | 1 |
| 0 | 0,0  | 0,-2 |   |
| 1 | -2,0 | 1,1  |   |

Figure 3: Simple coordination game

Harsanyi and Selten (1988) refer to  $(1, 1)$  as the *payoff-dominant equilibrium* and to  $(0, 0)$  as the *risk-dominant equilibrium*:  $(0, 0)$  is strictly risk-dominant in this example. Payoff-dominance is equivalent to the traditional definition of efficiency and, as noted in van Damme (1995), is a notion of collective rationality.

In contrast, risk-dominance is about individual rationality. In the example above,  $(0, 0)$  is risk-dominant because each player will play 0 if she believes that her opponent is equally likely to play either 0 or 1. More generally, i.e. without symmetry, let  $\Lambda_i$  denote the probability with which player  $3 - i$  plays 0 according to player  $i$ 's beliefs and let  $\Lambda_i^*$  be the value of  $\Lambda_i$  for which player  $i$  is indifferent between playing 0 and 1. Then  $(0, 0)$  is risk-dominant if it is a Nash equilibrium and  $\Lambda_1^* + \Lambda_2^* \leq 1$ . In this example,  $\Lambda_1^* + \Lambda_2^* = 2/3 < 1$ . Note that the same equilibrium can be both risk-dominant and payoff-dominant and need not be either.

There is an extensive literature on equilibrium selection in this context in game theory (e.g. Foster and Young, 1990; Kandori, Mailath, and Rob, 1993). A common question in this literature is whether or not a payoff-dominant equilibrium is more appealing than a risk-dominant one. For instance, Kandori, Mailath, and Rob (1993) propose an evolutionary model and they discuss the probability of a risk-dominant equilibrium being played. There, the same  $2 \times 2$  stage game is repeated indefinitely, players are chosen at random from a finite set, and a player from the set is replaced with a new uninformed player at random times according to a probability Kandori, Mailath, and Rob call the 'mutation rate.' They show that the probability that the risk-dominant equilibrium is played in the stage game is in the long run independent of initial conditions (the initial strategies) and can be any number between 0 and 1, with the number depending on both the payoffs and the mutation rate.<sup>18</sup> Indeed, the probability that the payoff-dominant equilibrium is played is increasing in the mutation rate. Thus, the payoff values in figure 3 are insufficient to determine which equilibrium is more likely to occur at a given point in time.<sup>19</sup>

In our context, then,  $v$  is the mutation rate and  $p(e, v, x)$  the probability that  $(1, 1)$  is selected. In order to determine which equilibrium is played in the stage game we need an

<sup>18</sup>We focus on coordination games here, but Kandori, Mailath, and Rob's analysis applies more generally. In particular, it also applies to entry games, as will become apparent in section 2.3.2.

<sup>19</sup>Note that these are unconditional probabilities. The probability that a particular equilibrium occurs at time  $t + 1$  if the same equilibrium occurred at time  $t$  is much higher than the unconditional probability if the mutation rate is low.

additional standard uniform random variable  $\mathbf{u}$  such that in the multiplicity region,

$$\mathbf{u} \leq p(\mathbf{e}, \mathbf{v}, \mathbf{x}) \Leftrightarrow \mathbf{y} = (1, 1).$$

So the mutation rate  $\mathbf{v}$  and the payoff variables  $\mathbf{e}, \mathbf{x}$  determine the probability that  $(1, 1)$  is selected in the steady state if  $\mathbf{e} \in S_m(\mathbf{x})$  and  $\mathbf{u}$  determines which equilibrium is selected in the stage game at the time of ‘observation.’ We may call  $\mathbf{u}$  a ‘sunspot’ in that it is an extrinsic random component that affects the actual outcome of the game (e.g. Cass and Shell, 1983; Aguirregabiria and Mira, 2018).

The example described above is just one of several conceivable examples generating two errors  $\mathbf{u}, \mathbf{v}$  that are unrelated to the payoffs. In any such example, the variables  $\mathbf{u}, \mathbf{v}$  can vary across markets.

The fact that  $\mathbf{u}$  and  $\mathbf{v}$  have different interpretations does not mean that both are required to describe equilibrium selection. Indeed, one could feed  $\mathbf{u}$  and  $\mathbf{v}$  to a single unobserved variable and represent the selection probability as if the payoff variables were its sole determinants. Put differently, one could define the selection probability without conditioning on  $\mathbf{v}$ . In fact, this is the route chosen by other authors, including Grieco (2014) and Aguirregabiria and Mira (2018).

However, if one only had one of these random variables, say  $\mathbf{u}$ , then one could not attribute meanings to each of  $\mathbf{u}, \mathbf{v}$  as in the Kandori, Mailath, and Rob (1993) example above. This is not necessarily a limitation when the purpose of the analysis is something other than the identification of counterfactuals. However, the focus of our paper *is* on the counterfactuals and the nature of the counterfactual experiments that we consider is therefore important. Indeed, if we only had a single market heterogeneity unobservable then (in the above example) we could not differentiate between holding constant the mutation rate  $\mathbf{v}$  and holding constant the state  $\mathbf{u}$  that the market is in: the mutation rate is a long run object whereas  $\mathbf{u}$  is specific to the stage game.

More generally, without  $\mathbf{v}$  there are only two counterfactuals that can be considered in  $S_m(\mathbf{x})$  if the payoff variables are the same in the counterfactual as in the ‘data:’ (1) the counterfactual outcome is the same as the ‘observed’ outcome; (2) the counterfactual outcome is, conditional on the values of the payoff variables, independent of the ‘observed’ outcome. The first case arises if  $\mathbf{u}^*$  (‘sunspot’ in the counterfactual) is the same as  $\mathbf{u}$  (‘sunspot’ in the ‘data’) and the second case if  $\mathbf{u}^*$  is a new draw. In other words, by omitting  $\mathbf{v}$  one rules out the possibility that the payoff-irrelevant unobservables contain *some but not all* information helpful for the counterfactual outcome over and above that which is contained in the payoff variables  $\mathbf{e}, \mathbf{x}$ . An explicit example of such a situation can be found at the end of section 2.3.2. While it is true that this limitation can be relaxed by allowing for imperfect dependence between  $\mathbf{u}^*$  and  $\mathbf{u}$ , doing so would introduce a new random variable, also.<sup>20</sup>

**2.3.2 Entry games:** If  $\tau_i(0, \mathbf{x}) > \tau_i(1, \mathbf{x})$  for  $i = 1, 2$  in figure 1 then we have strategic substitutability, an ‘entry’ game. In the multiplicity region, the two pure strategy Nash equilibria are now  $(1, 0)$  and  $(0, 1)$ . In contrast to the case discussed in section 2.3.1, in

---

<sup>20</sup>For instance, one could write  $\mathbf{u}^* = \mathbf{u} + \mathbf{v}$ .

the entry game version of figure 1 each player prefers a different equilibrium and the two equilibria cannot be ranked.<sup>21</sup> Since there are again two pure strategy Nash equilibria, equilibrium selection is still an issue.

Consider the example depicted in figure 4. This example differs from figure 3 because

|   |      |        |   |
|---|------|--------|---|
| 1 | 2    | 0      | 1 |
| 0 | 0, 0 | 0, 1   |   |
| 1 | 1, 0 | -1, -1 |   |

Figure 4: An entry game

the two Nash equilibria, i.e.  $(1, 0)$  and  $(0, 1)$ , cannot be ranked: each equilibrium is neither payoff-dominant nor strictly risk-dominant. The same interpretations based on types of market heterogeneity apply here:  $v$  is long-run or player-originated heterogeneity such as the mutation rate or a commitment signal from a player and  $u$  corresponds to short-run stage idiosyncrasies. Indeed, in the context of Kandori, Mailath, and Rob (1993)  $u, v$  play the same role as they did in section 2.3.1.

|   |      |        |   |
|---|------|--------|---|
| 1 | 2    | 0      | 1 |
| 0 | 0, 0 | 0, 2   |   |
| 1 | 1, 0 | -1, -1 |   |

Figure 5: Another entry game

A second example concerns collusion. Suppose that in the game depicted in figure 5, nonparticipants' payoffs are increasing in the number of 'entrants,' and that it is illegal for players to coordinate. Thus, total surplus can be less if players coordinate. Suppose further that side payments are infeasible due to the risk of discovery. One can think of participants in a procurement withholding a bid to benefit another ring member. Bid rigging on milk contracts for schools is a well-known example; see Porter and Zona (1999).

The players could assign the right to bid randomly, resulting in outcome  $(1, 0)$  with probability  $p$ . If  $v$  denotes the bargaining weight of player 1 then  $p$  depends on both the conflict payoffs and on  $v$ .<sup>22</sup>

For example, in figure 5 there are two pure strategy Nash equilibria plus a mixed strategy equilibrium (probabilities  $1/3, 2/3$  and  $1/2, 1/2$ , respectively). Taking the mixed strategy

<sup>21</sup>Note that some of these additional features arise because playing 0 yields 0 regardless of what the opponent chooses to do. For instance, in more general games with strategic substitutability, i.e. ones that do not conform to figure 1, it can easily be the case that the players prefer different equilibria.

<sup>22</sup>In the standard Nash bargaining solution, the conflict payoffs are the payoffs players obtain when they fail to reach agreement. The bargaining weights represent the fraction of the surplus generated by the agreement that each player obtains. For instance, if absent agreement payoffs are  $(c_1, c_2)$  and with agreement total surplus is  $a$ , then for bargaining weights  $v, 1 - v$  for players 1 and 2 respectively, player 1 receives  $c_1 + (a - c)v$  and player 2 receives  $c_2 + (a - c)(1 - v)$ , where  $c = c_1 + c_2$ .

payoffs 0,0 to be the ‘conflict payoffs,’ we have  $p = 2v / (v + 1)$ .<sup>23</sup> In the normal form in figure 1 with strategic substitutability, we have in the multiplicity region that

$$p = p(e, v, x) = \frac{\{\tau_2(x, 0) + e_2\}v}{\{\tau_1(x, 0) + e_1\}(1 - v) + \{\tau_2(x, 0) + e_2\}v},$$

noting that  $p(e, v, x)$  is (if  $e \in S_m(x)$ ) increasing in  $v$  and between zero and one.

Like in section 2.3.1, in both of the above examples, the structural interpretation is lost and the range of counterfactuals that can be studied is diminished if  $u, v$  are collapsed into a single payoffs-independent unobservable random variable. In the collusion example, for instance, it requires two random variables to consider a counterfactual in which the bargaining weights are held fixed but in which different outcomes can arise for the same values of the payoff variables.

*2.3.3 Comments:* In both the coordination game and entry game examples, if one were willing to make assumptions about the specific equilibrium selection mechanism then such information could be exploited and confidence bands on counterfactual outcome probabilities could be tightened. We have chosen not to go this route to maintain generality.

In view of the similarity between the prediction problems for the coordination and entry games, from here on we only discuss the entry game case, noting that our results equally apply to the coordination game example and indeed to multiplayer examples that possess features of both games.<sup>24</sup>

### 3. Counterfactuals

We now consider thought experiments in which we consider what will happen if the game is played again under various scenarios and with (potentially) different covariate values. We denote the ex post variables by  $(e^*, p^*, u^*, v^*, x^*, y^*)$  which, except where otherwise noted, will be an independent copy of  $(e, p, u, v, x, y)$ .<sup>25</sup> The exceptions are that in different scenarios different combinations of the input variables are assumed to stay unchanged, which we explain in more detail below.

For given  $e, u, v$ , recall from (7) that  $y(e, u, v, x)$  is the value  $y$  would take if  $e = e, u = u, v = v, x = x$ . Thus,  $p^* = p(e^*, v^*, x^*)$  and  $y^* = y(e^*, u^*, v^*, x^*)$ . Using this notation, we can now consider various counterfactual outcomes, which differ depending on which combination of the conditions  $e^* = e, u^* = u, v^* = v$ , is applied: e.g. if we keep the unobserved payoff shifter unchanged but redraw all the other unobservables, then the

<sup>23</sup>Indeed, player 1’s expected payoff is then  $2v / (v + 1)$  and player 2’s is  $\{1 - 2v / (v + 1)\}2 = (2 - 2v) / (v + 1)$ , such that player 1’s share of the pie is  $v$ .

<sup>24</sup>For instance, it is conceivable that  $(1, 1, 0)$  is a Nash equilibrium in a three player game.

<sup>25</sup>So,  $x^*$  has the same support as  $x$ . If the counterfactual of interest is not in the support of  $x$ , then we would have to rely on extrapolation.



counterfactual of interest is  $y^* = y(e, \mathbf{u}^*, \mathbf{v}^*, \mathbf{x}^*)$ .<sup>26</sup> So there are up to eight different scenarios to consider, but we will focus on two cases, which we believe to be the most meaningful: keeping  $\mathbf{v}$  unchanged and keeping both  $\mathbf{e}$  and  $\mathbf{v}$  unchanged. The rationale for emphasizing these two cases is that we would like to know what will happen in a similar market, or indeed in the same market under different circumstances. In other words, we think of these two cases as representative. In the interest of completeness, we provide results on other counterfactuals, i.e. fixing other combinations of  $\mathbf{e}$ ,  $\mathbf{u}$ ,  $\mathbf{v}$  in appendix B,<sup>27</sup> and provide a brief summary thereof in section 3.3.

For each of the counterfactuals, we wish to construct a prediction of  $\mathbf{y}^*$ . We define  $q_c(y^* | x^*, x, y)$  to be the conditional probability that  $\mathbf{y}^* = y^*$  given  $\mathbf{x}^* = x^*$ ,  $\mathbf{x} = x$ ,  $\mathbf{y} = y$ , where the subscript  $c$  indicates which of  $\mathbf{e}^*$ ,  $\mathbf{u}^*$ ,  $\mathbf{v}^*$  are fixed. For instance,

$$\begin{cases} q(y^* | x^*, x, y) = \mathbb{P}\{y(\mathbf{e}^*, \mathbf{u}^*, \mathbf{v}^*, x^*) = y^* | \mathbf{x}^* = x^*, \mathbf{x} = x, \mathbf{y} = y\}, \\ q_v(y^* | x^*, x, y) = \mathbb{P}\{y(\mathbf{e}^*, \mathbf{u}^*, \mathbf{v}, x^*) = y^* | \mathbf{x}^* = x^*, \mathbf{x} = x, \mathbf{y} = y\}, \\ q_{ev}(y^* | x^*, x, y) = \mathbb{P}\{y(\mathbf{e}, \mathbf{u}^*, \mathbf{v}, x^*) = y^* | \mathbf{x}^* = x^*, \mathbf{x} = x, \mathbf{y} = y\}, \end{cases}$$

where  $q$  represents the case with no constraints and  $q_v, q_{ev}$  represent the cases where (some of) the unobservables are unchanged. We use both  $\mathbf{x}$  and  $\mathbf{y}$  to predict  $\mathbf{y}^*$ , because the values of both  $\mathbf{x}$  and  $\mathbf{y}$  contain information about  $\mathbf{p}$ .

The quantity  $q(y^* | x^*, x, y)$  is identified because it is equal to

$$\tilde{\pi}_{y^*}(x^*) = \mathbb{P}(y^* = y^* | \mathbf{x}^* = x^*) = \mathbb{P}(y = y^* | \mathbf{x} = x^*).$$

The *regression prediction*  $\tilde{\pi}_{y^*}(x^*)$  is simple, but it represents the case where none of the unobservables in the environment remains the same and is hence less interesting as a counterfactual exercise.

In sections 3.1 and 3.2 we discuss  $q_v$  and  $q_{ev}$  in greater detail. Define

$$\mu_m(x) = \mathbb{E}\{\mathbf{p} | \mathbf{e} \in S_m(x), \mathbf{x} = x\}. \quad (8)$$

Since  $\pi_y(x)$  and  $\pi_m(x)$  can be recovered from the payoff structure,  $\mu_m(x)$  can be identified from

$$\tilde{\pi}_{10}(x) = \pi_{10}(x) + \pi_m(x)\mu_m(x),$$

provided that  $\pi_m(x) > 0$ : if  $\pi_m(x) = 0$  then the value of  $\mu_m(x)$  is immaterial. In fact,  $\mu_m$  is the only identifiable feature of the conditional distribution of  $\mathbf{p}$  given  $\mathbf{e}$ ,  $\mathbf{x}$ , because the probability mass function  $\tilde{\pi}_y(x)$  depends on  $\rho(\cdot, \cdot, x)$  only through  $\mu_m(x)$ .

Throughout the remainder of the paper, we will frequently use the shorthand

$$\delta_y = \mathbb{1}\{y = (1, 0)\} - \mathbb{1}\{y = (0, 1)\}. \quad (9)$$

<sup>26</sup>Therefore, the ‘condition’  $\mathbf{e}^* = \mathbf{e}$  does not mean that we are ‘conditioning.’ It means that we keep  $\mathbf{e}$  *unchanged*, i.e. use the *same* random variable, in defining the counterfactual outcome. These two concepts are generally not the same. For instance, suppose that  $\xi_1, \xi_2$  are random variables with a standard exponential distribution. If  $\xi_2$  is the same random variable as  $\xi_1$  (which is the case we are considering) then  $\mathbb{E}\xi_2$  is by definition equal to  $\mathbb{E}\xi_1 = 1$ . But if e.g.  $\xi_1, \xi_2$  are independent random variables then  $\mathbb{E}(\xi_2 | \xi_2 = \xi_1) = 1/2$ . So there is a difference between the random variables being the same and conditioning on two different random variables having the same value.

<sup>27</sup>There are eight potential counterfactuals of interest: the regression prediction  $q$ ; the leading examples  $q_v$  and  $q_{ev}$ ; the predictions considered in appendix B  $q_e, q_{uv}$ , and  $q_{ev}$ ; and predictions that are omitted because of their similarity to ones that are discussed in the paper, namely  $q_u$  and  $q_{eu}$ .

**3.1 Case 1:  $v^* = v$ :** We now study identification of  $q_v(y^* | x^*, x, y)$ . Recall that this corresponds to the case in which the probability  $\mathbf{p}$  of selecting a particular equilibrium in the multiplicity region only varies because of changes in the values of the payoff variables.

Let

$$\begin{aligned} \rho(x, x^*) &= \text{Cov}\{\rho(\mathbf{e}, \mathbf{v}, x), \rho(\mathbf{e}^*, \mathbf{v}, x^*) \mid \mathbf{e} \in S_m(x), \mathbf{e}^* \in S_m(x^*)\} \\ &= \text{Cov}\{h(\mathbf{v}, x), h(\mathbf{v}, x^*)\}, \end{aligned} \quad (10)$$

where

$$h(\mathbf{v}, x) = \mathbb{E}\{\mathbf{p} \mid \mathbf{e} \in S_m(x), \mathbf{v} = v, \mathbf{x} = x\} = \mathbb{E}\{\rho(\mathbf{e}, \mathbf{v}, x) \mid \mathbf{e} \in S_m(x)\}. \quad (11)$$

The function  $\rho$  describes the predictability of  $\mathbf{p}^*$  by  $\mathbf{p}$ : if  $\rho(x, x^*)$  (which is nonnegative by construction) is large then we can learn a lot about the value of  $\mathbf{p}^*$  from the value of  $\mathbf{p}$ . Although we do not observe  $\mathbf{p}$ , the values of  $\mathbf{x}$ ,  $\mathbf{y}$  (which are observed) provide us with information about  $\mathbf{p}$ . It is intuitive that the only relevance of  $\mathbf{e}$  in this context is whether or not it belongs to the multiplicity region: since here  $\mathbf{e}^*$  is an independent copy of  $\mathbf{e}$  and  $\mathbf{e}$  is independent of everything else, we cannot learn anything about the effect of  $\mathbf{e}$  on the value of  $\mathbf{p}$  other than whether or not  $\mathbf{e} \in S_m(\mathbf{x})$ . Hence,  $\mathbf{e}$  is averaged out. Theorem 1 formalizes this intuition.

**Theorem 1.** Let

$$\tilde{\rho}_{yy^*}(x, x^*) = \delta_y \delta_{y^*} \frac{\pi_m(x) \pi_m(x^*)}{\tilde{\pi}_y(x)} \rho(x, x^*),$$

where  $\delta_y$  was defined in (9). Then,  $q_v(y^* | x^*, x, y) = \tilde{\pi}_{y^*}(x^*) + \tilde{\rho}_{yy^*}(x, x^*)$ .  $\square$

The results of theorem 1 are intuitive. The correction term  $\tilde{\rho}_{yy^*}(x, x^*)$  reflects what we can learn about  $\mathbf{v}$  from observing  $\mathbf{x}$ ,  $\mathbf{y}$ . Indeed, if  $y = (0, 0)$  or  $y = (1, 1)$  then we do not learn anything about the value of  $\mathbf{v}$  by observing  $\mathbf{x}$ ,  $\mathbf{y}$  since the value of  $\mathbf{v}$  is only of relevance in the multiplicity region: hence the correction term is equal to zero. Likewise, if  $y^* = (0, 0)$  or  $y^* = (1, 1)$  then any knowledge obtained about  $\mathbf{v}$  by observing  $\mathbf{x}$ ,  $\mathbf{y}$  is useless. In the remaining cases, if the new multiplicity region is large (and hence  $\pi_m(x^*)$  is large) then the correction term is large, also. If the probability that we are in the multiplicity region is large relative to the probability that  $\mathbf{y} = y$ , then again the correction term is large. Finally, if  $\rho(\mathbf{e}, \mathbf{v}, x)$  is highly correlated with  $\rho(\mathbf{e}^*, \mathbf{v}, x^*)$  conditional on  $\mathbf{e}$ ,  $\mathbf{e}^*$  belonging to their respective multiplicity regions, then again the correction term is large.

If the policy maker is only interested in the number of ‘entrants,’ not their identity, as in e.g. Bresnahan and Reiss (1991a), then  $q_v\{(1, 0) | x^*, x, y\} + q_v\{(0, 1) | x^*, x, y\}$  is the object of interest. For the present counterfactual, the prediction of the number of entrants is point-identified since theorem 1 implies that

$$q_v\{(1, 0) | x^*, x, y\} + q_v\{(0, 1) | x^*, x, y\} = \tilde{\pi}_{10}(x^*) + \tilde{\pi}_{01}(x^*). \quad (12)$$

So, if one is concerned whether or not there is a monopoly then the object of interest is point-identified. However, if the identity of the monopolist is of interest then theorem 1 implies that there is a (partially identified) correction term that needs to be dealt with. In

the counterfactual experiment of section 3.2 there will be a correction term even if one is only concerned with the number of entrants.

For the remainder of section 3.1, we will focus on the individual counterfactual  $q_v(y^* | x^*, x, y)$  rather than the sum (12). The size of the correction term in theorem 1 depends on a number of factors, but it can be bounded. The bound depends on  $B(x, x^*)$  given in theorem 2.

**Theorem 2.** For all  $x, x^*, 0 \leq \rho(x, x^*) \leq B(x, x^*)$ , where  $B(x, x^*) = \min\{\mu_m(x), \mu_m(x^*)\} - \mu_m(x)\mu_m(x^*)$ . Both bounds are sharp.  $\square$

An immediate implication of theorem 2 is that the regression prediction coincides with the lower bound of the sharp identified interval for  $q_v$ , which is not true for  $q_{ev}$ , as will be shown in section 3.2. Therefore, the maximum difference between the correct counterfactual object and the regression prediction coincides with the length of the sharp identified set of  $q_v\{(1, 0) | x^*, x, (1, 0)\}$ , which is equal to

$$\tilde{\rho}_{10,10}(x, x^*) = \frac{\pi_m(x)\pi_m(x^*)}{\tilde{\pi}_{10}(x)} B(x, x^*).$$

Note that  $\tilde{\rho}_{10,10}(x, x^*)$  can be made arbitrarily close to one by letting  $\pi_m(x), \pi_m(x^*) \rightarrow 1$  and  $\mu_m(x^*) = \mu_m(x) \rightarrow 0$ . Indeed,  $\tilde{\rho}_{10,10}(x, x^*)$  approaches  $\min\{1, \mu_m(x^*) / \mu_m(x)\} - \mu_m(x^*)$  as  $\pi_m(x), \pi_m(x^*) \rightarrow 1$ . More generally,  $\tilde{\rho}_{10,10}(x, x^*)$  varies with the ratio  $\pi_{10}(x) / \pi_m(x)$  as depicted in figure 6. Much the same applies to other combinations of  $y, y^* \in \{(1, 0), (0, 1)\}$ .

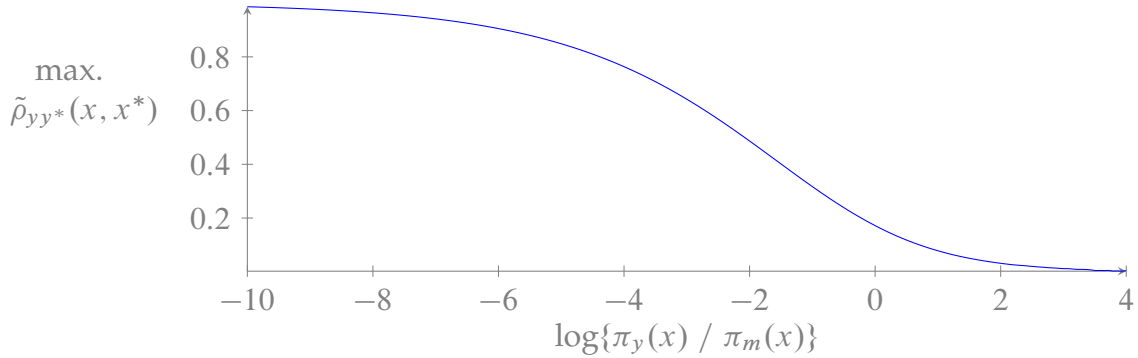


Figure 6: Maximum possible value of  $\tilde{\rho}_{yy^*}(x, x^*)$  as a function of  $\log\{\pi_y(x) / \pi_m(x)\}$  if  $\mu_m(x^*) = \mu_m(x)$  and  $y^* = y = (1, 0)$

**3.2 Case 2:  $e^* = e$  and  $v^* = v$ :** We now turn to the case where both  $e$  and  $v$  are fixed, i.e. the probability  $p$  of selecting a particular equilibrium under multiplicity varies only because of changes in the values of the observable payoff variables. In this case, the value of  $e$  is important for prediction.<sup>28</sup> For instance, if  $x^* = x$  then the  $S$ -regions are unchanged and  $p^* = p$ .<sup>29</sup> If  $x^* \neq x$  then the  $S$ -regions will be different and  $p^*$  can be different from

<sup>28</sup>The value of  $e$  would also be important for prediction if  $v$  were not fixed. This case is simpler and is discussed in theorem 10 in appendix B.

<sup>29</sup>By 'S-region' we mean one of the sets  $S(\cdot)$ .

$\mathbf{p}$ , but the values of  $\mathbf{x}$ ,  $\mathbf{y}$  still contain more information to aid in the prediction of  $\mathbf{y}^*$  than in the case discussed in section 3.1.

In section 3.1  $\mathbf{e}^*$  was an independent copy of  $\mathbf{e}$  and  $\pi_m(\mathbf{x})$ ,  $\pi_m(\mathbf{x}^*)$  were relevant for prediction. Now,  $\mathbf{e}^*$  is the same random variable as  $\mathbf{e}$  and hence the probability that  $\mathbf{e}$  belongs to the intersection of  $S_m(\mathbf{x})$  and  $S_m(\mathbf{x}^*)$  becomes relevant. We now introduce further notation to facilitate our analysis.

Let  $S_{my}(\mathbf{x}, \mathbf{x}^*) = S_m(\mathbf{x}) \cap S_y(\mathbf{x}^*)$ ,  $\pi_{my}(\mathbf{x}, \mathbf{x}^*) = \mathbb{P}\{\mathbf{e} \in S_{my}(\mathbf{x}, \mathbf{x}^*)\}$ , and let  $S_{mm}$ ,  $S_{yy^*}$ ,  $\pi_{mm}$ ,  $\pi_{yy^*}$ , etcetera, be analogously defined.<sup>30</sup> We further define  $S_r(\mathbf{x}) = S_r^y(\mathbf{x}) = \{S_y(\mathbf{x}) \cup S_m(\mathbf{x})\}^c$  and  $S_{r^*}(\mathbf{x}) = S_{r^*}^{y^*}(\mathbf{x}) = \{S_{y^*}(\mathbf{x}) \cup S_m(\mathbf{x})\}^c$ . In order to determine the bounds on  $q_{ev}$ , it matters where the  $S$ -regions intersect. Let

$$\rho(\mathbf{x}, \mathbf{x}^*) = \mathbb{E}\left[\text{Cov}\{\rho(\mathbf{e}, \mathbf{v}, \mathbf{x}), \rho(\mathbf{e}, \mathbf{v}, \mathbf{x}^*) \mid \mathbf{e}\} \mid \mathbf{e} \in S_{mm}(\mathbf{x}, \mathbf{x}^*)\right], \quad (13)$$

which is similar to (10) but imposes that  $\mathbf{e}$  and  $\mathbf{e}^*$  are the same random variable in both  $\rho$  functions. We use the same symbol  $\rho$  in both (10) and (13) for several reasons even though the expressions differ. First,  $\rho$  serves the same role of correcting for the information contained in  $\mathbf{x}$ ,  $\mathbf{y}$  about the value of  $\mathbf{p}^*$ : it is always clear which one is meant from the context. Further, if  $\rho(\mathbf{e}, \mathbf{v}, \mathbf{x})$  is flat in  $\mathbf{e}$ , then imposing the constraint of  $\mathbf{e}^* = \mathbf{e}$  makes the two expressions coincide. As will be shown in the following section, the maximum entropy solution for  $\rho$  is indeed flat in  $\mathbf{e}$ .

But we need an additional correction term to address the fact that  $\mathbf{x}$ ,  $\mathbf{y}$  contain information about the  $S(\mathbf{x}^*)$ -region to which  $\mathbf{e}$  belongs. This correction term can once again be expressed in terms of a covariance of constructed random variables. Theorem 3 contains the result, for which we need to introduce notation. Let  $c_y(\mathbf{e}, \mathbf{x}) = \mathbb{1}\{\mathbf{e} \in S_y(\mathbf{x})\}$ ,  $c_m(\mathbf{x}) = \mathbb{1}\{\mathbf{e} \in S_m(\mathbf{x})\}$ , and

$$b_y(\mathbf{e}, \mathbf{v}, \mathbf{x}) = \begin{cases} \rho(\mathbf{e}, \mathbf{v}, \mathbf{x}), & \mathbf{y} = (1, 0), \\ 1 - \rho(\mathbf{e}, \mathbf{v}, \mathbf{x}), & \mathbf{y} = (0, 1), \\ 0, & \text{otherwise.} \end{cases}$$

**Theorem 3.** Let  $a_y(\mathbf{e}, \mathbf{x}) = c_y(\mathbf{e}, \mathbf{x}) + c_m(\mathbf{e}, \mathbf{x})\mathbb{E}b_y(\mathbf{e}, \mathbf{v}, \mathbf{x})$  and

$$\begin{cases} \alpha_{yy^*}(\mathbf{x}, \mathbf{x}^*) = \frac{\text{Cov}\{a_y(\mathbf{e}, \mathbf{x}), a_{y^*}(\mathbf{e}, \mathbf{x}^*)\}}{\tilde{\pi}_y(\mathbf{x})}, \\ \tilde{\rho}_{yy^*}(\mathbf{x}, \mathbf{x}^*) = \delta_y \delta_{y^*} \frac{\pi_{mm}(\mathbf{x}, \mathbf{x}^*)}{\tilde{\pi}_y(\mathbf{x})} \rho(\mathbf{x}, \mathbf{x}^*). \end{cases}$$

Then, in view of (6),  $q_{ev}(y^* \mid x^*, x, y) = \tilde{\pi}_{y^*}(x^*) + \alpha_{yy^*}(\mathbf{x}, \mathbf{x}^*) + \tilde{\rho}_{yy^*}(\mathbf{x}, \mathbf{x}^*)$ .  $\square$

The intuition for the  $\tilde{\rho}$ -term in theorem 3 is similar to that in theorem 1: it reflects what we can learn about  $\mathbf{v}$  from observing  $\mathbf{y}$ ,  $\mathbf{x}$ . Likewise, the  $\alpha$  term captures the information for  $\mathbf{e}$  contained in  $\mathbf{y}$ ,  $\mathbf{x}$ . For example, consider first the case that  $\mathbf{y} = \mathbf{y}^* = (0, 0)$ . Then the

<sup>30</sup>Please note that the orders of the input arguments and the subindices are relevant here.

$\alpha_{yy^*}$  term simply increases the probability that  $y^* = (0, 0)$  because we know that  $y = (0, 0)$  implies that  $e_1, e_2$  are both small. In each of the other cases the intuition is merely a more complicated version of this argument.

Define

$$\gamma_y(x) = \mathbb{E}\{\ell_y(\mathbf{e}, \mathbf{v}, x) \mid \mathbf{e} \in S_m(x)\} \quad \text{and} \quad \phi_y(x) = \pi_m(x)\gamma_y(x).$$

So,  $\gamma_y(x)$  is  $\mu_m(x)$ ,  $1 - \mu_m(x)$ , or 0, depending on the value of  $y$ .

Now, suppose again that one is interested in the number of entrants rather than their identity. Then, theorem 3 implies that

$$\begin{aligned} q_{ev}\{(1, 0) \mid x^*, x, y\} + q_{ev}\{(0, 1) \mid x^*, x, y\} \\ = \tilde{\pi}_{10}(x^*) + \tilde{\pi}_{01}(x^*) + \frac{\text{Cov}\{a_y(\mathbf{e}, x), a_{10}(\mathbf{e}, x^*) + a_{01}(\mathbf{e}, x^*)\}}{\tilde{\pi}_y(x)}. \end{aligned} \quad (14)$$

The covariance on the right hand side in (14) is equal to

$$C_y(x, x^*) + \sum_{\tilde{y} \in \{(1,0), (0,1), m\}} [\pi_{y\tilde{y}}(x, x^*) - \pi_{\tilde{y}}(x^*)\{\pi_y(x) + \phi_y(x)\}],$$

where

$$C_y(x, x^*) = \sum_{\tilde{y} \in \{(1,0), (0,1), m\}} \pi_{m\tilde{y}}(x, x^*) \mathbb{E}\{\ell_y(\mathbf{e}, \mathbf{v}, x) \mid \mathbf{e} \in S_{m\tilde{y}}\}.$$

Therefore, unlike in the case discussed in section 3.1, the counterfactual probability of monopoly is not necessarily identified because  $C_y(x, x^*)$  can depend on  $p(\mathbf{e}, \mathbf{v}, x)$ . Note however that  $p(\mathbf{e}, \mathbf{v}, x^*)$  is irrelevant for the counterfactual probability of monopoly. The following theorem provides sharp bounds of  $C_y(x, x^*)$ .

**Theorem 4.** For all  $x, x^*$ , we have

$$\begin{aligned} \max\{0, \phi_y(x) - \pi_{m,00}(x, x^*) - \pi_{m,11}(x, x^*)\} &\leq C_y(x, x^*) \\ &\leq \min\{\phi_y(x), \pi_m(x) - \pi_{m,00}(x, x^*) - \pi_{m,11}(x, x^*)\}. \end{aligned}$$

The bounds are sharp. □

If  $y$  is either  $(0, 0)$  or  $(1, 1)$  then point identification obtains: e.g. the probability of monopoly in the counterfactual conditional on no monopoly ex ante is point-identified. Further, if the  $S$ -regions do not change (e.g. if  $x = x^*$ ), then we have point identification, also.

However, the situation is less clear if there is a monopoly ex ante. It is then not generally possible to determine the counterfactual probability of monopoly if  $x$  is changed to  $x^*$  (e.g. by a policy maker). For instance, suppose that  $y = (1, 0)$  and that  $\pi_m(x) = 1$ . Suppose further that the policy is such that  $\pi_{m,00}(x, x^*) + \pi_{m,11}(x, x^*) = 1 - \nu < 1$ . Then, the length of the identified interval of  $q_{ev}\{(1, 0) \mid x^*, x, (1, 0)\} + q_{ev}\{(0, 1) \mid x^*, x, (1, 0)\}$  equals

$$\min\{1, \nu / \mu_m(x)\} - \max\{0, 1 - (1 - \nu) / \mu_m(x)\},$$

which converges to 1 as  $\mu_m(x)$  approaches 0.

We now turn to the problem of determining sharp identified bounds for  $q_{ev}(y^* | x^*, x, y)$  for the balance of this section. Define

$$L_{yy^*}(x, x^*) = \max. \text{ of } \begin{cases} 0, \\ \phi_y(x) - \pi_{mr^*}(x, x^*) - \pi_{mm}(x, x^*), \\ \phi_{y^*}(x^*) - \pi_{rm}(x, x^*) - \pi_{mm}(x, x^*), \\ \phi_y(x) + \phi_{y^*}(x^*) - \pi_{mr^*}(x, x^*) - \pi_{rm}(x, x^*) - \pi_{mm}(x, x^*), \end{cases}$$

and

$$U_{yy^*}(x, x^*) = \min. \text{ of } \begin{cases} \phi_y(x) + \phi_{y^*}(x^*), \\ \phi_y(x) + \pi_{ym}(x, x^*), \\ \phi_{y^*}(x^*) + \pi_{my^*}(x, x^*), \\ \pi_{mm}(x, x^*) + \pi_{ym}(x, x^*) + \pi_{my^*}(x, x^*), \end{cases} \quad (15)$$

which will enter into the lower and upper bound formulas for  $q_{ev}$ .  $L_{yy^*}(x, x^*)$  and  $U_{yy^*}(x, x^*)$  depend on identified objects only. Indeed,  $L_{yy^*}(x, x^*)$  and  $U_{yy^*}(x, x^*)$  are determined by a combination of  $\pi$ -values and the values of  $\phi_y(x)$  and  $\phi_{y^*}(x^*)$ , as is illustrated in figure 7. Depending on the values of  $\phi_y(x)$  and  $\phi_{y^*}(x^*)$ , different bounds are binding.

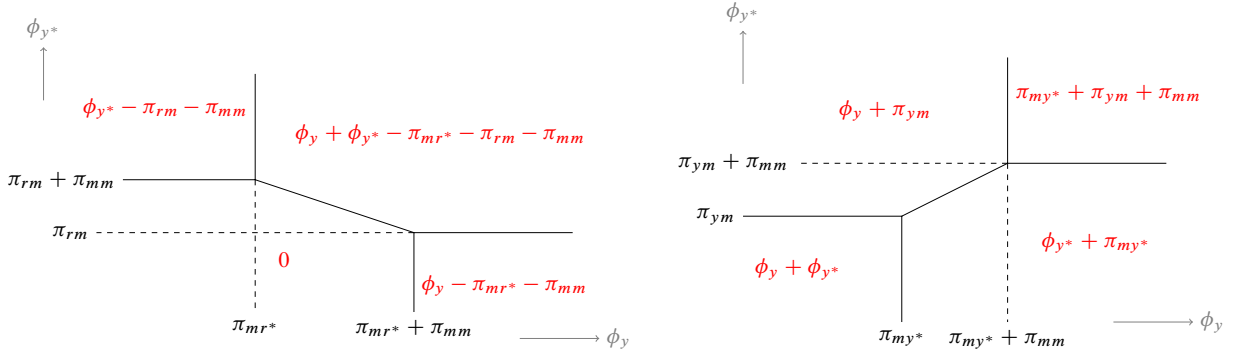


Figure 7: Regions of  $(\phi_y(x), \phi_{y^*}(x^*))$  and the corresponding values of  $L_{yy^*}(x, x^*)$ ,  $U_{yy^*}(x, x^*)$ .

We are now ready to describe the sharp identified bounds of  $q_{ev}(y^* | x^*, x, y)$ .

**Theorem 5.** If  $x \neq x^*$ , then

$$\pi_{yy^*}(x, x^*) + L_{yy^*}(x, x^*) \leq q_{ev}(y^* | x^*, x, y) \tilde{\pi}_y(x) \leq \pi_{yy^*}(x, x^*) + U_{yy^*}(x, x^*).$$

If  $x = x^*$ , then

$$q_{ev}(y^* | x, x^*, y) \in \begin{cases} \{\mathbb{1}(y = y^*)\}, & \text{if } \delta_y \delta_{y^*} = 0, \\ [\{\pi_y(x) + \pi_m(x) \gamma_y^2(x)\} / \tilde{\pi}_y(x), 1], & \text{if } \delta_y \delta_{y^*} = 1, \\ [0, \{\pi_m(x) \gamma_y(x) \gamma_{y^*}(x)\} / \tilde{\pi}_y(y)], & \text{if } \delta_y \delta_{y^*} = -1. \end{cases}$$

The bounds are sharp in both cases. □

The formulas for the bounds in theorem 5 are different depending on whether or not  $x = x^*$  and the bounds for  $x = x^*$  do not obtain as a limit if  $x \rightarrow x^*$ : if  $\delta_y \delta_{y^*} = 1$  (-1) then the lower (upper) bound is discontinuous at  $x = x^*$ . This fact can be explained as follows. If  $x = x^*$  then  $\mathbf{p} = \mathbf{p}^*$ . If  $x \neq x^*$  then  $\mathbf{p}$  and  $\mathbf{p}^*$  are the outputs of *different* functions of the *same* random variables  $\mathbf{e}, \mathbf{v}$ :  $\mathbf{p} = \rho(\mathbf{e}, \mathbf{v}, x)$  and  $\mathbf{p}^* = \rho(\mathbf{e}, \mathbf{v}, x^*)$ . Even if  $\rho(e, v, x)$  were continuous in  $x$  (which we do not assume) then  $|\rho(e, v, x^*) - \rho(e, v, x)| / \|x - x^*\|$  can still be arbitrarily large. Since  $\rho$  need not be monotonic in  $e$ , the covariance between  $\mathbf{p}$  and  $\mathbf{p}^*$  can be negative, unlike the variance of  $\mathbf{p}$ .

If the function  $\rho$  is flat in  $e$ , then  $\mathbf{p}$  and  $\mathbf{p}^*$  cannot be negatively correlated, and consequently the formulas for the bounds in theorem 5 for the cases  $x = x^*$  and  $x \neq x^*$  then coincide. However, even in that case,  $q_{ev}(y^* | x^*, x, y)$  and  $q_v(y^* | x^*, x, y)$  generally have different values, because for  $q_{ev}(y^* | x^*, x, y)$  it matters where the  $S$ -regions intersect.

Please note that the bounds in theorem 5 do not generally contain the regression prediction  $\tilde{\pi}_{y^*}(x^*)$ : see section 7.3 for examples. Therefore, the regression prediction is a poor choice if the object of interest is  $q_{ev}$ .

**3.3 Other cases:** As announced, appendix B contains rigorous results for the remaining counterfactual predictions, namely  $q_u, q_{uv}, q_e, q_{eu}, q_{euv}$ . The infeasible predictions themselves generally differ from those that we derived in sections 3.1 and 3.2. Indeed, if  $y = y^* = (1, 0)$  then

$$\begin{aligned} \tilde{\pi}_y(x)q.(y^* | x^*, x, y) &= \pi_y(x)\pi_{y^*}(x^*) + \pi_y(x)\pi_m(x^*)\mu_m(x^*) + \pi_{y^*}(x^*)\pi_m(x)\mu_m(x) + \\ \pi_m(x)\pi_m(x^*) &\times \begin{cases} \mu_m(x)\mu_m(x^*), & q. = q, \\ \mathbb{E}\{k(\mathbf{v}, x)k(\mathbf{v}, x^*)\}, & q. = q_v, \\ \mathbb{E}[\min\{\rho(\mathbf{e}, \mathbf{v}, x), \rho(\mathbf{e}^*, \mathbf{v}^*, x^*)\} | \mathbf{e} \in S_m(x), \mathbf{e}^* \in S_m(x^*)], & q. = q_u, \\ \mathbb{E}[\min\{\rho(\mathbf{e}, \mathbf{v}, x), \rho(\mathbf{e}^*, \mathbf{v}, x^*)\} | \mathbf{e} \in S_m(x), \mathbf{e}^* \in S_m(x^*)], & q. = q_{uv}, \end{cases} \end{aligned}$$

and

$$\begin{aligned} \tilde{\pi}_y(x)q.(y^* | x^*, x, y) &= \pi_{yy^*}(x, x^*) + \pi_{ym}(x, x^*)\mathbb{E}\{\rho(\mathbf{e}, \mathbf{v}^*, x^*) | \mathbf{e} \in S_{ym}(x, x^*)\} + \\ &\pi_{my^*}(x, x^*)\mathbb{E}\{\rho(\mathbf{e}, \mathbf{v}, x) | \mathbf{e} \in S_{my^*}(x, x^*)\} + \\ \pi_{mm}(x, x^*) &\times \begin{cases} \mathbb{E}\{\bar{p}(\mathbf{e}, x)\bar{p}(\mathbf{e}, x^*) | \mathbf{e} \in S_{mm}(x, x^*)\}, & q. = q_e, \\ \mathbb{E}[\rho(\mathbf{e}, \mathbf{v}, x)\rho(\mathbf{e}, \mathbf{v}, x^*) | \mathbf{e} \in S_{mm^*}(x, x^*)], & q. = q_{ev}, \\ \mathbb{E}[\min\{\rho(\mathbf{e}, \mathbf{v}, x), \rho(\mathbf{e}, \mathbf{v}^*, x^*)\} | \mathbf{e} \in S_{mm}(x, x^*)], & q. = q_{eu}, \\ \mathbb{E}[\min\{\rho(\mathbf{e}, \mathbf{v}, x), \rho(\mathbf{e}, \mathbf{v}, x^*)\} | \mathbf{e} \in S_{mm}(x, x^*)], & q. = q_{euv}, \end{cases} \end{aligned}$$

where  $\bar{p}(e, x) = \mathbb{E}\rho(e, \mathbf{v}, x)$ . The reason for the dichotomy between the cases in which  $\mathbf{e}$  is or is not fixed is that  $\mathbf{e}$  is vector-valued, that  $\rho$  is not necessarily monotonic in  $e$ , and that the multiplicity regions live in  $\mathbf{e}$ -space.

One notable conclusion from the above two displayed equations is that the prediction

probabilities can be ordered, e.g.

$$\begin{cases} q(y^* | x^*, x, y) \leq q_v(y^* | x^*, x, y) \leq q_{uv}(y^* | x^*, x, y), \\ q_e(y^* | x^*, x, y) \leq q_{ev}(y^* | x^*, x, y) \leq q_{euv}(y^* | x^*, x, y), \end{cases} \quad (16)$$

which is intuitive. However, if  $x \neq x^*$ , as is shown in theorems 10 to 14 in appendix B, then the bounds on  $q_u, q_{uv}$  coincide with those on  $q_v$  and the bounds on  $q_e, q_{eu}, q_{euv}$  coincide with those on  $q_{ev}$ .<sup>31</sup> Since the bounds coincide, so do the midpoint predictions. Thus, despite the fact that the  $q$ -functions can be ordered as indicated in (16), the midpoint predictions are the same in each case.<sup>32</sup> The maximum entropy solution proposed in section 4 does not have this unfortunate feature and neither does the minmax solution with entropy constraint in section 5 for most constraint values.

## 4. Maximum entropy

The principle of maximum entropy is the notion that among the probability distributions that satisfy all testable restrictions available the probability distribution that best represents the current state of knowledge is the one that maximizes the entropy (Jaynes, 1957a; Jaynes, 1957b). The entropy measures the amount of uncertainty that a probability distribution represents.<sup>33</sup>

In our case, this entails finding the joint density of  $e, p, x$  which maximizes the entropy: recall from (6) that  $p$  is the quantile function corresponding to the conditional density  $f$  of  $p$  given  $e, x$ . Our notation suggests that  $p$  is continuously distributed, but discrete distributions obtain as limit cases: we do *not* assume any regularity on  $f$  such as smoothness or boundedness. For ease of exposition, we will further take  $x$  to be continuously distributed but the nature of the distribution of  $x$  is immaterial.

We continue using the same environment as before, i.e. we take the distributions of  $e, x$  as given and take  $e, x$  to be independent. Further, recall that  $\mu_m(x)$ , defined in (8), is identified for all  $x$ , and hence it imposes another constraint on  $f$ . Thus, using all information available, we consider maximizing the *conditional entropy*<sup>34</sup>

$$\begin{aligned} f^* = \operatorname{argmax}_f & - \iiint_0^1 f(p | e, x) \log f(p | e, x) \, dp \, f_e(e) f_x(x) \, de \, dx \\ \text{subject to} & \begin{cases} \forall e, x : \int_0^1 f(p | e, x) \, dp = 1, \\ \forall x : \int_{S_m(x)} \int_0^1 p f(p | e, x) \, dp \, f_e(e) \, de = \mu_m(x) \pi_m(x). \end{cases} \end{aligned} \quad (17)$$

<sup>31</sup>We focus on the case  $x \neq x^*$  because  $q_{euv}(y^* | x, x, y) = \mathbb{1}(y^* = y)$ .

<sup>32</sup>The midpoint predictions are different if one compares fixing  $e$  and not fixing  $e$ , however.

<sup>33</sup>Admittedly, one could define alternative such measures; our use of entropy is motivated by it being the dominant choice in information theory.

<sup>34</sup>The problem can be equivalently formulated by the unconditional entropy of the joint density of  $(e, p, x)$ .



The optimization problem in (17) can be solved using standard constrained optimization techniques (see e.g. Cover and Thomas, 2012). The result is provided in theorem 6 below. Let

$$I(\lambda) = \int_0^1 \exp(p\lambda) dp, \quad \mathcal{L}(\lambda) = \log I(\lambda). \quad (18)$$

**Theorem 6.** The solution to (17) is given by

$$f^*(p | e, x) = \begin{cases} \mathbb{1}(0 \leq p \leq 1), & e \notin S_m(x), \\ A\{p, \lambda_m(x)\} \mathbb{1}(0 \leq p \leq 1), & e \in S_m(x), \end{cases} \quad (19)$$

where

$$\lambda_m(x) = \underset{\lambda}{\operatorname{argmin}} \{ \mathcal{L}(\lambda) - \mu_m(x)\lambda \}, \quad A(p, \lambda) = \exp(p\lambda) / I(\lambda). \quad (20)$$

□

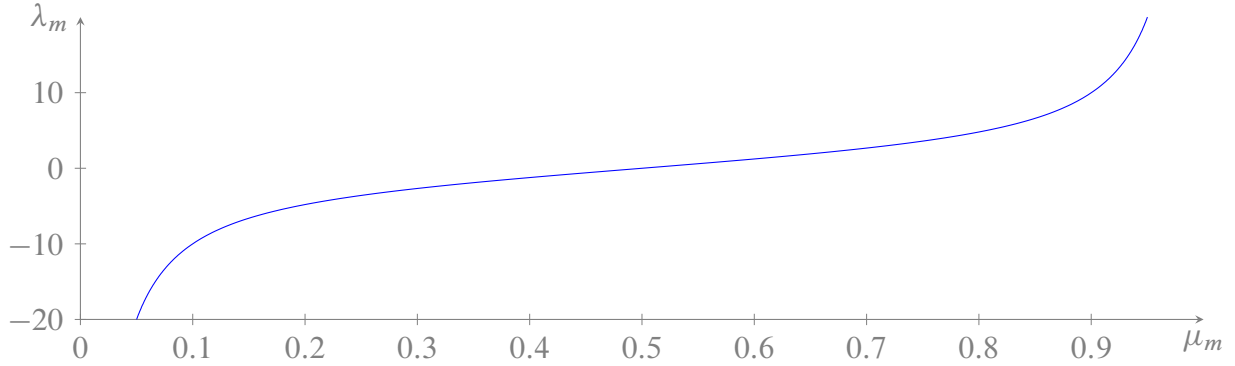


Figure 8:  $\lambda_m$  as a function of  $\mu_m$ .

The relationship between  $\lambda_m$  and  $\mu_m$  is depicted in figure 8. We can see from the first order condition of the minimization problem in (20) that the solution  $\lambda_m$  satisfies

$$\mu_m = \mathcal{L}'(\lambda_m) = \begin{cases} \frac{1}{2}, & \lambda_m = 0, \\ \frac{1}{1 - \exp(-\lambda_m)} - \frac{1}{\lambda_m}, & \lambda_m \neq 0, \end{cases}$$

where it should be noted that  $\mathcal{L}'$  is positive and well-behaved at zero.<sup>35</sup>

Theorem 6 implies that  $f^*(p | e, x)$  is different depending on whether or not  $e \in S_m(x)$ , but does not otherwise depend on the value of  $e$ . For  $e \notin S_m(x)$ , we have no information about  $f^*(p | e, x)$  and hence maximum entropy produces a uniform distribution. For  $e \in S_m(x)$ , we only have information about the value of  $\mu_m(x)$ . If  $\mu_m(x) = 0.5$  then there is nothing to suggest that the conditional distribution is not uniform. Otherwise, the density must be adjusted to accommodate the value of  $\mu_m(x)$ , as depicted in figure 9. In all cases,

<sup>35</sup>For instance,  $\mathcal{L}'$  is continuous and differentiable at zero.

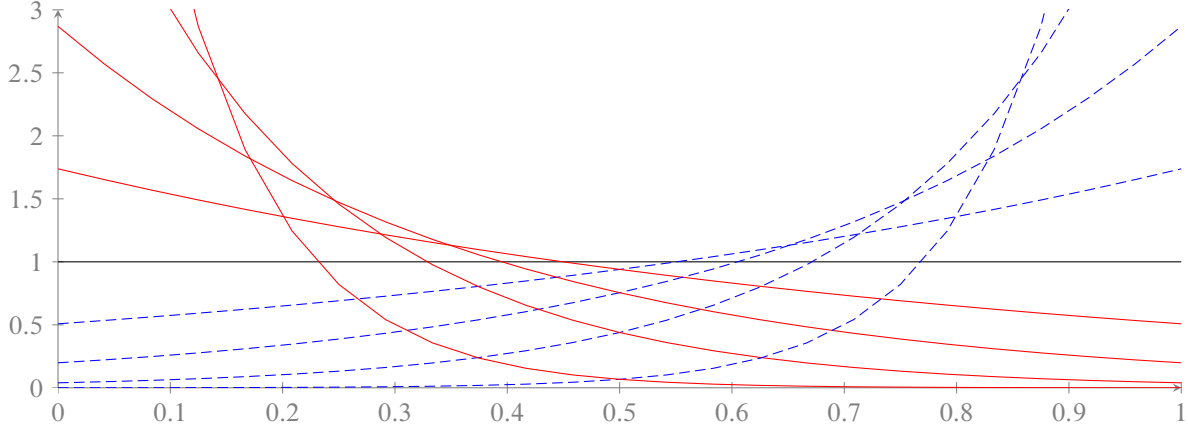


Figure 9: Conditional densities  $f^*(p | e, x)$  for  $e \in S_m(x)$  for values of  $\mu_m = 0.1, \dots, 0.9$ , with blue dashed curves corresponding to  $\mu_m > 0.5$  and red solid curves to  $\mu_m < 0.5$ .

the maximum entropy solution picks  $f^*$  such that the conditional density of  $p$  given  $e, x$  is closest to a uniform as measured by the entropy. The conditional density function becomes steeper as  $|\mu_m(x) - 0.5|$  increases.

Recall that  $p$  is the quantile function corresponding to  $f_{p|e, x}$ .

**Theorem 7.** The function  $p$  corresponding to  $f = f^*$  is the function  $p^*$  given by

$$p^*(e, v, x) = \begin{cases} v, & e \notin S_m(x), \\ k^*(v, x), & e \in S_m(x), \end{cases}$$

where the function  $k^*$  corresponds to  $k$  defined in (11) and is given by

$$k^*(v, x) = \begin{cases} \frac{\log(1 + v[\exp\{\lambda_m(x)\} - 1])}{\lambda_m(x)}, & \lambda_m(x) \neq 0, \\ v, & \lambda_m(x) = 0. \end{cases} \quad \square$$

It is straightforward to calculate the correction terms in theorems 1 and 3 that are implied by  $p^*$  (and  $k^*$ ). Indeed, since  $p^*$  is flat in  $e$  when  $e \in S_m(x)$ , the maximum entropy correction terms corresponding to  $\tilde{\rho}_{yy^*}(x, x^*)$  in theorems 1 and 3 coincide. We provide the general formulas in appendix C.

The function  $p^*$  is not generally continuous in  $e$  or  $x$  since the behavior of  $p^*$  is different inside and outside of  $S_m(x)$ . However, for  $e$  in the interior of  $S_m(x)$ ,  $p^*(e, v, x)$  is continuous in  $x$  and flat in  $e$ .

To see how the maximum entropy predictions differ across counterfactual experiments, we now provide the maximum entropy analogs to the counterfactual prediction formulas presented in section 3.3. If again  $y = y^* = (1, 0)$  and  $x \neq x^*$  then

$$\tilde{\pi}_y(x)q.(y^* | x^*, x, y) = \pi_y(x)\pi_{y^*}(x^*) + \pi_y(x)\pi_m(x^*)\mu_m(x^*) + \pi_m(x)\pi_{y^*}(x^*)\mu_m(x) +$$

$$\pi_m(x)\pi_m(x^*) \times \begin{cases} \mu_m(x)\mu_m(x^*), & q. = q, \\ \mathbb{E}\{k^*(\mathbf{v}, x)k^*(\mathbf{v}, x^*)\}, & q. = q_v, \\ \mathbb{E}[\min\{k^*(\mathbf{v}, x), k^*(\mathbf{v}^*, x^*)\}], & q. = q_u, \\ \min\{\mu_m(x), \mu_m(x^*)\}, & q. = q_{uv}. \end{cases}$$

and

$$\begin{aligned} \tilde{\pi}_y(x)q.(y^* | x^*, x, y) &= \pi_{yy^*}(x, x^*) + \pi_{ym}(x, x^*)\mu_m(x^*) + \pi_{my^*}(x, x^*)\mu_m(x) \\ &+ \pi_{mm}(x, x^*) \times \begin{cases} \mu_m(x)\mu_m(x^*), & q. = q_e, \\ \mathbb{E}\{k^*(\mathbf{v}, x)k^*(\mathbf{v}, x^*)\}, & q. = q_{ev}, \\ \mathbb{E}[\min\{k^*(\mathbf{v}, x), k^*(\mathbf{v}^*, x^*)\}], & q. = q_{eu}, \\ \min\{\mu_m(x), \mu_m(x^*)\}, & q. = q_{euv}. \end{cases} \end{aligned}$$

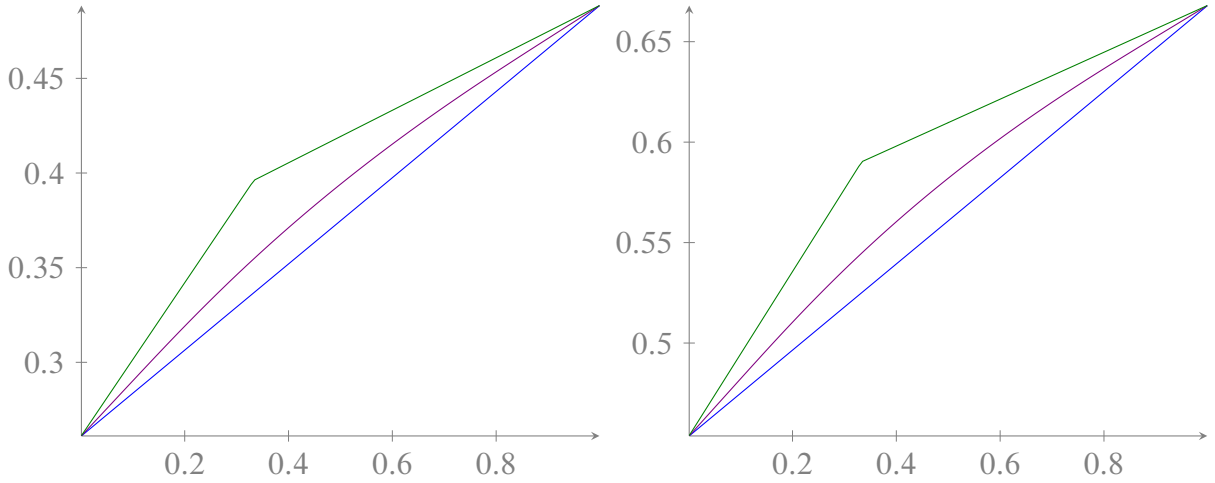


Figure 10: Maximum entropy solutions for  $q.$ ,  $q_v$ ,  $q_{uv}$  (left graph, bottom to top) and  $q_e$ ,  $q_{ev}$ ,  $q_{euv}$  (right graph, bottom to top) as a function of  $\mu_m(x^*)$  if  $y = y^* = (1, 0)$ ,  $\mu_m(x) = 1/3$ ,  $e_1, e_2$  independent  $N(0, 1)$ , shift in  $x$  value moves multiplicity region from  $[-1, -1]^2$  to  $[0, 2]^2$ .

For one simple design, the maximum entropy predictions are depicted in figure 10. The predictions are ordered the way they should be (see section 3.3), unlike the midpoint solutions. The kink in both graphs is due to the minimum function and located where  $\mu_m(x^*) = \mu_m(x)$ .

In section 7, we will compare the maximum entropy solutions with various alternatives using a number of stylized examples.

## 5. Minmax with entropy

As mentioned in the introduction, the usual minmax approach with a symmetric loss function yields an uninteresting solution which is moreover difficult to defend. We now

consider a generalized minmax problem, where we minimize the maximum risk subject to the corresponding entropy value being bounded below by a specified value. It will become clear that the midpoint and maximum entropy predictions are special cases.

We focus our attention on  $q_v(y^* | x, x^*, y; \rho)$ , where our notation makes it explicit that the counterfactual prediction depends on  $\rho$  which is only partially identified.<sup>36</sup> The problem is then to choose a function of ‘identifiable objects’ from the identified set of  $\rho$ . This function can then be used to produce  $q_v(y^* | x, x^*, y; \rho)$ . The exercise is at the population level in that the objective is to find a function of (population) distributions, not a sample. So there is no difference between risk and loss here.

Fix  $\chi = (x, x^*, y, y^*)$ . The aggregate entropy of  $\rho$  at  $x, x^*$  is

$$\mathcal{E}(\rho) = \sum_{\tilde{x} \in \{x, x^*\}} \int_0^1 \log \partial_v \rho(e, v, \tilde{x}) dv f_e(e) de. \quad (21)$$

Since  $\rho(e, v, \tilde{x})$  is the  $v$ -conditional quantile of  $\mathbf{p}$  given  $\mathbf{e} = e$  and  $\mathbf{x} = \tilde{x}$ ,  $\mathcal{E}(\rho)$  is simply sum of two conditional entropy values of  $\mathbf{p}$ , fixing  $x$  and  $x^*$ , i.e.

$$\mathcal{E}(\rho) = - \sum_{\tilde{x} \in \{x, x^*\}} \int_0^1 f_{\mathbf{p}|e, \mathbf{x}}(p | e, \tilde{x}) \log f_{\mathbf{p}|e, \mathbf{x}}(p | e, \tilde{x}) dp f_e(e) de. \quad (22)$$

The expression in (21) is more convenient for us than the one in (22) because the risk is more easily expressed as a function of  $\rho$ . Let  $\mathcal{P}$  be the class of functions  $\rho$  that is nondecreasing in  $v$  and satisfies the mean restrictions

$$\int_{S_m(\tilde{x})} \int_0^1 \rho(e, v, \tilde{x}) dv f_e(e) de = \mu_m(\tilde{x}) \pi_m(\tilde{x}), \quad \tilde{x} = x, x^*$$

as well as  $0 \leq \rho(e, v, \tilde{x}) \leq 1$  for all  $e, v$  and  $\tilde{x} = x, x^*$ .

Now let  $\ell$  be a symmetric loss function:  $\ell(d, \tilde{d}) = (d - \tilde{d})^2$  is a common choice. Consider the following optimization problem: for some chosen value  $\mathcal{E}^*$ ,

$$\min_d \max_{\substack{\rho \in \mathcal{P} \\ \mathcal{E}(\rho) \geq \mathcal{E}^*}} \ell\{d, q_v(\rho)\}, \quad (23)$$

where  $q_v(\rho)$  is shorthand for  $q_v(y^* | x, x^*, y; \rho)$ . As noted, if  $\mathcal{E}^* = -\infty$ , then the solution to (23) is simply the midpoint of the bounds in theorems 1 and 2. If  $\mathcal{E}^*$  is chosen to be the maximum achievable value, then the solution to (23) is equal to the maximum entropy prediction, i.e. the counterfactual prediction using  $\rho$  given in theorem 7. Thus, the midpoint and the maximum entropy predictions are special cases of the present approach.

Since the loss function is symmetric, the solution to (23) will be the midpoint of the identified set for  $q_v(\rho)$  restricted to the class of functions  $\rho \in \mathcal{P}$  for which  $\mathcal{E}(\rho) \geq \mathcal{E}^*$ . We therefore present the functions  $\rho_{-1}$  and  $\rho_1$  that correspond to the maximum and minimum attainable values for  $q_v$ , respectively: we can then use  $\rho_\sigma$ ,  $\sigma = \pm 1$ , together with (10) and theorem 1 to compute the solution. We now describe the solutions  $\rho_{-1}, \rho_1$ .

<sup>36</sup>We provide results for  $q_v$ : results for the other counterfactual predictions are messier.

Define

$$\epsilon_\sigma(z) = \int_0^z \exp(-\sigma t^2) dt.$$

Further, for  $j = 0, 1, 2$  and  $\kappa_0 < \kappa_1$ , let

$$I_{\sigma j}(\kappa_0, \kappa_1) = \int_{\kappa_0}^{\kappa_1} t^j \exp(-\sigma t^2) dt / \{\epsilon_\sigma(\kappa_1) - \epsilon_\sigma(\kappa_0)\}.$$

Finally, let  $\kappa_{\sigma 0}^\circ, \kappa_{\sigma 1}^\circ$  be the maximizers ( $\sigma = -1$ ) or the minimizers ( $\sigma = 1$ ) of

$$\{I_{\sigma 2}(\kappa_0, \kappa_1) - I_{\sigma 1}^2(\kappa_0, \kappa_1)\} \exp\{\mathcal{E}^* - 2\sigma I_{\sigma 2}(\kappa_0, \kappa_1)\} / \{\epsilon_\sigma(\kappa_1) - \epsilon_\sigma(\kappa_0)\}^2$$

subject to

$$\begin{cases} 2\sigma I_{\sigma 2}(\kappa_0, \kappa_1) + 2 \log[\{\epsilon_\sigma(\kappa_1) - \epsilon_\sigma(\kappa_0)\} / \{\kappa_1 - 2\sigma I_{\sigma 2}(\kappa_0, \kappa_1)\}] = \mathcal{E}^* - \log(1 - \mu_m)(1 - \mu_m^*), \\ \{\kappa_1 - I_{\sigma 1}(\kappa_0, \kappa_1)\} / \{I_{\sigma 1}(\kappa_0, \kappa_1) - \kappa_0\} \geq \max\{(1 - \mu_m)/\mu_m, (1 - \mu_m^*)/\mu_m^*\}, \end{cases}$$

where  $\mu = \mu_m(x)$  and  $\mu^* = \mu_m(x^*)$ .

**Theorem 8.** If a solution to (23) exists and  $\mathcal{E}^* > -\infty$ , then the solution is given by

$$d^* = \tilde{\pi}_y(x^*) + \delta_y \delta_{y^*} \pi_m(x) \pi_m(x^*) \{\rho_{-1}^\circ(x, x^*) + \rho_1^\circ(x, x^*)\} / \{2\tilde{\pi}_y(x)\},$$

where  $\rho_\sigma^\circ(x, x^*) = \int_0^1 p_\sigma(v) \rho_\sigma^*(v) dv - \mu_m(x) \mu_m(x^*)$  with

$$\begin{cases} p_\sigma(v) = \mu_m(x) + \{1 - \mu_m(x)\} \{\kappa_\sigma^\circ(v) - I_{\sigma 1}^\circ\} / (\kappa_{\sigma 1}^\circ - I_{\sigma 1}^\circ), \\ p_\sigma^*(v) = \mu_m(x^*) + \{1 - \mu_m(x^*)\} \{\kappa_\sigma^\circ(v) - I_{\sigma 1}^\circ\} / (\kappa_{\sigma 1}^\circ - I_{\sigma 1}^\circ), \end{cases}$$

where  $I_{\sigma 1}^\circ = I_{\sigma 1}(\kappa_{\sigma 0}^\circ, \kappa_{\sigma 1}^\circ)$  and  $\kappa_\sigma^\circ(v) = \epsilon_\sigma^{-1}(c_{\sigma 1}^\circ v + c_{\sigma 2}^\circ)$  with  $c_{\sigma 1}^\circ = \epsilon_\sigma(\kappa_{\sigma 1}^\circ) - \epsilon_\sigma(\kappa_{\sigma 0}^\circ)$  and  $c_{\sigma 2}^\circ = \epsilon_\sigma(\kappa_{\sigma 0}^\circ)$ .  $\square$

Theorem 8 provides a general solution to (23) but it does not establish the uniqueness of the function  $p_\sigma$  that delivers the solution. Note further that we have fixed  $x, x^*$  and take them to be different.

Figure 11 depicts the results of a numerical example and shows  $\rho_\sigma^\circ(x, x^*)$  as a function of  $\mathcal{E}^*$ . The entropy constraint is not binding if  $\mathcal{E}^* = 0$  and the solution is then simply the midpoint prediction. Otherwise the entropy constrained minmax procedure yields the midpoint of the entropy-constrained identified set, which is an interval whose length shrinks as  $\mathcal{E}^*$  increases. Note that the upper bound adjusts faster than the lower bound; this is consistent with our finding in section 7 that the midpoint predictions tend to exceed the maximum entropy ones.

## 6. Dirichlet approach

**6.1 Overview:** In this section we discuss yet another possibility involving Dirichlet processes. The approach has a decision-theoretic foundation, but its implementation is considerably more complicated than maximum entropy and it does not do much to mitigate the

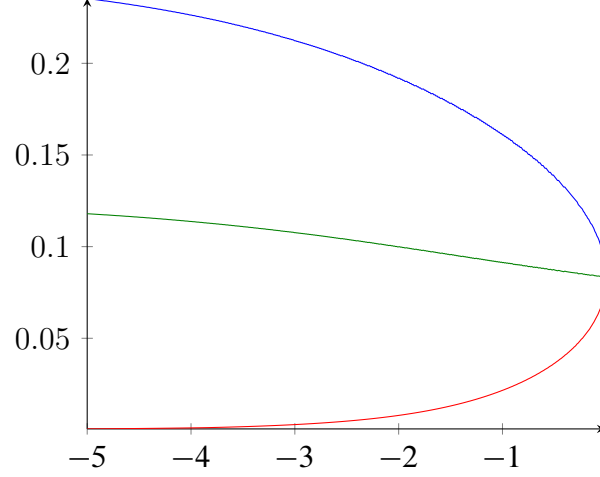


Figure 11:  $\rho_\sigma^\circ(x, x^*)$  as a function of  $\mathcal{E}^*$  with  $\mu(x) = 0.5$  and  $\mu(x^*) = 0.25$ . The red line corresponds to  $\sigma = 1$  and the blue to  $\sigma = -1$ . The green line is the midpoint between the blue and red lines.

inherent arbitrariness of the problem at hand. We will compare the results to the maximum entropy ones in section 7.

Recall that the prediction of interest is a deterministic function of  $p$ ; we will be explicit about this fact in our discussion below. We will treat  $p$  as a parameter, consider a prior distribution for  $p$ , and look for the prediction that minimizes average risk. Since  $p(e, \cdot, x)$  is a conditional quantile function (satisfying a mean constraint), we need to assign a probability measure on the (constrained) space of (conditional) distributions; Dirichlet processes will be used for this purpose.

Let  $T_p(x, x^*, y, y^*)$  be the object of interest, i.e. either  $q_v(y^* | x^*, x, y)$  or  $q_{ev}(y^* | x^*, x, y)$ . Suppose that we have a prior on the parameter space that  $p^{-1}$  belongs to: we will say more about the parameter space later.<sup>37</sup> Suppose further that  $\mathbf{G}$  is a draw from the prior. Then, our prediction based on the draw  $\mathbf{G}$  will be

$$T_{\mathbf{G}^{-1}}(x, x^*, y, y^*), \quad (24)$$

which is a random object. The approach considered in this section is based on the probability distribution of (24). We show in section 6.2 that the mean of (24) is the average risk optimal prediction when a quadratic loss function is used.

We emphasize that there is no Bayesian updating here, because the only information available about the parameter  $p$  in the distribution of observables is the regression function  $\mu_m$ , which we assume to be known (throughout this paper).

We now obtain the likelihood function (in terms of the parameter  $p$ ). Note from (6) that for  $e \in S_m(x)$ ,

$$\mathbb{P}\{\mathbf{y} = (1, 0) | e = e, \mathbf{v} = v, \mathbf{x} = x; p\} = p(e, v, x),$$

<sup>37</sup> $p^{-1}$  is the (generalized) inverse of the quantile function  $p$  with respect to the argument  $v$ . For a quantile function  $Q$ , its generalized inverse  $G^{-1}$  is defined by  $Q^{-1}(p) = \sup\{v : Q(v) \leq p\}$ , which is a distribution function. Similarly, for a distribution function  $G$ , the generalized inverse  $G^{-1}$  is defined by  $G^{-1}(v) = \inf\{x : G(x) \geq v\}$ , which is a quantile function.

which implies

$$\frac{1}{\pi_m(x)} \int_{S_m(x)} \mathbb{P}\{\mathbf{y} = (1, 0) \mid \mathbf{e} = e, \mathbf{x} = x; \rho\} f_e(e) \, de = \mu_m(x). \quad (25)$$

Equation (25) shows that the likelihood function depends on  $\rho$  only through  $\mu_m$ , which is known. Therefore, since  $\mu_m$  is known and fixed, the likelihood is also known and it does not depend on  $\rho$  other than via  $\mu_m$ , which is a different way of saying (as we have before) that  $\mu_m$  is the only identifiable object about behavior in the multiplicity region.

**6.2 Average risk optimality:** We now consider an alternative approach, where we formulate the problem as a decision problem of making the ‘best’ guess for  $T_\rho(x, x^*, y, y^*)$  as a function of  $x, x^*, y, y^*$ . To be consistent with our previous discussion, we study the decision problem at the ‘population level,’ i.e. in terms of the ‘population’ distribution: there is no sample in the problem considered.

Fixing  $x, x^*, y, y^*$ , denote our decision by  $d(x, x^*, y, y^*)$ . For a given loss function  $\ell$ , the loss (risk) of the decision  $d$  is

$$R(d, \rho) = \ell\{d(x, x^*, y, y^*), T_\rho(x, x^*, y, y^*)\}. \quad (26)$$

Like in section 5 our analysis is at the population level as: the loss is not random and therefore the loss and the risk are the same thing.<sup>38</sup> However, since the prediction of interest depends on  $\rho$ , the risk in (26) depends on  $\rho$ , also. Therefore, we consider average risk, i.e. risk averaged over  $\rho$ , for which we need a prior for  $\rho$ .<sup>39</sup>

Now consider the parameter space that the function  $\rho^{-1}$  belongs to. For given values of  $x, x^*$ , the function  $\rho^{-1}$  belongs to  $\mathcal{F}_{xx^*} = \mathcal{F}_x \cup \mathcal{F}_{x^*}$ , where

$$\mathcal{F}_x = \left\{ \rho^{-1}(\cdot, \cdot, x) : \frac{1}{\pi_m(x)} \int_{S_m(x)} \int_0^1 \rho(e, v, x) f_e(e) \, de \, dv = \mu_m(x) \right\}. \quad (27)$$

Therefore, once we have priors  $\omega_x, \omega_{x^*}$  on  $\mathcal{F}_x, \mathcal{F}_{x^*}$ , we can induce a prior  $\omega_{xx^*}$  on  $\mathcal{F}_{xx^*}$ . Since  $\rho^{-1}$  is a (conditional) distribution function, Dirichlet processes provide a natural way of choosing  $\omega_x$ : Dirichlet processes are the infinite-dimensional extension of the Dirichlet distribution, which is a conjugate prior for the multinomial distribution. We propose an algorithm to do this in section 6.3.

Once we have  $\omega_{xx^*}$ , the average risk is given by

$$\bar{R}(d) = \int_{\mathcal{F}_{xx^*}} R(d, \rho) \, d\omega_{xx^*}(\rho^{-1}) \geq \min_s \int_{\mathcal{F}_{xx^*}} \ell\{s, \rho(x, x^*, y, y^*)\} \, d\omega_{xx^*}(\rho^{-1}),$$

where the inequality shows that the function  $d^*$  that minimizes  $\bar{R}$  can be found by separate minimization for each combination of  $x, x^*, y, y^*$ . Indeed, the decision  $d^*$  that minimizes

<sup>38</sup>We could average  $y, y^*, x, x^*$  by using the likelihood function, which is independent of  $\rho$  (other than via  $\mu_m$ ) and hence is known. But it does not make any difference from simply fixing them as we do here.

<sup>39</sup>As we discussed earlier, a simple (unconstrained) minmax approach leads to the midpoint prediction if  $\ell$  is symmetric. But the midpoint prediction has flaws as will become apparent in section 7.1.

$\bar{R}$  is

$$d^*(x, x^*, y, y^*) = \operatorname{argmin}_{s \in [0,1]} \int_{\mathcal{F}_{xx^*}} \ell\{s, T_\rho(x, x^*, y, y^*)\} d\omega_{xx^*}(\rho^{-1}).$$

Thus, when the loss function is quadratic,  $d^*(x, x^*, y, y^*)$  is simply a mean as we explained below (24).

**6.3 Dirichlet prior:** We now describe a way of drawing probability distributions from  $\omega_{xx^*}$ , for which it suffices that we can draw probability distributions from  $\omega_x$  on support  $\mathcal{F}_x$ . To keep things simple, we will restrict  $\mathcal{F}_x$  by considering elements  $\rho^{-1}$  that do not depend on  $e$  for all  $e \in S_m(x)$ . With maximum entropy, the fact that the  $\rho$  used is flat in  $e$  is a *result* of the maximum entropy procedure, but here it is a *restriction*. Dropping the argument  $e$  from our notation, we have

$$\mathcal{F}_x = \left\{ \rho^{-1}(\cdot, x) : \int_0^1 \rho(v, x) dv = \mu_m(x) \right\}, \quad (28)$$

which is a collection of distribution functions constrained to have mean  $\mu_m(x)$ .

A Dirichlet process  $\mathcal{D}$  is commonly used as a probability distribution for probability distributions. It is a stochastic process whose finite marginals are described by the Dirichlet distribution. It is characterized by a (Dirichlet process) prior  $H$ , not to be confused with  $\omega_x$ , and a hyperparameter  $\zeta > 0$ .  $H$  and  $\zeta$  correspond to the mean and the precision of the Dirichlet process, respectively, in the following sense. If  $\mathbf{G} \sim \mathcal{D}(\zeta, H)$ , then for an arbitrary event  $A$  we have

$$\mathbb{E}\mathbf{G}(A) = H(A), \quad \mathbb{V}\mathbf{G}(A) = H(A)\{1 - H(A)\} / (1 + \zeta).$$

In the discussion below we will use the uniform distribution  $U[0, 1]$  for  $H$  and let  $\zeta = 1$ .

A probability distribution  $\mathbf{G}$  drawn from a Dirichlet process is always a discrete distribution with infinitely many mass points. A standard procedure for drawing  $\mathbf{G}$  from e.g.  $\mathcal{D}(1, U[0, 1])$  is the following. Let  $\beta_1, \tau_1, \beta_2, \tau_2, \dots$  be independent draws from  $U[0, 1]$ , let  $\beta_1^* = \beta_1$  and  $\beta_j^* = \beta_j \prod_{t=1}^{j-1} (1 - \beta_t)$  for  $j > 1$ . Then, the distribution  $\mathbf{G}$  assigning probability masses  $\beta_1^*, \beta_2^*, \dots$  to the mass points  $\tau_1, \tau_2, \dots$ , respectively, can be shown to be a draw from  $\mathcal{D}(1, U[0, 1])$ . This procedure is known as the ‘stick-breaking’ construction. See e.g. Teh (2010) for details.

In our case, however, the standard stick-breaking construction does not provide what we want because the parameter space  $\mathcal{F}_x$  has a constraint on the mean. Below, we explain how we modify the standard procedure to draw a probability distribution from  $\mathcal{F}_x$ . Since  $\mathcal{F}_x$  depends on  $x$  only through  $\mu_m(x)$ , we will drop the dependence on  $x$  in our notation and simply use  $\mu_m$  in lieu of  $\mu_m(x)$  for the remainder of this section.

Let the  $\beta_j$ ’s and  $\beta_j^*$ ’s be defined as above. Suppose that the  $\tau_j$ ’s are i.i.d.  $U[0, 1]$  and (unconditionally) independent of the  $\beta_j$ ’s. So, if we would not impose a mean condition then our procedure below would be equivalent to standard stick-breaking. We now determine the locations of the mass points  $\tau_j$ ’s to ensure that

$$W_1 = \sum_{j=1}^{\infty} \beta_j^* \tau_j = \mu_m,$$



for which we will draw  $\tau_j$ 's from appropriate conditional distributions.

First, note that the distribution of

$$W_j = \frac{\beta_j}{\beta_j^*} \sum_{i=j}^{\infty} \beta_i^* \tau_i,$$

does not depend on  $j$  and that

$$\forall j : W_j = \beta_j \tau_j + (1 - \beta_j) W_{j+1}. \quad (29)$$

Further, the Bayes rule says that

$$f_{\tau_j | \beta_j, W_j}(\tau | \beta, \mu) = \frac{f_{W_j | \tau_j, \beta_j}(\mu | \tau, \beta)}{f_{W_j | \beta_j}(\mu | \beta)}. \quad (30)$$

Here, by (29) and the fact that the distribution of  $W_j$  does not depend on  $j$ , it follows that

$$\mathbb{P}(W_j \leq w | \tau_j = \tau_j, \beta_j = \beta_j) = \mathbb{P}\{\beta_j \tau_j + (1 - \beta_j) \tau_j W_{j+1} \leq w\} = F^\circ\left(\frac{w - \beta_j \tau_j}{1 - \beta_j}\right), \quad (31)$$

where  $F^\circ$  is the distribution function of  $W_1$ . We discuss further down how to compute  $F^\circ$  efficiently.

Finally, note from (30) and (31) that  $F_{\tau_j | \beta_j, W_j}(\tau | \beta, \mu) = \tilde{F}^\circ(\tau; \beta, \mu)$ , where

$$\tilde{F}^\circ(t; \tilde{\beta}, \tilde{\mu}) = \frac{F^\circ\left(\frac{\tilde{\mu}}{1 - \tilde{\beta}}\right) - F^\circ\left(\frac{\tilde{\mu} - \tilde{\beta}t}{1 - \tilde{\beta}}\right)}{F^\circ\left(\frac{\tilde{\mu}}{1 - \tilde{\beta}}\right) - F^\circ\left(\frac{\tilde{\mu} - \tilde{\beta}}{1 - \tilde{\beta}}\right)}, \quad \max\left\{0, \frac{\tilde{\mu} - (1 - \tilde{\beta})}{\tilde{\beta}}\right\} \leq t \leq \min\left(1, \frac{\tilde{\mu}}{\tilde{\beta}}\right). \quad (32)$$

This motivates the following procedure.

**Procedure 1.** Do the following:

1. Draw  $\beta_1, \beta_2, \dots \sim U(0, 1)$ ;
2. Draw  $\tau_1$  from  $\tilde{F}^\circ(\cdot; \beta_1, \mu_m)$ ;
3. Draw  $\tau_2$  from  $\tilde{F}^\circ\{\cdot; \beta_2, (\mu_m - \beta_1 \tau_1) / (1 - \beta_1)\}$ ;
4. Draw  $\tau_3$  from  $\tilde{F}^\circ[\cdot; \beta_3, \{\mu_m - \beta_1 \tau_1 - (1 - \beta_1)\beta_2 \tau_2\} / \{(1 - \beta_1)(1 - \beta_2)\}]$ ;
5. Continue ad nauseam.
6. Let  $G$  be the probability distribution with mass points  $\beta_1^*, \beta_2^*, \beta_3^*, \dots$  at  $\tau_1, \tau_2, \tau_3, \dots$ , where  $\beta_j^*$  is as defined in the text.

□

Then,  $\mathbf{G}$  is a draw from  $\mathcal{D}(1, U[0, 1])$  unconditionally. But once we fix  $\mathbf{G}$ 's mean at  $\mu_m$ , its distribution is different for which we use the notation  $\mathbf{G} \sim \mathcal{D}_{\mu_m}^* = \mathcal{D}_{\mu_m}^*(1, U[0, 1])$ .

From (32) it follows that for any  $j > 1$ ,

$$\mu_m - \prod_{t=1}^{j-1} (1 - \beta_t) \leq \sum_{t=1}^j \beta_t^* \tau_t \leq \mu_m,$$

which ensures that  $\mathbf{G}$  has mean  $\mu_m$ , as required: taking  $j \rightarrow \infty$  leads to  $W_1 = \mu_m$ .

Our procedure requires us to implement a draw from  $\tilde{F}^\circ(\cdot; \tilde{\beta}, \tilde{\mu})$ . All one has to do for this purpose is to compute

$$\tau = \frac{\tilde{\mu} - (1 - \tilde{\beta})F^{\circ-1}\left\{(1 - \tau^*)F^\circ\left(\frac{\tilde{\mu}}{1 - \tilde{\beta}}\right) + \tau^*F^\circ\left(\frac{\tilde{\mu} - \tilde{\beta}}{1 - \tilde{\beta}}\right)\right\}}{\tilde{\beta}},$$

where  $\tau^*$  is a draw from a  $U(0, 1)$ .

It now remains to be shown how to compute  $F^\circ$ . From (29) it follows that

$$\forall w : F^\circ(w) = \int_0^1 \int_0^1 F^\circ\left(\frac{w - \beta\tau}{1 - \beta}\right) d\beta d\tau,$$

which allows us to solve for  $F^\circ$  numerically.

## 7. Comparison of methods

In this section we compare the maximum entropy approach with other possibilities. Section 7.1 shows that the midpoint prediction method is inconsistent in the sense that there generally is no single function  $\rho$  that produces the midpoint predictions in all cases. Section 7.2 uses two simple designs to compare the sharp bounds on  $q_v$  to regression predictions, midpoint predictions, and maximum entropy predictions. Section 7.3 does the same for  $q_{ev}$ . Finally, section 7.4 compares and contrasts all methods to the Dirichlet–process–based idea of section 6.

**7.1 Inconsistency of the midpoint method:** As promised, we now show that there generally exists no single function  $\rho$  that is consistent with the midpoint prediction for all values of  $x^*$ . In other words, we can construct examples in which any function  $\rho$  that rationalizes the midpoint prediction for one combination of  $x, x^*$  does not rationalize the midpoint prediction for a different combination, as evidenced by example 1.

**Example 1.** Suppose that  $x \neq x^*$ ,  $y = y^* = (1, 0)$ ,  $\pi_m(x) = \pi_m(x^*) = 1$ , and  $\mu_m(x) = \mu_m(x^*) = \tilde{\mu}$  for some  $0 < \tilde{\mu} \leq 1$ . Suppose that there exists a function  $\rho^\circ$  such that the predictions based on  $\rho^\circ$  are the midpoint predictions. Then, for any  $\bar{x}$  in the support of  $\mathbf{x}$ , the predictions based on  $\rho^\circ$  must satisfy

$$q_v^\circ(y^* | \bar{x}, \bar{x}, y) = \frac{\mathbb{E}\{\rho^\circ(\mathbf{e}, \mathbf{v}, \bar{x})\rho^\circ(\mathbf{e}^*, \mathbf{v}, \bar{x})\}}{\mu_m(\bar{x})} = q_{ev}^\circ(y^* | \bar{x}, \bar{x}, y) = \frac{\mathbb{E}\rho^{\circ 2}(\mathbf{e}, \mathbf{v}, \bar{x})}{\mu_m(\bar{x})},$$

because the midpoint predictions of  $q_v(y^* | \bar{x}, \bar{x}, y)$  and  $q_{ev}(y^* | \bar{x}, \bar{x}, y)$  coincide by theorems 2 and 5. Therefore, we have  $\mathbb{E} \nabla \{p^\circ(e, v, x) | v\} = 0$  a.s., from which we conclude that  $p^\circ$  is flat in  $e$ . From hereon we omit the argument  $e$  without loss of generality, i.e.  $p^\circ(e, \cdot, \cdot) = p^\circ(\cdot, \cdot)$ .

Now, suppose that  $\tilde{\mu} = 1/2$ , in which case the midpoint prediction of  $q_{ev}(y^* | x^*, x, y)$  equals  $1/2$ , and the midpoint predictions of  $q_{ev}(y^* | x, x, y)$  and  $q_{ev}(y^* | x^*, x^*, y)$  are both equal to  $3/4$  by theorem 5. Therefore, if  $p^\circ$  is to be consistent with the midpoint predictions we must have

$$q_{ev}^\circ(y^* | x^*, x, y) = \frac{\mathbb{E}\{p^\circ(v, x)p^\circ(v, x^*)\}}{\tilde{\mu}} = \frac{1}{2}, \quad (33)$$

$$q_{ev}^\circ(y^* | x, x, y) = \frac{\mathbb{E}p^{\circ 2}(v, x)}{\tilde{\mu}} = \frac{3}{4}, \quad (34)$$

$$q_{ev}^\circ(y^* | x^*, x^*, y) = \frac{\mathbb{E}p^{\circ 2}(v, x^*)}{\tilde{\mu}} = \frac{3}{4}. \quad (35)$$

Equation (33) implies that  $\text{Cov}\{p^\circ(v, x), p^\circ(v, x^*)\} = 0$ , which means that either  $p^\circ(v, x)$  or  $p^\circ(v, x^*)$  should be equal to  $\tilde{\mu} = 1/2$  with probability one, because  $p^\circ$  is monotonic in  $v$ . But that conclusion contradicts either (34) or (35). Therefore, there exists no function  $p^\circ$  that rationalizes the midpoint predictions.  $\square$

There are further arguments against the use of midpoint predictions, as will become apparent in the remainder of section 7.

**7.2 Case 1:  $v^* = v$ :** We now turn to a comparison of the prediction methods for the case discussed in section 3.1, i.e. if  $q_v(y^* | x^*, x, y)$  is the object of interest. Both here and in section 7.3 we focus on the case in which  $\mu_m(x^*) = \mu_m(x)$ , which simplifies the comparison while still allowing us to convey the main issues at hand.

Recall from theorem 7 that the maximum entropy method produces a solution  $p^*$  which (in  $S_m(x)$ ) depends on  $x$  only through  $\mu_m(x)$ . Therefore,  $\mu_m(x) = \mu_m(x^*)$  implies that  $p^*(e, v, x) = p^*(e, v, x^*)$  for all  $e \in S_m(x)$ ,  $v$ , and hence it follows that

$$\begin{aligned} \rho^*(x, x^*) &= \nabla \{k^*(v, x) | x = x\} \\ &= \int_0^1 p^2 A\{p, \lambda_m(x)\} dp - \left\{ \int_0^1 p A\{p, \lambda_m(x)\} dp \right\}^2 = \mathcal{L}''\{\lambda_m(x)\}, \end{aligned} \quad (36)$$

where  $A$  is defined in (20). Example 2 is based on (36).

**Example 2.** Suppose that  $\pi_{10}(x) = \pi_{10}(x^*) = \pi_{01}(x) = \pi_{01}(x^*)$ ,  $\pi_m(x) = \pi_m(x^*)$ ,  $\pi_{00}(x) = \pi_{00}(x^*)$ , and  $\mu_m(x) = \mu_m(x^*) = \tilde{\mu}$  for some  $0 < \tilde{\mu} \leq 1$ . We consider two cases, namely  $\pi_{10}(x) = 1/4$ ,  $\pi_m(x) = 4/9$  and  $\pi_{10}(x) = 0$ ,  $\pi_m(x) = 1$ . Figure 12 depicts prediction probabilities as a function of  $\tilde{\mu}$  in each case.

In both cases, the midpoint prediction method yields higher predictions than the maximum entropy method. In the extreme case  $\pi_m(x) = 1$  (depicted in the right panel), the difference is especially pronounced when  $\tilde{\mu}$  is close to zero. The upper bound there

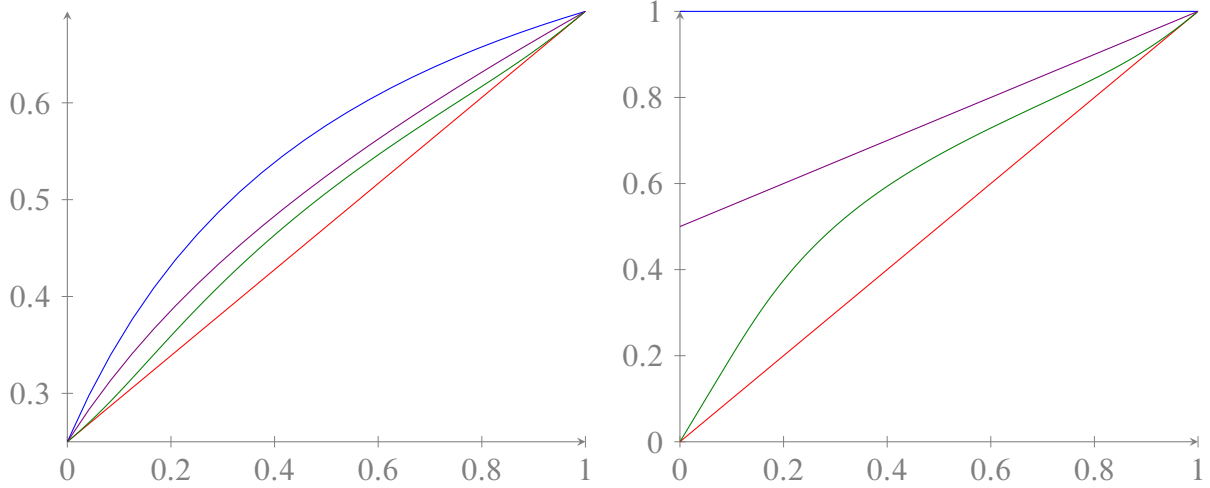


Figure 12: Infeasible predictions as a function of  $\mu_m(x) = \mu_m(x^*)$  if  $y = y^* = (1, 0)$  using regression (red; also smallest), largest (blue), midpoint (purple), and maximum entropy (green) methods. In both panels  $\pi_m(x) = \pi_m(x^*)$  and  $\pi_y(x) = \pi_y(x^*)$ . In the left panel,  $\pi_m(x) = 4/9$  and  $\pi_y(x) = 1/4$ , whereas in the right panel  $\pi_m(x) = 1$  and  $\pi_y(x) = 0$ .

equals one for all  $\tilde{\mu} \in (0, 1]$  since the largest possible value of  $\rho(x, x^*)$  is  $\tilde{\mu}(1 - \tilde{\mu})$ , which corresponds to

$$\rho(e, v, x) = \mathbb{1}(v > 1 - \tilde{\mu}), \quad e \in S_m(x), \quad (37)$$

From theorem 1, we know that (in the right panel case)

$$q_v\{(1, 0) \mid x^*, x, (1, 0)\} = \frac{1}{\tilde{\mu}} \int_0^1 \hat{k}(v, x) \hat{k}(v, x^*) dv. \quad (38)$$

For most functions  $\rho$ , (38) will be close to zero if  $\tilde{\mu}$  is close to zero. However, for  $\rho$  in (37), the value of  $q_v$  is one for any value of  $\tilde{\mu}$ , which is an extreme possibility. In a way, then, the midpoint method is overly conservative since it puts an inordinate amount of weight on an extreme choice of  $\rho$  in (37); see also example 5. The maximum entropy method, in contrast, does not have this problem.  $\square$

**7.3 Case 2:  $e^* = e$  and  $v^* = v$ :** Since  $p^*(e, v, x)$  defined in theorem 7 does not vary with  $e$  over  $S_m(x)$ , the function  $\rho$  in theorem 3 using  $p^*$  is given by

$$\rho^*(x, x^*) = \mathbb{E}[\text{Cov}\{p^*(e, v, x)p^*(e, v, x^*) \mid e\} \mid e \in S_{mm}(x, x^*)] = \text{Cov}\{\hat{k}^*(v, x), \hat{k}^*(v, x^*)\}.$$

We continue to focus on the case in which  $\mu_m(x) = \mu_m(x^*)$  and hence  $\hat{k}^*(v, x) = \hat{k}^*(v, x^*)$  for all  $v$ , as explained in section 7.2. Consequently, the formula for  $\rho^*(x, x^*)$  is again given by (36), i.e.  $\rho^*(x, x^*) = \mathcal{L}''\{\lambda_m(x^*)\}$ . The function  $a_y$  in theorem 3 using  $p^*$  is provided in (54) in appendix C, which can be used to determine the maximum entropy value of  $\alpha_{yy}(x, x^*)$ . Indeed, if the  $S$  regions at  $x$  coincide with those at  $x^*$  and  $\mu_m(x) = \mu_m(x^*)$ , then the maximum entropy choice of  $\alpha_{yy^*}(x, x^*)$  is for  $y^* = y = (1, 0)$  given by

$$[\pi_y(x)\{1 - \pi_y(x)\} + \mu_m^2(x)\pi_m(x)\{1 - \pi_m(x)\} - 2\pi_y(x)\pi_m(x)\mu_m(x)] / \tilde{\pi}_y(x). \quad (39)$$

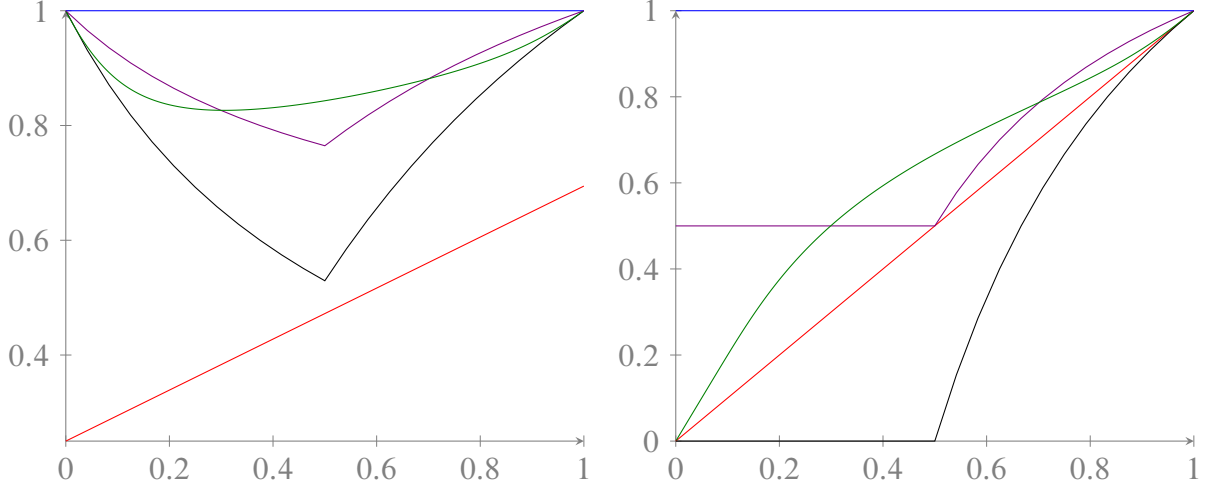


Figure 13: Infeasible predictions as a function of  $\mu_m(x) = \mu_m(x^*)$  if  $x$  and  $x^*$  correspond to the same  $S$ -regions with  $x \neq x^*$ , and  $y = y^* = (1, 0)$ . The predictions are produced using the regression (red), smallest (black), midpoint (purple), largest (blue), and maximum entropy (green) methods. In the left panel,  $\pi_m(x) = 4 / 9$  and  $\pi_y(x) = 1 / 4$ , whereas in the right panel  $\pi_m(x) = 1$  and  $\pi_y(x) = 0$ .

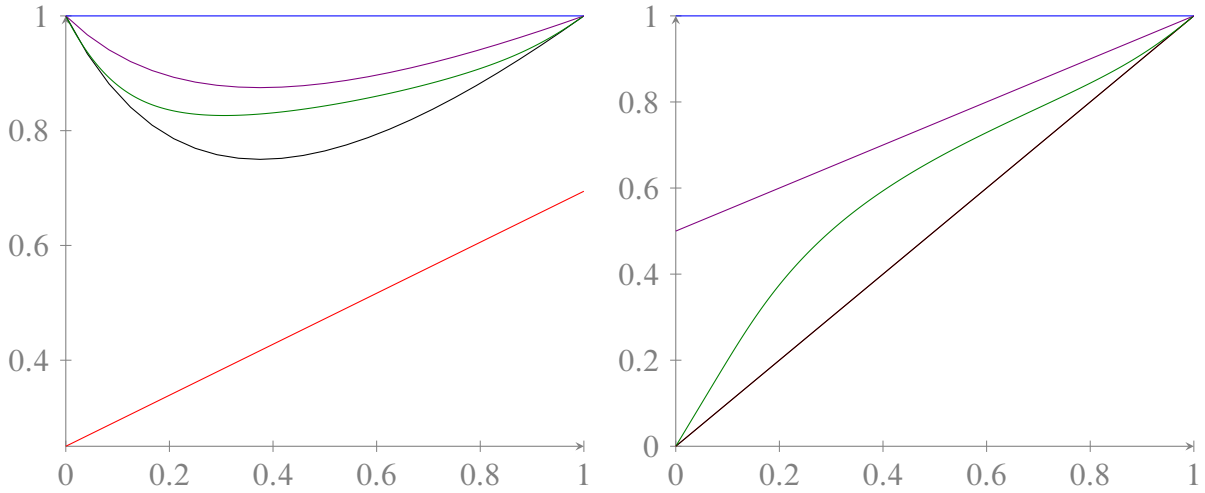


Figure 14: Infeasible predictions as a function of  $\mu_m(x)$  if  $x = x^*$  and  $y = y^* = (1, 0)$ , using regression (red), smallest (black), midpoint (purple), largest (blue), and maximum entropy (green) methods. In the left panel,  $\pi_m(x) = 4 / 9$  and  $\pi_y(x) = 1 / 4$ , whereas in the right panel  $\pi_m(x) = 1$  and  $\pi_y(x) = 0$ .

Example 3 below is based on (36) and (39).

**Example 3.** Consider  $x, x^*$  such that the  $S$ -regions coincide and  $\mu_m(x) = \mu_m(x^*) = \tilde{\mu}$  for some  $0 < \tilde{\mu} \leq 1$ . We consider four cases: each of the two cases considered in example 2 where we now distinguish between  $x \neq x^*$  and  $x = x^*$ . Holding  $e, v, x$  fixed does not necessarily yield the same outcome because of the incompleteness of the model: a different equilibrium can be selected in the multiplicity region. Figures 13 and 14 depict the cases

$x^* \neq x$  and  $x^* = x$ , respectively, where in both figures  $y^* = y = (1, 0)$ . Recall that the regression prediction  $\tilde{\pi}_{10}(x^*)$  need not belong to the identified set for  $q_{ev}$ . Therefore, the regression prediction is a poor choice if the object of interest is  $q_{ev}(y^* | x, x^*, y)$ .

First consider the right panel in figure 13. The midpoint method yields 0.5 for all  $\tilde{\mu} < 0.5$ . So the ‘conservative prediction’ problem mentioned in example 2 arises here, also. Further, the maximum entropy predictions behave more smoothly than e.g. the midpoint predictions: the midpoint predictions have a noticeable kink at  $\tilde{\mu} = 0.5$ .

Comparing figures 13 and 14 shows that the midpoint method yields drastically different predictions depending on whether or not  $x^* = x$ . For instance, in the extreme case  $\pi_m(x) = \pi_m(x^*) = 1$ , if  $\tilde{\mu} = 1/2$  then theorem 5 implies that the sharp identified interval is given by

$$[\mathbb{1}(x^* = x) / 2, 1].$$

So a minute change in  $x^*$  can result in a 25 percentage point jump in the midpoint prediction, which is undesirable.  $\square$

Intuition for the discontinuity problem illustrated in example 3 was provided below theorem 5. Recall that the sharp bounds change *continuously* if  $p(e, \mathbf{v}, x)$  and  $p(e, \mathbf{v}, x^*)$  are restricted to have nonnegative correlation. For instance, if the function  $p$  is flat in  $e$ , then the discontinuity issue disappears.

Since all available information about the function  $p$  is contained in  $\mu_m(x)$  and  $S_m(x)$  (and hence in  $\pi_m(x)$ ), restricting  $p$  to be flat in  $e$  does not lead to a loss of ‘information’ in the maximum entropy sense: the restriction that  $p$  is flat in  $e$  is not binding in the maximum entropy optimization problem and hence the maximum entropy solution is flat in  $e$ . With the Dirichlet-based procedure of section 6, we restricted  $p$  to be flat in  $e$ . However, since the parameter space  $\mathcal{F}_x$  in (27) only depends on  $x$  via  $\mu_m(x)$  and  $S_m(x)$ , any prior that depends on  $x$  only through the identifiable objects  $\mu_m(x)$  and  $S_m(x)$  will not impose ‘extra’ information about how  $p$  depends on  $e$ , and therefore will produce predictions that are continuous in  $x$ .

**7.4 Dirichlet-based predictions, midpoints, and maximum entropy:** We now consider the average risk optimal predictions using the prior we proposed in section 6. Instead of choosing a particular loss function, we consider the distribution of the prediction when  $p$  is drawn by the algorithm described in section 6.3. For instance, the mean and median are the average risk optimal predictions using a quadratic and an absolute deviation loss function, respectively.

**Example 4.** Consider again the scenario of examples 2 and 3, specifically the cases depicted in the right panels of figures 12 and 13, i.e.  $\pi_m(x) = \pi_m(x^*) = 1$  and  $y = y^* = (1, 0)$ . Recall that by theorem 1 the infeasible prediction  $q_v(y^* | x^*, x, y)$  in this case is given by

$$\frac{\mathbb{E}\{k(\mathbf{v}, x)k(\mathbf{v}, x^*)\}}{\mu_m(x)}. \quad (40)$$

Likewise, the infeasible prediction for  $q_{ev}(y^* | y, x, x^*)$  is given by

$$\frac{\mathbb{E}\{p(e, \mathbf{v}, x)p(e, \mathbf{v}, x^*)\}}{\mu_m(x)}. \quad (41)$$

Now, consider the restricted parameter space defined in (28). Then, in our current design, (40) and (41) coincide and equal

$$T_p(x, x^*, y, y^*) = \frac{\mathbb{E}\{p(\mathbf{v}, x)p(\mathbf{v}, x^*)\}}{\mu_m(x)}. \quad (42)$$

Note that the  $S$ -regions do not change in the current setup. If the parameter space for  $p$  is not restricted to  $\mathcal{F}_x$  then (42) is an upper bound of both (40) and (41).

Figure 15 depicts the distribution of (42) when the function  $p$  is a random draw from the Dirichlet-like-process  $\mathcal{D}_{\mu_m}^*$  described in section 6.3. Recall from section 6 that this entails pretending that each probability distribution is a possible distribution of  $\mathbf{p}$  for given  $\mu_m$ -value with the caveat that the draws from  $\mathcal{D}_{\mu_m}^*$  are discrete distributions whereas the distribution of  $\mathbf{p}$  is in most cases continuous.

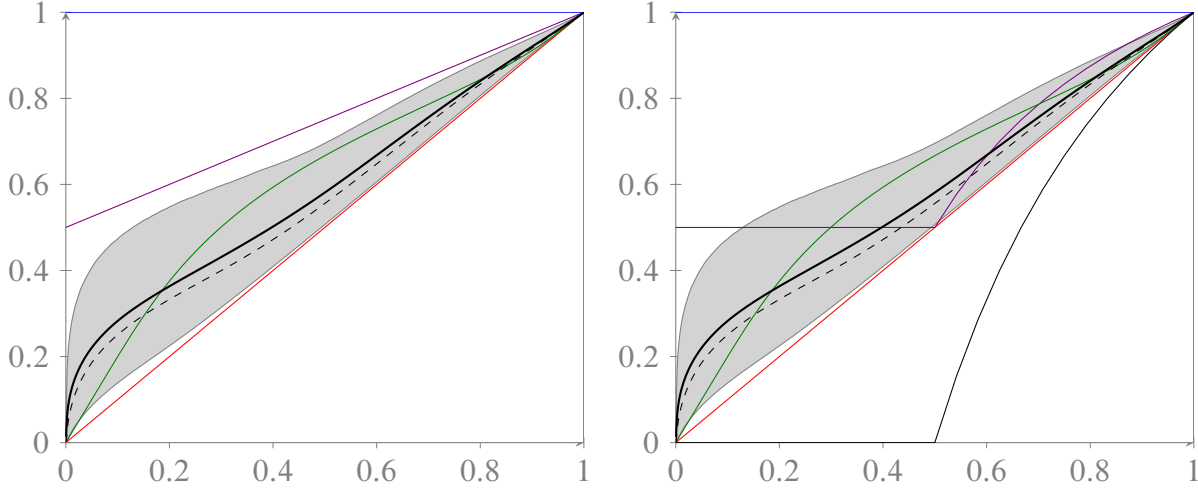


Figure 15: Infeasible predictions as a function of  $\mu_m(x) = \mu_m(x^*)$  if  $y = y^* = (1, 0)$  using regression (red), largest (blue), midpoint (purple), and maximum entropy (green) methods. The grey regions depict the area between 10% and 90% quantiles of the Dirichlet process-based experiment described in section 6.3 as a function of its mean with the dashed line indicating the median. In both panels  $\pi_m(x) = \pi_m(x^*) = 1$  and  $x \neq x^*$ . The left panel depicts  $q_v$  and the right panel depicts  $q_{ev}$ . The maximum entropy prediction and the Dirichlet-based predictions are the same in both panels while the bounds (and hence the midpoints) are different.

For each value of  $\mu_m$ , each draw from  $\mathcal{D}_{\mu_m}^*$  produces a probability distribution for  $\mathbf{p}$ , whose quantile function corresponds to (a draw of) the function  $p$ . Therefore, each draw from  $\mathcal{D}_{\mu_m}^*$  produces the prediction  $T_p(x, x^*, y, y^*)$  displayed in (42) for a different function  $p$ .

The results of our experiments are depicted in figure 15. The left panel in figure 15 is identical to the right panel in figure 12 and the right panel in figure 15 is identical to the right panel in figure 13, except that there are some additions. The additions in the two panels are identical. The thick solid black line represents the mean of the distribution of

$T_p(x, x^*, y, y^*)$  as a function of  $\mu_m$ , the dashed line the median, and the bottom and top of the gray shaded area the bottom and top decile, respectively.

The Dirichlet graphs are consistent with the maximum entropy predictions (green curves). But the Dirichlet graphs are inconsistent with the midpoint predictions. As we explained below (42), the Dirichlet graphs are generally upper bounds to the infeasible predictions, with equality if  $p$  is flat in  $e$ . They moreover depend on  $x$  only through  $\mu_m(x)$ . The fact that the midpoint prediction in the left panel of figure 15 is for essentially all  $\mu_m$  above the 90% quantile of the distribution of  $T_p(x, x^*, y, y^*)$  produced by the Dirichlet experiment is troubling. The results in the right panel are less problematic, but the midpoint predictions nevertheless look at odds with the Dirichlet draws.

The Dirichlet experiment has limitations. First, Dirichlet process draws are, as noted, discrete probability distributions whereas  $\mathbf{p}$  is usually (though not necessarily) continuously distributed. Further, we chose a uniform prior and pseudocount hyperparameter equal to one for the Dirichlet process for convenience.<sup>40</sup> Different choices for the input parameters produce different results. Finally, we only draw distributions for  $\mathbf{p}$  (and hence functions  $p$ ) that do not depend on the value of  $e \in S_m(x)$  and only depend on  $x$  via  $\mu_m(x)$ . But that actually helps the midpoint prediction because (42) is only an upper bound to the infeasible prediction: the infeasible prediction can in fact be lower still.  $\square$

While acknowledging its limitations, example 4 provides further evidence against the use of the midpoint method. We conclude with an example that emphasizes the fact that the midpoint prediction puts too much weight on extreme possibilities.

**Example 5.** Consider the case in which  $x \neq x^*$ ,  $y = y^* = (1, 0)$ ,  $\pi_m(x) = \pi_m(x^*) = 1$ , but where  $\mu_m(x^*)$  may or may not equal  $\mu_m(x)$ . Then, by theorem 1,

$$q_v(y^* | x^*, x, y) = \mu_m(x^*) + \frac{\text{Cov}\{k(v, x), k(v, x^*)\}}{\mu_m(x)} = \frac{\mathbb{E}\{k(v, x)k(v, x^*)\}}{\mu_m(x)},$$

where  $k(v, x) = \mathbb{E}p(e, v, x)$ , as before. The highest attainable value of  $q_v(y^* | x^*, x, y)$  is one. Here is what needs to happen to get  $q_v(y^* | x^*, x, y)$  equal to one. Since  $k$  cannot exceed one, it must be true that  $\mu_m(x^*) \geq \mu_m(x)$  and that  $k(v, x^*) = 1$  whenever  $k(v, x) > 0$ . Since  $k$  is nondecreasing in  $v$ , it must be true that  $k(v, x) = \mathbb{1}\{v \geq 1 - S_m(x)\}$  and that  $k(v, x^*) = k(v, x)$  for all  $v \geq 1 - S_m(x)$ .

The midpoint solution assigns weight 0.5 to this case and weight 0.5 to the case in which  $k(v, x) = \mu_m(x)$  for all values of  $v$ .  $\square$

## 8. Inference

So far we have focused on counterfactuals in games, for which we assumed that both the structure and probability distributions were known. However, in practice payoffs need

<sup>40</sup>The prior determines how likely it is that the mass points of a probability distribution drawn are in particular locations. A uniform prior means that those mass points can be anywhere in the unit interval with equal probability: the distributions generated by the Dirichlet process are themselves not uniform. Moreover, recall that we generate distributions conditional on the distributions having mean  $\mu_m$ , not unconditionally. The pseudocount hyperparameter, determines the relative size of mass points: the value chosen by us (one) implies that the probability mass at the first mass point is on average 1/2, at the second 1/4, etcetera.



to be estimated which produces estimation error. In this section we discuss inferential issues for maximum entropy (ME) counterfactual predictions, focusing on  $q_v$ . Like before, payoffs are not the main focus of this paper, so we make the following assumption.

**Assumption A.** For some  $0 < r \leq 1/2$  and some  $\mathbf{X}$ ,  $n^r(\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}) \xrightarrow{d} \mathbf{X}$ , where  $n$  is the sample size,  $\boldsymbol{\mu} = [\mu_m(x), \mu_m(x^*)]^\top$ , and  $\hat{\boldsymbol{\mu}}$  is an estimator of  $\boldsymbol{\mu}$ .  $\square$

As noted in the introduction, finding  $\hat{\boldsymbol{\mu}}$  and  $\mathbf{X}$  has been discussed by many authors with their exact form depending on the assumptions made. In Tamer (2003),  $r = 1/2$  and  $\mathbf{X}$  is normally distributed.

In the interest of space conservation we assume that  $x \neq x^*$ : the case  $x = x^*$  is analogous. Further, we focus on doing inference on the nontrivial portion of the counterfactual prediction, i.e. the object  $\rho(x, x^*)$  (or more precisely, the one that corresponds to the ME solution of  $\rho$ ), which only depends on  $\boldsymbol{\mu}$ . We now translate knowledge of the large sample behavior of  $\hat{\boldsymbol{\mu}}$  into a result on the limit properties on  $\hat{\boldsymbol{\rho}} = \hat{\boldsymbol{\rho}}(x, x^*)$ , an estimator of  $\rho = \rho(x, x^*)$  which is defined in (10) and used in theorem 1; theorems 6 and 7 provide the ME solution to be used in theorem 1.

Indeed, theorem 7 tells us the explicit functional form of  $k^*(\cdot, x)$  and  $k^*(\cdot, x^*)$ : they depend on  $\boldsymbol{\lambda} = [\lambda_m(x), \lambda_m(x^*)]^\top$ , which is related to  $\boldsymbol{\mu}$  via (20). The function  $\mathcal{L}$  linking  $\boldsymbol{\lambda}$  and  $\boldsymbol{\mu}$  is defined in (18), which is well-behaved as will be apparent below. Define  $\hat{\boldsymbol{\lambda}} = [\hat{\lambda}_m(x), \hat{\lambda}_m(x^*)]^\top$  with  $\hat{\lambda}_m(\tilde{x}) = \mathcal{L}'^{-1}\{\hat{\boldsymbol{\mu}}_m(\tilde{x})\}$  for  $\tilde{x} = x, x^*$ .

**Theorem 9.** Under assumption A,  $n^r(\hat{\boldsymbol{\lambda}} - \boldsymbol{\lambda}) \xrightarrow{d} \mathbb{D}_{\boldsymbol{\lambda}}^{-1} \mathbf{X}$ , where  $\mathbb{D}_{\boldsymbol{\lambda}}$  is a diagonal matrix with diagonal elements  $\mathcal{L}'\{\lambda_m(x)\}$  and  $\mathcal{L}'\{\lambda_m(x^*)\}$ .  $\square$

In theorem 9 invertibility of  $\mathbb{D}_{\boldsymbol{\lambda}}$  is implicit. Theorem 9 establishes that the limit distribution of  $\hat{\boldsymbol{\lambda}}$  follows from that of  $\hat{\boldsymbol{\mu}}$ , which was assumed known in assumption A. Since the counterfactual predictions are simple functions of  $k^*$ , conducting inference on them is just a matter of applying the Delta method. The above results do not depend on which counterfactual one considers but, as mentioned, we focus on  $q_v$  from here on.

Making the dependence of the ME solution of  $\rho$  on  $\boldsymbol{\lambda}$  explicit in our notation, we obtain

$$\rho(x, x^*; \boldsymbol{\lambda}) = \int_0^1 k^*\{v, x; \lambda_m(x)\} k^*\{v, x^*; \lambda_m(x^*)\} dv - \mathcal{L}'\{\lambda_m(x)\} \mathcal{L}'\{\lambda_m(x^*)\},$$

where  $\hat{\boldsymbol{\rho}}$  is defined by replacing  $\boldsymbol{\lambda}$  with  $\hat{\boldsymbol{\lambda}}$ . The Delta method then produces

$$n^r\{\hat{\boldsymbol{\rho}}(x, x^*) - \rho(x, x^*)\} = \mathbb{D}_{\boldsymbol{\rho}}(\boldsymbol{\lambda}) n^r(\hat{\boldsymbol{\lambda}} - \boldsymbol{\lambda}) + o_p(1),$$

where  $\mathbb{D}_{\boldsymbol{\rho}}$  is the Jacobian matrix of  $\rho$ .

## 9. Conclusions

We have shown in the paper that the problem of counterfactual prediction in an incomplete model has issues that are distinct from the familiar problems with identification and

estimation of model parameters. We have considered five approaches in the specific context of a complete information binary decision game with pure strategies and Nash equilibria: bounding the counterfactual prediction probabilities, using the midpoint prediction, taking a minmax approach with an entropy constraint, minimizing the average risk, and maximum entropy prediction. On balance, we prefer the maximum entropy approach since it does not have some of the unattractive features that other approaches have and is comparatively straightforward to implement.

Our results can be applied to other (counterfactual) prediction problems with partially identified parameters. It may further be of interest to study the relative merits of maximum entropy for the purpose of selecting a single point from the identified set in a partially identified world.

## A. Proofs

*Proof of Theorem 1.* First,

$$q_v(y^* | x^*, x, y) = \frac{\mathbb{P}\{y(\mathbf{e}^*, \mathbf{u}^*, \mathbf{v}^*, x^*) = y^*, y(\mathbf{e}, \mathbf{u}, \mathbf{v}^*, x) = y\}}{\tilde{\pi}_y(x)} =$$

$$\tilde{\pi}_{y^*}(x^*) + \frac{\pi_m(x)\pi_m(x^*)}{\tilde{\pi}_y(x)} \times$$

$$\text{Cov}[\mathbb{1}\{y(\mathbf{e}^*, \mathbf{u}^*, \mathbf{v}^*, x^*) = y^*\}, \mathbb{1}\{y(\mathbf{e}, \mathbf{u}, \mathbf{v}^*, x) = y\} | \mathbf{e} \in S_m(x), \mathbf{e}^* \in S_m(x^*)]. \quad (43)$$

The covariance in (43) equals zero unless both  $y$  and  $y^*$  belong to  $\{(1, 0), (0, 1)\}$ . For  $y = y^* = (1, 0)$  the covariance in (43) equals (by the law of iterated expectations)

$$\text{Cov}\{k(\mathbf{v}, x^*), k(\mathbf{v}, x)\} = \rho(x, x^*).$$

The remaining three cases follow analogously.  $\square$

*Proof of Theorem 2.* Since  $k(\mathbf{v}, x)$  is for all  $x$  nondecreasing in  $v$ , it follows that  $k(\mathbf{v}, x)$  and  $k(\mathbf{v}, x^*)$  have nonnegative covariance, which can be made to equal zero, e.g. by making  $k(\mathbf{v}, x) = \mu_m(x)$  for all values of  $v$ . This establishes the lower bound.

For the upper bound, note that the covariance is maximized if  $k(\mathbf{v}, x^*) = k(\mathbf{v}, x)\mu_m(x^*) / \mu_m(x)$ , resulting in the stated upper bound. This upper bound is attained by setting  $k(\mathbf{v}, x) = \mathbb{1}\{v > 1 - \mu_m(x)\}$ .  $\square$

*Proof of Theorem 3.* We have

$$\mathbb{P}\{y(\mathbf{u}^*, \mathbf{e}, \mathbf{v}, x^*) = y^*, y(\mathbf{u}, \mathbf{e}, \mathbf{v}, x) = y\}$$

$$= \mathbb{E}\{c_y(\mathbf{e}, x)c_{y^*}(\mathbf{e}, x^*)\} + \mathbb{E}\{c_y(\mathbf{e}, x)c_m(\mathbf{e}, x^*)b_{y^*}(\mathbf{e}, \mathbf{v}, x^*)\}$$

$$+ \mathbb{E}\{c_m(\mathbf{e}, x)c_{y^*}(\mathbf{e}, x^*)b_y(\mathbf{e}, \mathbf{v}, x)\} + \mathbb{E}\{c_m(\mathbf{e}, x)c_m(\mathbf{e}, x^*)b_y(\mathbf{e}, \mathbf{v}, x)b_{y^*}(\mathbf{e}, \mathbf{v}, x^*)\}$$

$$= \pi_{mm}\text{Cov}\{b_y(\mathbf{e}, \mathbf{v}, x), b_{y^*}(\mathbf{e}, \mathbf{v}, x^*) | \mathbf{e} \in S_{mm}(x, x^*)\}$$

$$+ \text{Cov}\{a_y(\mathbf{e}, x), a_{y^*}(\mathbf{e}, x^*)\} + \mathbb{E}a_y(\mathbf{e}, x)\mathbb{E}a_{y^*}(\mathbf{e}, x^*). \quad (44)$$

The stated result follows if one divides the right hand side in (44) by  $\tilde{\pi}_y(x) = \mathbb{E}a_y(\mathbf{e}, x)$ .  $\square$

*Proof of Theorem 4.* Partition  $S_m(x)$  into four disjoint regions:  $S_A = S_m(x) \cap S_{10}(x^*)$ ,  $S_B = S_m(x) \cap S_{01}(x^*)$ ,  $S_C = S_m(x) \cap S_m(x^*)$ , and  $S_D = S_m(x) \cap \{S_{00}(x^*) \cup S_{11}(x^*)\}$ . Let  $\pi_A = \mathbb{P}(\mathbf{e} \in S_A)$  and let  $\mu_A = \mathbb{E}\{b_y(\mathbf{e}, \mathbf{v}, x) \mid \mathbf{e} \in S_A\}$ . Let  $\pi_B, \mu_C$ , etcetera, be similarly defined. We then solve

$$\begin{aligned} \min_{\mu_A, \mu_B, \mu_C, \mu_D} \quad & (\pi_A \mu_A + \pi_B \mu_B + \pi_C \mu_C) \\ \text{subject to} \quad & \begin{cases} 0 \leq \mu_A, \mu_B, \mu_C, \mu_D \leq 1, \\ \pi_A \mu_A + \pi_B \mu_B + \pi_C \mu_C + \pi_D \mu_D = \phi_y(x), \end{cases} \end{aligned}$$

to obtain the lower bound, noting that the solution has

$$\pi_A \mu_A + \pi_B \mu_B + \pi_C \mu_C = \max\{0, \phi_y(x) - \pi_D\},$$

because we are trying to make  $\mu_D$  as large as possible. Maximizing the same objective function subject to the same constraints yields the upper bound. Since  $\mu_A, \mu_B, \mu_C, \mu_D$  are otherwise unconstrained, the bounds are sharp.  $\square$

*Proof of Theorem 5.* When  $x = x^*$ , we have

$$\begin{aligned} \mathbb{P}\{y(\mathbf{e}, \mathbf{u}^*, \mathbf{v}, x^*) = y^*, y(\mathbf{e}, \mathbf{u}, \mathbf{v}, x) = y\} &= \pi_{yy^*}(x, x) \\ &+ \pi_m(x) \mathbb{E}\{\bar{b}_y(\mathbf{e}, x) \bar{b}_{y^*}(\mathbf{e}, x) \mid \mathbf{e} \in S_m(x)\} + \delta_y \delta_{y^*} \pi_m(x) \mathbb{E}\{\mathbb{V}(p(\mathbf{e}, \mathbf{v}, x) \mid \mathbf{e}) \mid \mathbf{e} \in S_m(x)\}, \end{aligned}$$

where  $\bar{b}_y(\mathbf{e}, x) = \mathbb{E}b_y(\mathbf{e}, \mathbf{v}, x)$ . The result follows from the fact that

$$0 \leq \mathbb{V}\{p(\mathbf{e}, \mathbf{v}, x) \mid \mathbf{e}\} \leq \mathbb{E}\{p(\mathbf{e}, \mathbf{v}, x) \mid \mathbf{e}\} [1 - \mathbb{E}\{p(\mathbf{e}, \mathbf{v}, x) \mid \mathbf{e}\}],$$

where the upper bound is attained when  $p(\mathbf{e}, \mathbf{v}, x)$  is binary taking one with probability  $\mathbb{E}p(\mathbf{e}, \mathbf{v}, x)$ , and the lower bound is attained when  $p(\mathbf{e}, \mathbf{v}, x) = \mu_m(x)$  with probability one.

So, we focus on the case  $x \neq x^*$ . Some tedious but simple mathematical manipulations show that  $\tilde{\pi}_y(x) q_{ev}(y^* \mid x^*, x, y)$  can alternatively be expressed as

$$\begin{aligned} & \mathbb{E}\left[\{c_y(\mathbf{e}, x) + c_m(\mathbf{e}, x) b_y(\mathbf{e}, \mathbf{v}, x)\} \{c_{y^*}(\mathbf{e}, x^*) + c_m(\mathbf{e}, x^*) b_{y^*}(\mathbf{e}, \mathbf{v}, x^*)\}\right] \\ &= \pi_{yy^*}(x, x^*) + \pi_{my^*}(x, x^*) \mathbb{E}\{\bar{b}_y(\mathbf{e}, x) \mid \mathbf{e} \in S_{my^*}(x, x^*)\} \\ & \quad + \pi_{ym}(x, x^*) \mathbb{E}\{\bar{b}_{y^*}(\mathbf{e}, x^*) \mid \mathbf{e} \in S_{ym}(x, x^*)\} \\ & \quad + \pi_{mm}(x, x^*) \mathbb{E}\{b_y(\mathbf{e}, \mathbf{v}, x) b_{y^*}(\mathbf{e}, \mathbf{v}, x^*) \mid \mathbf{e} \in S_{mm}(x, x^*)\}. \quad (45) \end{aligned}$$

Note that  $b_y(\mathbf{e}, \mathbf{v}, x)$  is unrestricted in  $\mathbf{e} \in S_m(x)$  but is restricted (i.e. monotonic) in  $\mathbf{v}$ , albeit that we have the condition  $\mathbb{E}\{p(\mathbf{e}, \mathbf{v}, x) \mid \mathbf{e} \in S_m(x)\} = \mu_m(x)$  and  $b_y(\mathbf{e}, \mathbf{v}, x)$  must belong to  $[0, 1]$ . Therefore, the bounds we seek can be attained for functions  $p$  that are flat in  $\mathbf{v}$ , which means that there is no loss of generality in dropping  $\mathbf{v}$  from the notation, which we do from hereon.

Now, let  $S_r^y(x) = [0, 1]^2 - S_m(x) - S_y(x)$ . Further, let  $S_{mr^*}(x, x^*) = S_m(x) \cap S_r^{y^*}(x^*)$  and  $S_{rm}(x, x^*) = S_r^y(x) \cap S_m(x^*)$ . We then define

$$\bar{z}_{my^*}^{xy}(x, x^*) = \mathbb{E}\{b_y(\mathbf{e}, x) \mid \mathbf{e} \in S_{my^*}(x, x^*)\},$$

and let  $\bar{z}$ 's with different subscript and superscript combinations be defined analogously. Then the right hand side in (45) reduces to

$$\begin{aligned} \pi_{yy^*}(x, x^*) + \pi_{my^*}(x, x^*)\bar{z}_{my^*}^{xy}(x, x^*) + \pi_{ym}(x, x^*)\bar{z}_{ym}^{x^*y^*}(x, x^*) \\ + \pi_{mm}(x, x^*)\mathbb{E}\{b_y(e, x)b_{y^*}(e, x^*) \mid e \in S_{mm}(x, x^*)\}. \end{aligned} \quad (46)$$

Therefore, we seek to minimize/maximize (46) subject to the (mean) restrictions, namely (dropping the  $(x, x^*)$  arguments)

$$\begin{cases} \pi_{ym}\bar{z}_{ym}^{x^*y^*} + \pi_{mm}\bar{z}_{mm}^{x^*y^*} + \pi_{rm}\bar{z}_{rm}^{x^*y^*} = \phi_{y^*}(x^*), \\ \pi_{my^*}\bar{z}_{my^*}^{xy} + \pi_{mm}\bar{z}_{mm}^{xy} + \pi_{mr^*}\bar{z}_{mr^*}^{xy} = \phi_y(x). \end{cases} \quad (47)$$

Since  $b_y(e, x)$  must belong to  $[0, 1]$ , (46) implies that for the lower bound we want to strive to make  $b_y(e, x)b_{y^*}(e, x^*)$  equal to zero for  $e \in S_{mm}(x, x^*)$  and equal to one for the upper bound. This implies that there is no loss of generality in assuming that  $b_y(e, x)$  is binary for all  $e, y$ . So, we assume that  $b_y(e, x)$  is either one or zero hereafter. Define

$$S_+^y(x) = \{e \in S_m(x) : b_y(e, x) = 1\} \quad \text{and} \quad S_-^y(x) = \{e \in S_m(x) : b_y(e, x) = 0\},$$

which forms a partition of  $S_m(x)$ .

The optimization of (46) subject to (47) is illustrated in figure 16: minimization at the top and maximization at the bottom.

Consider the lower bound first. Define  $a, b, c, d$  to be the probabilities illustrated at the top of figure 16 so that  $\mathbb{P}(y = y, y^* = y^* | x = x, x^* = x^*) = a + b + c + c$ . Then, note that  $a, b, c, d$  correspond to the four terms in (46), (e.g.  $\pi_{my^*}\bar{z}_{my^*}^{xy} = b$ ), because by construction  $b_y(e, x)$  is equal to one when  $e \in S_+^y(x)$  and equal to zero, otherwise. The constraints in (47) can likewise be expressed in terms of the probabilities illustrated at the top of figure 16. So, using the shorthand  $\phi = \phi_y(x)$  and  $\phi^* = \phi_{y^*}(x^*)$ , minimizing (46) subject to (47) can be formulated as

$$\begin{cases} \min_{a,b,c,d} (a + b + c + d) \\ \text{s.t. } c + d + f \geq \max(\phi^* - \pi_{rm}, 0) \\ \quad b + d + e \geq \max(\phi - \pi_{mr^*}, 0) \\ \quad d + e + f \leq \pi_{mm} \\ \quad a = \pi_{yy^*} \\ \quad a, b, c, d, e, f \geq 0 \end{cases}$$

Now,

$$\begin{aligned} b + c + d &= (b + d + e) + (c + d + f) - (d + e + f) \\ &\geq \max\{\max(\phi^* - \pi_{rm}, 0) + \max(\phi - \pi_{mr^*}, 0) - \pi_{mm}, 0\} \\ &= \max(\phi^* - \pi_{rm} - \pi_{mm}, \phi - \pi_{mr^*} - \pi_{mm}, \phi + \phi^* - \pi_{mr^*} - \pi_{rm} - \pi_{mm}, 0). \end{aligned}$$

Finding the upper bound can likewise be formulated as a maximization problem. Define  $a, b, c, d$  as the probabilities illustrated at the bottom of figure 16. Then, maximizing (46) subject to (47) can be formulated as

$$\begin{cases} \max_{a,b,c,d} (a + b + c + d) \\ \text{s.t. } b \leq \min(\phi - d, \pi_{ym}) \\ c \leq \min(\phi^* - d, \pi_{my^*}) \\ 0 \leq d \leq \pi_{mm}, \\ a = \pi_{yy^*}. \end{cases}$$

Now,

$$\begin{aligned} b + c + d &\leq \min(\phi - d, \pi_{my^*}) + \min(\phi^* - d, \pi_{ym}) + d \\ &\leq \min(\phi + \phi^*, \phi + \pi_{ym}, \phi^* + \pi_{my^*}, \pi_{my^*} + \pi_{ym} + \pi_{mm}). \end{aligned}$$

Since the proof is constructive, the bounds are sharp.  $\square$

*Proof of Theorem 6.* This is a Lagrangean optimization problem with a function-valued parameter. The first order conditions are

$$\begin{cases} f_e(e)f_x(x)\{1 + \log f^*(p | e, x)\} - v_1(e, x) - v_2(x)pf_e(e)\mathbb{1}\{e \in S_m(x)\} = 0, \\ \int f^*(p | e, x) dp = 1, \\ \int_{S_m(x)} \int_0^1 pf^*(p | e, x) dp f_e(e) de = \mu_m(x)\pi_m(x), \end{cases}$$

almost everywhere, where  $v_1, v_2$  are Lagrangean parameters. Thus, taking  $\lambda(x) = v_2(x) / f_x(x)$ , it follows that for all  $p \in [0, 1]$ ,

$$f^*(p | e, x) \propto \begin{cases} \exp\{\lambda(x)p\}, & e \in S_m(x), \\ 1, & e \notin S_m(x), \end{cases}$$

which yields (19).

Further, the second condition in (20) imposes that  $f^*$  integrate to one and the first that  $f^*$  has the correct mean since  $\mathcal{L}'(\lambda) = I'(\lambda) / I(\lambda) = \int_0^1 p \exp(p\lambda) dp / \int_0^1 \exp(p\lambda) dp$ .  $\square$

*Proof of Theorem 7.* From (19) it follows that

$$f^*(p | e, x) = \begin{cases} 1, & e \notin S_m(x), \\ A\{p, \lambda_m(x)\}, & e \in S_m(x). \end{cases}$$

Inverting the conditional distribution function  $\int_0^p f^*(s | e, x) ds$  yields the stated result.  $\square$

**Lemma 1.**  $\mathcal{L}'$  is continuous, increasing, and nonzero everywhere.

*Proof.* Continuity is trivial by L'Hôpital. Further,

$$\mathcal{L}''(\lambda) = \begin{cases} 1/12, & \text{if } \lambda = 0, \\ 1/\lambda^2 - \exp(\lambda)/\{\exp(\lambda) - 1\}^2, & \text{if } \lambda \neq 0, \end{cases}$$

which is continuous by L'Hôpital, also. Indeed, for  $\lambda \neq 0$ ,

$$\begin{aligned} \lambda^2 \{\exp(\lambda) - 1\}^2 \mathcal{L}''(\lambda) &= \exp(2\lambda) - (2 + \lambda^2) \exp(\lambda) + 1 \\ &= \exp(\lambda) \{\exp(\lambda) + \exp(-\lambda) - 2 - \lambda^2\} = 2 \exp(\lambda) \sum_{j=2}^{\infty} \frac{\lambda^{2j}}{(2j)!} > 0. \end{aligned}$$

Finally, for the last assertion, simply note that  $\lim_{\lambda \rightarrow -\infty} \mathcal{L}'(\lambda) = 0$  and that  $\mathcal{L}'$  is increasing.  $\square$

*Proof of Theorem 9.* By lemma 1, we know that  $\mathcal{L}'^{-1}$  is differentiable everywhere. So,

$$n^r \{\hat{\lambda}_m(\tilde{x}) - \lambda_m(\tilde{x})\} = \frac{n^r \{\hat{\mu}_m(\tilde{x}) - \mu_m(\tilde{x})\}}{\mathcal{L}'[\mathcal{L}'^{-1}\{\tilde{\mu}_m(\tilde{x})\}]}, \quad \tilde{x} = x, x^*,$$

where  $\tilde{\mu}_m(\tilde{x})$  is between  $\hat{\mu}_m(\tilde{x})$  and  $\mu_m(\tilde{x})$ . Then, the result follows from lemma 1 and assumption A.  $\square$

## B. Other cases

Below, we present some results on prediction probabilities other than  $q, q_{ev}, q_v$ . The proofs refer back to details in earlier proofs, especially that of theorem 5.

**Theorem 10.**  $q_e(y^* | x^*, x, y) = \tilde{\pi}_{y^*}(x^*) + \alpha_{yy^*}(x, x^*)$ . Further, the bounds on  $q_e(y^* | x^*, x, y)$  are identical to those on  $q_{ev}(y^* | x^*, x, y)$  given in theorem 5.

*Proof.* Recall that  $\bar{b}_y(e, x) = \mathbb{E} b_y(e, \mathbf{v}, x)$ . Since  $b_y(e, \mathbf{v}, x)$  and  $b_{y^*}(e, \mathbf{v}^*, x^*)$  are independent conditional on  $e$ , we have

$$\begin{aligned} &\mathbb{P}\{y(\mathbf{e}, \mathbf{u}^*, \mathbf{v}^*, x^*) = y^*, y(\mathbf{e}, \mathbf{u}, \mathbf{v}, x) = y\} \\ &= \mathbb{E}\{c_y(\mathbf{e}, x)c_{y^*}(\mathbf{e}, x^*)\} + \mathbb{E}\{c_y(\mathbf{e}, x)c_m(\mathbf{e}, x^*)\bar{b}_{y^*}(\mathbf{e}, x^*)\} \\ &\quad + \mathbb{E}\{c_m(\mathbf{e}, x)c_{y^*}(\mathbf{e}, x^*)\bar{b}_y(\mathbf{e}, x)\} + \mathbb{E}\{c_m(\mathbf{e}, x)c_m(\mathbf{e}, x^*)\bar{b}_y(\mathbf{e}, x)\bar{b}_{y^*}(\mathbf{e}, x^*)\} \\ &= \tilde{\pi}_y(x)\tilde{\pi}_{y^*}(x^*) + \text{Cov}\{a_y(\mathbf{e}, x), a_{y^*}(\mathbf{e}, x^*)\} \\ &= \pi_{yy^*}(x, x^*) + \pi_{ym}(x, x^*)\mathbb{E}\{\bar{b}_{y^*}(\mathbf{e}, x^*) \mid \mathbf{e} \in S_{ym}(x, x^*)\} \\ &\quad + \pi_{my^*}(x, x^*)\mathbb{E}\{\bar{b}_y(\mathbf{e}, x) \mid \mathbf{e} \in S_{my^*}(x, x^*)\} \\ &\quad + \pi_{mm}(x, x^*)\mathbb{E}\{\bar{b}_y(\mathbf{e}, x)\bar{b}_{y^*}(\mathbf{e}, x^*) \mid \mathbf{e} \in S_{mm}(x, x^*)\}, \end{aligned}$$

which coincides with (45), except that  $b$  in the last term of (45) is replaced with  $\bar{b}$  here. Since the bounds in theorem 5 obtain for functions  $\rho$  that are flat in  $v$ , as noted in the proof of theorem 5, the bounds must coincide.  $\square$

**Theorem 11.**

$$q_u(y^* | x^*, x, y) = \tilde{\pi}_{y^*}(x^*) + \delta_y \delta_{y^*} \frac{\pi_m(x) \pi_m(x^*)}{\tilde{\pi}_y(x)} \rho_u(x, x^*),$$

where  $\rho_u(x, x^*) = \mathbb{E}\{\min\{p(\mathbf{e}, \mathbf{v}, x), p(\mathbf{e}^*, \mathbf{v}^*, x^*)\} | \mathbf{e} \in S_m(x), \mathbf{e}^* \in S_m(x^*)\} - \mu_m(x) \mu_m(x^*)$ . Further, the bounds on  $q_u$  are identical to the bounds of  $q_v$  given in theorem 2.

*Proof.* The expression for  $q_u$  follows from simple algebra. For the bounds, we focus on the case  $y = y^* = (1, 0)$ : the other cases are analogous. Since  $\mathbb{E}\{\min(\mathbf{p}, \mathbf{p}^*) | \cdot\} \leq \min\{\mathbb{E}(\mathbf{p} | \cdot), \mathbb{E}(\mathbf{p}^* | \cdot)\}$ , the upper bound of  $\rho_u(x, x^*)$  is given by  $\min\{\mu_m(x), \mu_m(x^*)\} - \mu_m(x) \mu_m(x^*)$ . This upper bound is attained when  $\mathbf{p}$  and  $\mathbf{p}^*$  are binary. For the lower bound, note that  $\mathbb{E}\{\min(\mathbf{p}, \mathbf{p}^*) | \cdot\} \geq \mathbb{E}(\mathbf{p} \mathbf{p}^* | \cdot)$ , where  $\mathbf{p}, \mathbf{p}^*$  are independent given  $\mathbf{e}, \mathbf{e}^*$ .  $\square$

**Theorem 12.**

$$q_{uv}(y^* | x^*, x, y) = \tilde{\pi}_{y^*}(x^*) + \delta_y \delta_{y^*} \frac{\pi_m(x) \pi_m(x^*)}{\tilde{\pi}_y(x)} \rho_{uv}(x, x^*), \quad (48)$$

where  $\rho_{uv}(x, x^*) = \mathbb{E}[\min\{p(\mathbf{e}, \mathbf{v}, x), p(\mathbf{e}^*, \mathbf{v}^*, x^*)\} | \mathbf{e} \in S_m(x), \mathbf{e}^* \in S_m(x^*)] - \mu_m(x) \mu_m(x^*)$ . Further, the bounds of  $q_{uv}$  are identical to those on  $q_v$  given in theorem 2.

*Proof.* Define

$$h_y(\mathbf{e}, u, \mathbf{v}, x) = \begin{cases} \mathbb{1}\{u \leq p(\mathbf{e}, \mathbf{v}, x)\}, & y = (1, 0), \\ \mathbb{1}\{u > p(\mathbf{e}, \mathbf{v}, x)\}, & y = (0, 1), \\ 0, & y \in \{(0, 0), (1, 1)\}. \end{cases}$$

Note that

$$\begin{aligned} & \mathbb{P}\{y(\mathbf{e}^*, \mathbf{u}, \mathbf{v}, x^*) = y^*, y(\mathbf{e}, \mathbf{u}, \mathbf{v}, x) = y\} \\ &= \tilde{\pi}_y(x) \tilde{\pi}_{y^*}(x, x^*) + \\ & \quad \pi_m(x) \pi_m(x^*) \text{Cov}\{h_y(\mathbf{e}, \mathbf{u}, \mathbf{v}, x), h_{y^*}(\mathbf{e}^*, \mathbf{u}, \mathbf{v}, x^*) | \mathbf{e} \in S_m(x), \mathbf{e}^* \in S_m(x^*)\}. \end{aligned} \quad (49)$$

Tedious but simple algebra shows that the covariance in (49) equals  $\delta_y \delta_{y^*} \rho_{uv}(x, x^*)$ , which yields (48).

Now we establish the bounds. We focus on the case  $y = y^* = (1, 0)$ ; the other cases follow analogously. For the upper bound, note that  $\mathbb{E}\{\min(\mathbf{p}, \mathbf{p}^*) | \cdot\} \leq \min\{\mathbb{E}(\mathbf{p} | \cdot), \mathbb{E}(\mathbf{p}^* | \cdot)\}$ , with equality for  $p(\mathbf{e}, \mathbf{v}, x) = \mathbb{1}\{v \geq 1 - \mu_m(x)\}$ . For the lower bound,

$$\mathbb{E}\{\min(\mathbf{p}, \mathbf{p}^*) | \mathbf{e}, \mathbf{e}^*\} \geq \mathbb{E}(\mathbf{p} \mathbf{p}^* | \mathbf{e}, \mathbf{e}^*) \geq \mathbb{E}(\mathbf{p} | \mathbf{e}) \mathbb{E}(\mathbf{p}^* | \mathbf{e}^*) \quad (50)$$

where the inequalities hold with equality for  $p(\mathbf{e}, \mathbf{v}, x) = \mathbb{1}\{\mathbf{e} \in S_m(x)\}$ . Take the expectations in (50) over  $S_m(x)$  and  $S_m(x^*)$  to obtain the sharp lower bound.  $\square$

For the remaining two cases ( $q_{eu}$  and  $q_{ev}$ ), we define

$$g_y(\mathbf{e}, \mathbf{v}, x) = c_y(\mathbf{e}, x) + c_m(\mathbf{e}, x) b_y(\mathbf{e}, \mathbf{v}, x).$$

**Theorem 13.**

$$q_{eu}(y^* | x^*, x, y) = \tilde{\pi}_{y^*}(x^*) + \frac{\text{Cov}\{g_y(\mathbf{e}, \mathbf{v}, x), g_{y^*}(\mathbf{e}, \mathbf{v}^*, x^*)\}}{\tilde{\pi}_y(x)} + \delta_y \delta_{y^*} \frac{\pi_{mm}(x, x^*)}{\tilde{\pi}_y(x)} \rho_{eu}(x, x^*),$$

where  $\rho_{eu}(x, x^*) = \mathbb{E}[\min\{p(\mathbf{e}, \mathbf{v}, x), p(\mathbf{e}, \mathbf{v}^*, x^*)\} - p(\mathbf{e}, \mathbf{v}, x)p(\mathbf{e}, \mathbf{v}^*, x^*) \mid \mathbf{e} \in S_{mm}(x, x^*)]$ . Further, the bounds on  $q_{eu}(y^* | x^*, x, y)$  are identical to those of  $q_{ev}(y^* | x^*, x, y)$  given in theorem 5.

*Proof.* The formula for  $q_{eu}$  follows from tedious algebra, so we focus on the bounds. We consider the case  $y = y^*(1, 0)$ : the other combinations follow analogously. Now,  $\tilde{\pi}_y(x)q_{eu}(y^* | x^*, x, y)$  equals

$$\begin{aligned} & \pi_{yy^*}(x, x^*) + \pi_{ym}(x, x^*)\mathbb{E}\{p(\mathbf{e}, \mathbf{v}^*, x^*) \mid \mathbf{e} \in S_{ym}(x, x^*)\} \\ & \quad + \pi_{my^*}(x, x^*)\mathbb{E}\{p(\mathbf{e}, \mathbf{v}, x) \mid \mathbf{e} \in S_{my^*}(x, x^*)\} \\ & \quad + \pi_{mm}(x, x^*)\mathbb{E}[\min\{p(\mathbf{e}, \mathbf{v}, x), p(\mathbf{e}, \mathbf{v}^*, x^*)\} \mid \mathbf{e} \in S_{mm}(x, x^*)]. \end{aligned} \quad (51)$$

If  $x = x^*$ , the second and third terms in (51) are equal to zero, so we only need to consider the last term. The upper bound follows from the fact that

$$\mathbb{E}\{\min\{\mathbf{p}, \mathbf{p}^*\} \mid \cdot\} \leq \min\{\mathbb{E}(\mathbf{p} \mid \cdot), \mathbb{E}(\mathbf{p}^* \mid \cdot)\}. \quad (52)$$

For the lower bound, note that by the Jensen inequality

$$\begin{aligned} \mathbb{E}\{\min(\mathbf{p}, \mathbf{p}^*) \mid \mathbf{e} \in S_m(x)\} &= \mathbb{E}\left[\int_0^1 \mathbb{P}\{\min(\mathbf{p}, \mathbf{p}^*) > t \mid \mathbf{e}\} dt \mid \mathbf{e} \in S_m(x)\right] \\ &\geq \left[\mathbb{E}\left\{\int_0^1 \mathbb{P}(\mathbf{p} > t \mid \mathbf{e}) dt \mid \mathbf{e} \in S_m(x)\right\}\right]^2 = \mu_m^2(x). \end{aligned}$$

If  $x \neq x^*$  then (52) and  $\mathbb{E}\{p(\mathbf{e}, \mathbf{v}^*, x^*) \mid \mathbf{e} \in S_m(x^*)\} = \mathbb{E}\{p(\mathbf{e}, \mathbf{v}, x^*) \mid \mathbf{e} \in S_m(x^*)\} = \mu_m(x^*)$  imply that (51) is bounded above by  $\pi_{yy^*}(x, x^*) + U_{yy^*}(x, x^*)$ , where  $U_{yy^*}$  is as defined in (15). Sharpness follows from the fact that (51) is bounded below by (45). For the lower bound, note that  $\mathbb{E}\{\min\{\mathbf{p}, \mathbf{p}^*\} \mid \mathbf{e} \in S_{mm}(x, x^*)\} \geq \mathbb{E}\{\bar{b}_y(\mathbf{e}, x)\bar{b}_{y^*}(\mathbf{e}, x^*) \mid \mathbf{e} \in S_{mm}(x, x^*)\}$ , which is identical to the second factor of the last term in (45), except that  $b_y$  is replaced with  $\bar{b}_y$ .  $\square$

**Theorem 14.**

(i) If  $x^* = x$  then  $q_{euv}(y^* | x^*, x, y) = \mathbb{1}(y^* = y)$ ;

(ii) If  $x^* \neq x$  then

$$q_{euv}(y^* | x^*, x, y) = \tilde{\pi}_{y^*}(x^*) + \frac{\text{Cov}\{g_{y^*}(\mathbf{e}, \mathbf{v}, x^*), g_y(\mathbf{e}, \mathbf{v}, x)\}}{\tilde{\pi}_y(x)} + \delta_y \delta_{y^*} \frac{\pi_{mm}(x, x^*)}{\tilde{\pi}_y(x)} \rho_{euv}(x, x^*),$$

where  $\rho_{euv}(x, x^*) = \mathbb{E}[\min\{p(\mathbf{e}, \mathbf{v}, x), p(\mathbf{e}, \mathbf{v}, x^*)\} - p(\mathbf{e}, \mathbf{v}, x)p(\mathbf{e}, \mathbf{v}, x^*) \mid \mathbf{e} \in S_{mm}(x, x^*)]$ ;



(iii) If  $x^* \neq x$  then the bounds on  $q_{euv}$  coincide with those on  $q_{ev}$ .

*Proof.* Part (i) is trivial and part (ii) tedious yet mechanical, so we focus on part (iii). We establish the result for the case  $y = y^* = (1, 0)$ , where the other combinations follow analogously.

Now, note from part (ii) that  $\tilde{\pi}_y(x)q_{euv}(y^* | x^*, x, y)$  equals

$$\begin{aligned} & \pi_{yy^*}(x, x^*) + \pi_{ym}(x, x^*)\mathbb{E}\{\rho(\mathbf{e}, \mathbf{v}, x^*) | \mathbf{e} \in S_{ym}(x, x^*)\} \\ & \quad + \pi_{my^*}(x, x^*)\mathbb{E}\{\rho(\mathbf{e}, \mathbf{v}, x) | \mathbf{e} \in S_{my^*}(x, x^*)\} \\ & \quad + \pi_{mm}(x, x^*)\mathbb{E}[\min\{\rho(\mathbf{e}, \mathbf{v}, x), \rho(\mathbf{e}, \mathbf{v}, x^*)\} | \mathbf{e} \in S_{mm}(x, x^*)]. \end{aligned} \quad (53)$$

Because  $\mathbb{E}\{\min(\mathbf{p}, \mathbf{p}^*) | \cdot\} \leq \min\{\mathbb{E}(\mathbf{p} | \cdot), \mathbb{E}(\mathbf{p}^* | \cdot)\}$ ,  $\mathbb{E}\{\rho(\mathbf{e}, \mathbf{v}, x) | \mathbf{e} \in S_m(x)\} = \mu_m(x)$ , and  $\mathbb{E}\{\rho(\mathbf{e}, \mathbf{v}, x^*) | \mathbf{e} \in S_m(x^*)\} = \mu_m(x^*)$ , (53) is bounded above by  $\pi_{yy^*}(x, x^*) + U_{yy^*}(x, x^*)$ , where  $U_{yy^*}$  is as defined in (15). For sharpness, note that (53) is bounded below by (45).

Now the lower bound. Since (53) is no less than (45), the lower bound for (45) is a lower bound here, also. It remains to be shown that it can be attained. Recall from the proof of theorem 5 that (45) attains the lower bound  $\pi_{yy^*}(x, x^*) + L_{yy^*}(x, x^*)$  when  $\mathbf{p}$  and  $\mathbf{p}^*$  are binary, in which case  $\mathbf{p}\mathbf{p}^* = \min(\mathbf{p}, \mathbf{p}^*)$ . Hence, (45) equals (53).  $\square$

## C. Useful formulas

In this appendix we provide the formulas for the correction terms in theorems 1 and 3 that are implied by the maximum entropy solution  $f^*$ . As we explained below theorem 7, the quantity  $\rho(x, x^*)$  entering the correction terms is for the maximum entropy solution the same for  $q_v$  and  $q_{ev}$ : i.e.

$$\begin{aligned} \rho^*(x, x^*) &= \text{Cov}\{k^*(\mathbf{v}, x)k^*(\mathbf{v}, x^*)\} \\ &= \mathbb{E}[\text{Cov}\{\rho^*(\mathbf{e}, \mathbf{v}, x)\rho^*(\mathbf{e}, \mathbf{v}, x^*) | \mathbf{e}\} | \mathbf{e} \in S_{mm}(x, x^*)] = \\ & \left\{ \begin{array}{ll} \int_0^1 \frac{\log(1 + v[\exp\{\lambda_m(x)\} - 1]) \log(1 + v[\exp\{\lambda_m(x^*)\} - 1])}{\lambda_m(x)\lambda_m(x^*)} dv & \lambda_m(x), \lambda_m(x^*) \neq 0, \\ \frac{1}{4\lambda_m(x)} + \frac{1}{2\lambda_m(x)\{e^{\lambda_m(x)} - 1\}} - \frac{e^{\lambda_m(x)}}{2\{e^{\lambda_m(x)} - 1\}^2}, & \lambda_m(x) \neq \lambda_m(x^*) = 0, \\ \frac{1}{4\lambda_m(x^*)} + \frac{1}{2\lambda_m(x^*)\{e^{\lambda_m(x^*)} - 1\}} - \frac{e^{\lambda_m(x^*)}}{2\{e^{\lambda_m(x^*)} - 1\}^2}, & \lambda_m(x^*) \neq \lambda_m(x) = 0, \\ \frac{1}{12}, & \lambda_m(x) = \lambda_m(x^*) = 0. \end{array} \right. \end{aligned}$$

Let  $a_y^*$  be the function  $a_y$  defined in theorem 3 corresponding to  $f^*$ . Then,

$$a_y^*(e, x) = c_y(e, x) + c_m(e, x)\mathbb{E}\beta_y^*(e, \mathbf{v}, x)$$

$$= \begin{cases} c_y(e, x) + c_m(e, x)\mu_m(x), & y = (1, 0), \\ c_y(e, x) + c_m(e, x)\{1 - \mu_m(x)\}, & y = (0, 1), \\ c_y(e, x), & y \in \{(0, 0), (1, 1)\}, \end{cases} \quad (54)$$

from which the formula for  $\text{Cov}\{a_y^*(e, x), \bar{a}_{y^*}^*(e, x^*)\}$  follows.

## D. Solving the minmax problem

For some constants  $c_{\sigma 1}, c_{\sigma 2}$ , let  $\kappa_\sigma(v) = \epsilon_\sigma^{-1}(c_{\sigma 1}v + c_{\sigma 2})$ ,  $\kappa_{\sigma 0} = \kappa_\sigma(0)$ ,  $\kappa_{\sigma 1} = \kappa_\sigma(1)$ . By construction  $c_{\sigma 2} = \epsilon_\sigma(\kappa_{\sigma 0})$  and  $c_{\sigma 1} = \epsilon_\sigma(\kappa_{\sigma 1}) - \epsilon_\sigma(\kappa_{\sigma 0})$ . For  $j = 0, 1, 2$ , define  $I_{\sigma j} = \int_0^1 \kappa_\sigma^j(t) dt$ . By substitution,  $I_{\sigma j} = \int_{\kappa_{\sigma 0}}^{\kappa_{\sigma 1}} t^j \exp(-\sigma t^2) dt / c_{\sigma 1}$ .

### D.1 Algebraic lemmas:

**Lemma 2.**  $\kappa'_\sigma(v) = c_{\sigma 1} \exp\{\sigma \kappa_\sigma^2(v)\}$ .

*Proof.* Trivial. □

**Lemma 3.**  $\forall v \in [0, 1]$ ,  $\int_v^1 \kappa_\sigma(t) dt = [\exp\{-\sigma \kappa_\sigma^2(v)\} - \exp\{-\sigma \kappa_\sigma^2(1)\}] / (2\sigma c_{\sigma 1})$ .

*Proof.* Substitute  $s = \kappa_\sigma(t)$  and integrate. □

**D.2 Proof of Theorem 8:** Since  $e$  is irrelevant, we drop it from the notation. Further, use the shorthand  $p = p(\cdot, x)$  and  $p^* = p(\cdot, x^*)$ . Similarly, we will write  $\mu_m, \mu_m^*$  for  $\mu_m(x)$  and  $\mu_m(x^*)$ , respectively. So, we consider

$$\begin{aligned} & \max_{p, p^*} \text{or min} \int_0^1 p(v)p^*(v) dv - \mu_m \mu_m^* \\ & \text{subject to} \left\{ \begin{array}{l} \int_0^1 p(v) dv = \mu_m, \\ \int_0^1 p^*(v) dv = \mu_m^*, \\ \int_0^1 \log a(v) dv + \int_0^1 \log a^*(v) dv \geq \mathcal{E}^*, \\ p(1), p^*(1) \leq 1, \\ p(0), p^*(0) \geq 0, \end{array} \right. \quad (55) \end{aligned}$$

where  $a$  and  $a^*$  denote  $\partial_v p$  and  $\partial_v p^*$ , respectively:  $a(v) > 0$  and  $a^*(v) > 0$  are implicit by the entropy inequality constraint. Noting that  $p(v) = \int_0^v a(t) dt$ , we have

$$\int_0^1 p(v)p^*(v) dv = \int_0^1 \int_0^v p^*(v)a(t) dt dv = \int_0^1 a(v) \int_v^1 p^*(t) dt dv,$$

$$\int_0^1 p(v) dv = \int_0^1 \int_0^v a(t) dt dv = \int_0^1 \int_t^1 a(t) dv dt = \int_0^1 a(v)(1-v) dv,$$

etcetera. Using this formulation, we solve (55) with respect to  $a$  and  $a^*$ , from which we obtain the first order conditions

$$\forall v : \begin{cases} \int_v^1 p^*(v) dv + \frac{\lambda}{a(v)} - \psi(1-v) - \omega = 0, \\ \int_v^1 p(v) dv + \frac{\lambda}{a^*(v)} - \psi^*(1-v) - \omega^* = 0, \end{cases} \quad (56)$$

where  $\lambda, \psi, \omega$  are Lagrange multipliers for the entropy constraint,  $\int_0^1 p(v) dv = \mu_m$ , and  $\int_0^1 a(v) dv \leq 1$ , respectively:  $\psi^*$  and  $\omega^*$  are the corresponding multipliers for  $a^*$ . The conditions in (56) are in addition to the constraints in (55). Note that in the maximization problem,  $\lambda, \omega, \omega^* \geq 0$ , and in the minimization problem  $\lambda, \omega, \omega^* \leq 0$ .

Differentiating (56) yields differential equations, from which we deduce that there are solutions of the form

$$\begin{cases} p_\sigma(v) = \mu_m + b_\sigma \{\kappa_\sigma(v) - I_{\sigma 1}\}, \\ p_\sigma^*(v) = \mu_m^* + b_\sigma^* \{\kappa_\sigma(v) - I_{\sigma 1}\}. \end{cases} \quad (57)$$

Indeed, with (57), the multipliers satisfy

$$\begin{aligned} \lambda = -\frac{b_\sigma b_\sigma^*}{2\sigma}, \quad \psi = \mu_m^* - b_\sigma^* I_{\sigma 1}, \quad \omega = -\frac{b_\sigma^* \exp(-\sigma \kappa_{\sigma 1}^2)}{2\sigma c_{\sigma 1}}, \\ \psi^* = \mu_m - b_\sigma I_{\sigma 1}, \quad \omega^* = -\frac{b_\sigma \exp(-\sigma \kappa_{\sigma 0}^2)}{2\sigma c_{\sigma 1}}, \end{aligned}$$

where we note that  $\sigma = \pm 1$  determines the signs of  $\lambda, \omega, \omega^*$  in accordance with the Kuhn Tucker requirements. So, all that remains to be done is to determine the constants  $b_1, b_{-1}, c_{11}, c_{-11}, c_{12}, c_{-12}$ . Below we discuss how to determine  $c_{\sigma 1}, c_{\sigma 2}$  or equivalently,  $\kappa_{\sigma 0}, \kappa_{\sigma 1}$ .

Now, note that if  $\lambda = 0$ , then either  $\forall v : a(v) = 0$  or  $\forall v : a^*(v) = 0$ , which would violate the entropy inequality constraint since  $\mathcal{E}^*$  was assumed to be finite. Since the entropy inequality constraint must be binding, we have

$$\log(b_\sigma b_\sigma^*) + 2 \log c_{\sigma 1} + 2\sigma I_{\sigma 2} = \mathcal{E}^*. \quad (58)$$

Further,  $\lambda \neq 0$  implies that  $\omega \neq 0$  and  $\omega^* \neq 0$ . Therefore, the boundary conditions at  $v = 1$  must be binding, also:

$$p_\sigma(1) = \mu_m + b_\sigma(\kappa_{\sigma 1} - I_{\sigma 1}) = 1 \quad \text{and} \quad p_\sigma^*(1) = \mu_m^* + b_\sigma^*(\kappa_{\sigma 1} - I_{\sigma 1}) = 1. \quad (59)$$

By combining (58) and (59), we obtain a relationship between  $\kappa_{\sigma 0}, \kappa_{\sigma 1}$ : i.e.

$$\log(1 - \mu_m)(1 - \mu_m^*) + 2 \log\{c_{\sigma 1}/(\kappa_{\sigma 1} - I_{\sigma 1})\} + 2\sigma I_{\sigma 2} = \mathcal{E}^*. \quad (60)$$

Note also that  $\kappa_{\sigma 0}$  and  $\kappa_{\sigma 1}$  must satisfy the boundary condition at  $v = 0$ . So,

$$\begin{cases} p_{\sigma}(0) = \mu_m - (1 - \mu_m)(I_{\sigma 1} - \kappa_{\sigma 0})/(\kappa_{\sigma 1} - I_{\sigma 1}) \geq 0, \\ p_{\sigma}^*(0) = \mu_m^* - (1 - \mu_m^*)(I_{\sigma 1} - \kappa_{\sigma 0})/(\kappa_{\sigma 1} - I_{\sigma 1}) \geq 0. \end{cases} \quad (61)$$

Finally, the objective function in (55) that corresponds to  $p_{\sigma}, p_{\sigma}^*$  can be written as

$$\int_0^1 p_{\sigma}(v)p_{\sigma}^*(v) dv - \mu_m\mu_m^* = b_{\sigma}b_{\sigma}^*(I_{\sigma 2} - I_{\sigma 1}^2) = (I_{\sigma 2} - I_{\sigma 1}^2) \exp(\mathcal{E}^* - 2\sigma I_{\sigma 2})/c_{\sigma 1}^2.$$

Therefore, we can determine  $\kappa_{\sigma 0}, \kappa_{\sigma 1}$  by maximizing ( $\sigma = -1$ ) or minimizing ( $\sigma = 1$ )

$$(I_{\sigma 2} - I_{\sigma 1}^2) \exp(\mathcal{E}^* - 2\sigma I_{\sigma 2})/c_{\sigma 1}^2$$

subject to (60) and (61). □

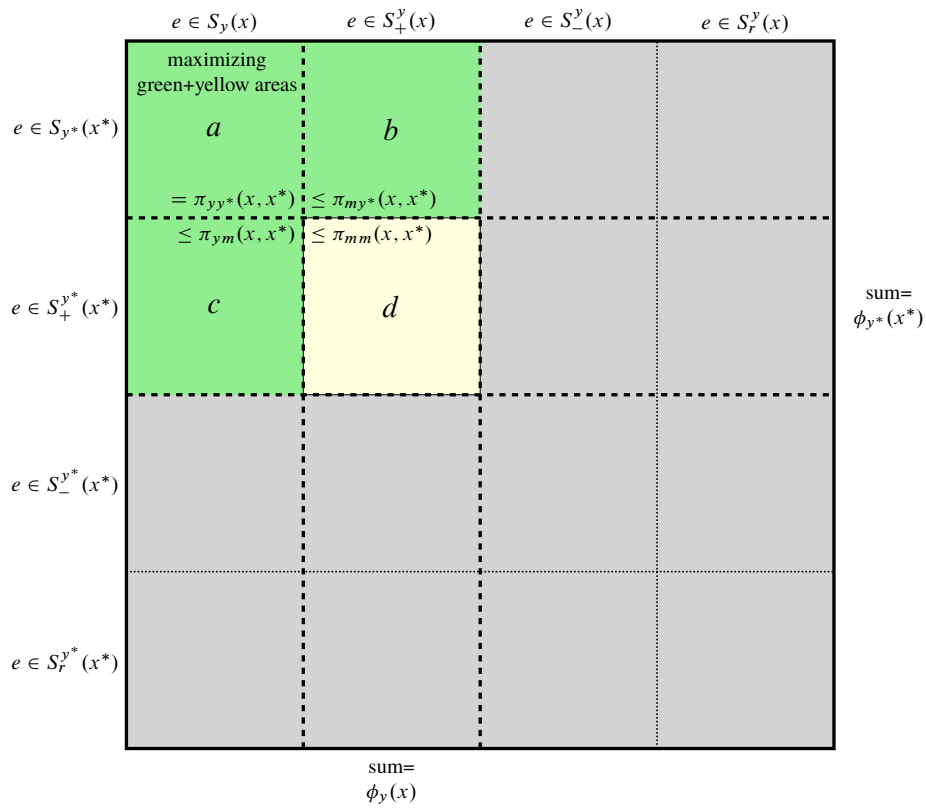
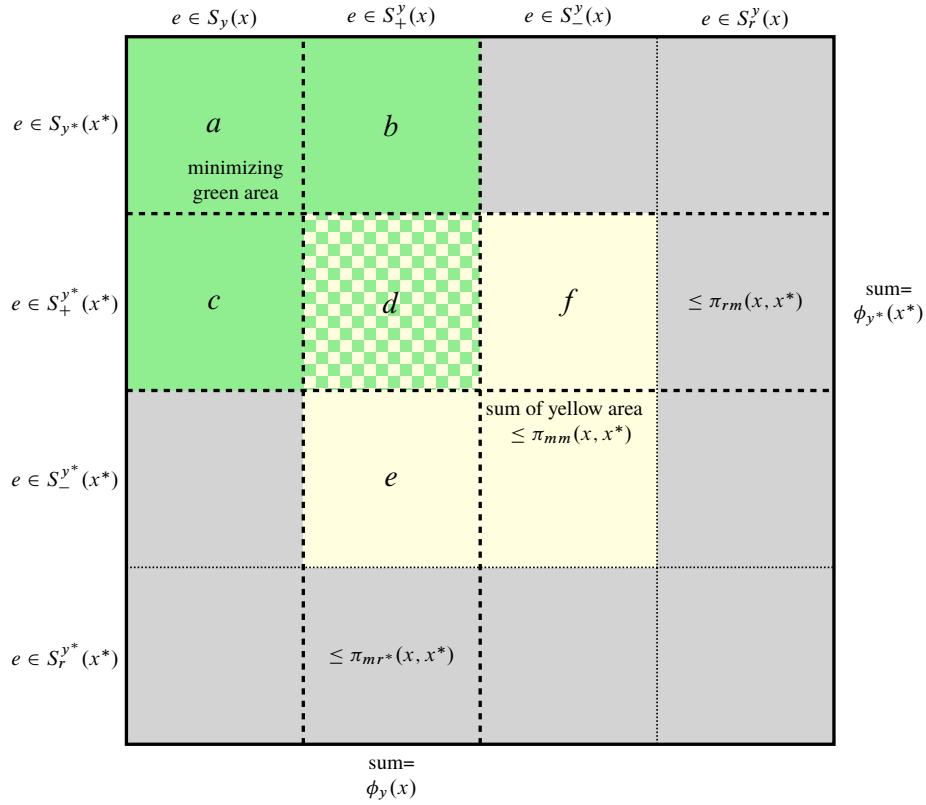


Figure 16: Finding the lower bound (top) and the upper bound (bottom)

## References

- Aguirregabiria, Victor and Pedro Mira (2007). "Sequential estimation of dynamic discrete games". *Econometrica* 75.1, pp. 1–53.
- (2018). "Identification of games of incomplete information with multiple equilibria and unobserved heterogeneity". *University of Toronto Department of Economics Working Paper*.
- Aradillas-Lopez, Andres (2011). "Nonparametric probability bounds for Nash equilibrium actions in a simultaneous discrete game". *Quantitative Economics* 2.2, pp. 135–171.
- Aradillas-López, Andres and Elie Tamer (2008). "The identification power of equilibrium in simple games". *Journal of Business & Economic Statistics* 26.3, pp. 261–283.
- Athey, Susan and Guido W Imbens (2007). "Discrete choice models with multiple unobserved choice characteristics". *International Economic Review* 48.4, pp. 1159–1192.
- Aumann, Robert J (1961). "Borel structures for function spaces". *Illinois Journal of Mathematics* 5.4, pp. 614–630.
- Bajari, Patrick, C Lanier Benkard, and Jonathan Levin (2007). "Estimating dynamic models of imperfect competition". *Econometrica* 75.5, pp. 1331–1370.
- Bajari, Patrick, Jinyong Hahn, Han Hong, and Geert Ridder (2011). "A note on semiparametric estimation of finite mixtures of discrete choice models with application to game theoretic models". *International Economic Review* 52.3, pp. 807–824.
- Bajari, Patrick, Han Hong, and Stephen P Ryan (2010). "Identification and estimation of a discrete game of complete information". *Econometrica* 78.5, pp. 1529–1568.
- Baker, Alan (2013). "Simplicity". *Stanford Encyclopedia of Philosophy*. Ed. by Edward N. Zalta. Fall 2013.
- Bjorn, Paul and Quang Vuong (1984). *Simultaneous equations models for dummy endogenous variables: a game theoretic formulation with an application to labor force participation*. Tech. rep. CalTech.
- Bresnahan, Timothy F and Peter C Reiss (1991a). "Empirical models of discrete games". *Journal of Econometrics* 48.1, pp. 57–81.
- (1991b). "Entry and competition in concentrated markets". *Journal of Political Economy*, pp. 977–1009.
- Briesch, Richard A, Pradeep K Chintagunta, and Rosa L Matzkin (2012). "Nonparametric discrete choice models with unobserved heterogeneity". *Journal of Business & Economic Statistics*.
- Bulow, Jeremy I, John D Geanakoplos, and Paul D Klemperer (1985). "Multimarket oligopoly: Strategic substitutes and complements". *Journal of Political economy* 93.3, pp. 488–511.
- Cass, David and Karl Shell (1983). "Do sunspots matter?" *Journal of Political Economy* 91.2, pp. 193–227.
- Chen, Xiaohong and Demian Pouzo (2012). "Estimation of nonparametric conditional moment models with possibly nonsmooth generalized residuals". *Econometrica* 80.1, pp. 277–321.
- Ciliberto, Federico and Elie Tamer (2009). "Market structure and multiple equilibria in airline markets". *Econometrica* 77.6, pp. 1791–1828.
- Cover, Thomas M and Joy A Thomas (2012). *Elements of information theory*. John Wiley & Sons.

- Foster, Dean and Peyton Young (1990). "Stochastic evolutionary game dynamics\*". *Theoretical population biology* 38.2, pp. 219–232.
- Fréchet, Maurice (1935). "Généralisations du théorème des probabilités totales". *Fundamenta Mathematica* 25, pp. 379–387.
- Galichon, Alfred and Marc Henry (2011). "Set identification in models with multiple equilibria". *Review of Economic Studies*, pp. 1264–1298.
- Golan, Amos (2017). *Foundations of info-metrics: modeling and inference with imperfect information*. Oxford University Press.
- Golan, Amos, George G Judge, and Douglas Miller (1996). *Maximum entropy econometrics: Robust estimation with limited data*. Wiley New York.
- Grieco, Paul LE (2014). "Discrete games with flexible information structures: An application to local grocery markets". *Rand Journal of Economics* 45.2, pp. 303–340.
- Grünwald, Peter D and A Philip Dawid (2004). "Game theory, maximum entropy, minimum discrepancy and robust Bayesian decision theory". *Annals of Statistics*, pp. 1367–1433.
- Haile, Philip A and Elie Tamer (2003). "Inference with an incomplete model of English auctions". *Journal of Political Economy* 111.1, pp. 1–51.
- Harremoës, Peter and Flemming Topsøe (2001). "Maximum entropy fundamentals". *Entropy* 3.3, pp. 191–226.
- Harsanyi, John C and Reinhard Selten (1988). "A general theory of equilibrium selection in games". *MIT Press Books* 1.
- Jaynes, E.T. (1957a). "Information theory and statistical mechanics". *Physical Review* 106, pp. 620–630.
- (1957b). "Information theory and statistical mechanics II". *Physical Review* 108, pp. 171–190.
- Jia, Panle (2008). "What happens when Wal-Mart comes to town: An empirical analysis of the discount retailing industry". *Econometrica*, pp. 1263–1316.
- Jun, Sung Jae and Joris Pinkse (2017). *Point decisions in partially identified auction models: an information–theoretic approach*. Tech. rep. The Pennsylvania State University.
- Kalai, Adam and Ehud Kalai (2012). "Cooperation in strategic games revisited". *Quarterly Journal of Economics*, qjs074.
- Kandori, Michihiro, George J Mailath, and Rafael Rob (1993). "Learning, mutation, and long run equilibria in games". *Econometrica*, pp. 29–56.
- Kashaev, Nail (2015). *Testing for Nash behavior in entry games with complete information*. Tech. rep. Penn State.
- Kashaev, Nail and Bruno Salcedo (2015). *Identification of solution concepts for discrete semi-parametric games with complete information*. Tech. rep. Penn State.
- Kasy, Maximilian (2011). "Identification in triangular systems using control functions". *Econometric Theory* 27.03, pp. 663–671.
- Kitagawa, Toru (2012). *Estimation and inference for set-identified parameters using posterior lower probability*. Tech. rep. University College London.
- Kline, Brendan (2015). "Identification of complete information games". *Journal of Econometrics* 189.1, pp. 117–131.
- Kline, Brendan and Elie Tamer (2012). "Bounds for best response functions in binary games". *Journal of Econometrics* 166.1, pp. 92–105.

- Kooreman, Peter (1994). "Estimation of econometric models of some discrete games". *Journal of Applied Econometrics* 9.3, pp. 255–268.
- Liu, Nianqing, Quang Vuong, and Haiqing Xu (2013). *Rationalization and identification of discrete games with correlated types*. Tech. rep. University of Texas.
- Magnolfi, Lorenzo and Camilla Roncoroni (2016). *Estimation of discrete games with weak assumptions on information*. Tech. rep. Yale.
- Manski, Charles (2015). *Interpreting point predictions*. Tech. rep. Northwestern University.
- Murray, Iain and Edward Snelson (2006). "A pragmatic Bayesian approach to predictive uncertainty". *Machine Learning Challenges. Evaluating Predictive Uncertainty, Visual Object Classification, and Recognising Textual Entailment*. Springer, pp. 33–40.
- Pakes, Ariel, Michael Ostrovsky, and Steven Berry (2007). "Simple estimators for the parameters of discrete dynamic games (with entry/exit examples)". *The Rand Journal of Economics* 38.2, pp. 373–399.
- Porter, Robert H. and J. Douglas Zona (1999). "Ohio school milk markets: an analysis of bidding". *RAND Journal* 30.2, pp. 263–288. ISSN: 07416261. URL: <http://www.jstor.org/stable/2556080>.
- Reguant, Mar (2016). *Bounding equilibria in counterfactual analysis*. Tech. rep. Northwestern University.
- Schennach, Susanne M (2014). "Entropic latent variable integration via simulation". *Econometrica* 82.1, pp. 345–385.
- Seim, Katja (2006). "An empirical model of firm entry with endogenous product-type choices". *Rand Journal of Economics* 37.3, pp. 619–640.
- Shore, John and Rodney Johnson (1980). "Axiomatic derivation of the principle of maximum entropy and the principle of minimum cross-entropy". *IEEE Transactions on information theory* 26.1, pp. 26–37.
- Soetevent, Adriaan R and Peter Kooreman (2007). "A discrete-choice model with social interactions: with an application to high school teen behavior". *Journal of Applied Econometrics* 22.3, pp. 599–624.
- Song, Kyungchul (2014). "Point decisions for interval-identified parameters". *Econometric Theory* 30.02, pp. 334–356.
- Tamer, Elie (2003). "Incomplete simultaneous discrete response model with multiple equilibria". *Review of Economic Studies* 70.1, pp. 147–165.
- Teh, Yee Whye (2010). "Dirichlet process". *Encyclopedia of machine learning*. Springer, pp. 280–287.
- Topsøe, Flemming (1979). "Information-theoretical optimization techniques". *Kybernetika* 15.1, pp. 8–27.
- van Damme, Eric (1995). *Equilibrium selection in team games*. Tech. rep. Tilburg University.
- Xu, Haiqing (2014). "Estimation of discrete games with correlated types". *The Econometrics Journal* 17.3, pp. 241–270.