

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

SUNG JAE JUN[†] AND JORIS PINKSE[‡]

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

August 9, 2016

We study the problem of counterfactual prediction in discrete decision games with complete information, pure strategies, and Nash equilibria. We show that the presence of multiple equilibria poses unique challenges for the problem of counterfactual prediction even if the payoff structure is known in its entirety. We show that multiple types of counterfactuals can be defined and that the prediction probabilities are not generally point-identified. We establish the sharp identified bounds of the prediction probabilities. We further propose, compare, and contrast various decision methods for the purpose of producing a point prediction, namely midpoint prediction, a decision-theoretic possibility using a Dirichlet-based prior, and a maximum-entropy approach. 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.

*We thank the Human Capital Foundation for their support of CAPCP. We thank Paul Grieco, Ken Hendricks, Andrés Aradillas-López, Joachim Freyberger, Ron Gallant, Bruce Hansen, Nail Kashaev, Francesca Molinari, Ulrich Müller, Jack Porter, Mar Reguant, Bruno Salcedo, Xiaoxia Shi, Guofu Tan, conference participants at the Penn State-Cornell econometrics conference at Ithaca, and seminar participants at Texas A&M, the University of Wisconsin at Madison, and Syracuse University for helpful comments and suggestions.

[†]619 Kern Bldg, University Park, PA 16802, sjun@psu.edu

[‡]joris@psu.edu

1. Introduction

We study the problem of counterfactual prediction in discrete decision games with complete information, pure strategies, and Nash equilibria. 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 problems associated with identification and estimation of payoffs in such games, which are mostly a consequence of the existence of multiple equilibria, have been well-studied and it is well-known how to address them (see e.g. Tamer, 2003; Kline and Tamer, 2012; Kline, 2015). However, the presence of multiple equilibria poses unique challenges for the problem of counterfactual prediction even if the payoff structure is known in its entirety. It is the counterfactual prediction problem that is the subject of this paper. What makes the counterfactual prediction problem studied here especially interesting from an econometric perspective is the *incompleteness* of the model because it features a partially identified function ρ that characterizes equilibrium selection.

To facilitate both the analysis and the exposition, we focus on the simplest possible (and somewhat hackneyed) case that has all relevant features: a single-shot game with two players, binary decisions, and strategic substitutability (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. 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.

Player payoffs are functions of unobservables \mathbf{e} and observables \mathbf{x} , both of which are known to the players.¹ The observables \mathbf{x} are exogenous in the sense that they are independent of all unobservables in the model. For some combinations of \mathbf{e} , \mathbf{x} there exists only a single Nash equilibrium in pure strategies and hence a single value y of the outcome variables \mathbf{y} . Because of strategic substitutability, for other combinations of \mathbf{e} , \mathbf{x} there exist multiple (in our case two) equilibria: (1, 0) and (0, 1). The model is hence ‘incomplete’ (Tamer, 2003) in the sense that the same values of the payoff variables \mathbf{e} , \mathbf{x} can lead to different outcomes. For given $\mathbf{x} = x$, let $S_m(x)$ be the ‘multiplicity region,’ i.e. the collection of \mathbf{e} -values for which there exist two equilibria, and let $S_y(x)$ the region of \mathbf{e} -values for which $\mathbf{y} = y$ is the unique equilibrium outcome. The reasons why a particular outcome arises in the multiplicity region is unknown to us. Since identification of the payoff structure does not require knowledge of such reasons (e.g. Tamer, 2003; Kline, 2015), we will not seek to *specify* players’ behavior in the multiplicity region, so we will *not* offer an equilibrium refinement. Instead, for an unknown function ρ , we *characterize* the behavior of players by a probability $\mathbf{p} = \rho(\mathbf{e}, \mathbf{v}, \mathbf{x})$ for which

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

where \mathbf{v} represents potential unobserved ‘market’ heterogeneity (in addition to the payoff shifters \mathbf{e} , \mathbf{x}) that can affect the probability of (1, 0) being selected under multiplicity. The only restriction

¹We use bold typeface to denote random variables.

that we impose on p is that it is *weakly* monotonic in v , which is without loss of generality.²

The characterization using p does not impose assumptions because there always exists a probability that a particular equilibrium is reached, albeit that it limits the scope of counterfactual experiments that can be conducted in a sense that is discussed in section 2. 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 more detailed discussion. We know of no other papers that allow for heterogeneity like v in the present context: its role is somewhat reminiscent of that of sunspots (Cass and Shell, 1983) in that v does not affect the payoffs but *can* affect the selection probability under multiplicity.

Now, 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 will be realized. Thus, 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 $e \notin S_m(x)$, there is only one equilibrium and hence the values of u, v are then irrelevant. For $e \in S_m(x)$, however, all four variables (e, u, v, x) are germane.

Note that we have two random variables, u and v , that together represent unobserved market heterogeneity not affecting payoffs: v affects the *probability* p that a particular equilibrium is reached in the multiplicity region and a combination of u and p produces the actual outcome. The reason for having both u and v , then, is that we want to allow for counterfactuals in which some, but not necessarily all, of the market heterogeneity is fixed.

We now turn our attention to the issue of counterfactual prediction. We denote the counterfactual outcomes of x, y by x^*, y^* , respectively, and are interested in predicting y^* .³ Since y^* is categorical, we focus on the conditional probabilities of the counterfactual outcome given the observables, i.e.

$$\mathbb{P}(y^* = y^* \mid x^* = x^*, x = x, y = y), \quad (2)$$

where $y, y^* \in \mathcal{S} = \{(0, 0), (0, 1), (1, 0), (1, 1)\}$. We also condition on x, y in (2) since their values contain information about the values of e, u, 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, $x = x$ is replaced by a new value $x^* = x^*$, where $x^* = x$ is allowed. The definition of the counterfactual outcome of interest y^* will further depend on how the unobservables e, u, v are treated. In the simplest (and least interesting) case, e, u, v are replaced with independent copies e^*, u^*, v^* (i.e. e, u, v are redrawn), in which case $y^* = y(e^*, u^*, v^*, x^*)$. Then, the *regression prediction*

$$q(y^* \mid x^*, x, y) = \tilde{\pi}_{y^*}(x^*) = \mathbb{P}(y = y^* \mid x = x^*) \quad (3)$$

is the same as (2), which is hence conveniently point-identified. However, if at least one of the unobservables e, u, v is not redrawn (i.e. is the same in the counterfactual), then the regression prediction (3) differs from the counterfactual prediction probability (2). In the main text, we focus on two scenarios that we believe are representative, i.e. fixing v and fixing both e and v : we provide results for the leading cases in sections 2 and 3 and for other possibilities in section 3.3 and appendix B. We denote the prediction probability (2) in the leading cases by

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

²In contexts with more than one multiplicity region, vector-valued p may be needed, but that is not an issue here.

³The type of counterfactual we consider in this paper is thus more commonly found in e.g. the treatment effects literature (Heckman, 2005) than in e.g. the industrial organization literature (Reguant, 2016).

respectively and 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 q_v, q_{ev} , 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 for q_v , the counterfactual probability of having a monopoly collapses to the corresponding regression prediction and is hence point-identified, although q_v itself is not point-identified.⁴ However, for q_{ev} , even the probability of monopoly is only partially identified.⁵

There are several solutions to our partial identification problem. One instinctive 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 *identified set* of counterfactual probability values. We determine the sharp identified sets for both $q_v(y^* | x^*, x, y)$ and $q_{ev}(y^* | x^*, x, y)$ in section 3: the bounds for other cases are provided in appendix B. 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. Second, 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. Third 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 6. In sum, it is useful to have point predictions *in addition to* bounds.

Paul: would like a distribution over outcomes (perhaps we can do the entropy loss thing I mentioned)

One intuitive point prediction choice is the midpoint θ_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 ranked⁶ 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 6.1, there does not generally exist a single function p that can generate midpoint predictions for different counterfactuals.

Another possibility is to use a decision-theoretic approach. This entails defining a loss function ℓ which measures the distance between the infeasible prediction (which depends on p) and the ‘decision’ (which does *not* depend on p). The integrated loss function, namely the risk, is then aggregated into an ‘average risk’ by integrating over a parameter space consisting of p -functions that are consistent with the distribution of observables. We can determine the function \tilde{p} that minimizes

⁴Point identification of the number of competitors does not extend to the case with more than two choices and/or players.

⁵We thank Ken Hendricks for suggesting this case.

⁶The order can be ascending or descending, depending on the values of y, y^* .

the average risk and use it to compute counterfactual probabilities.

We investigate the decision–theoretic approach in section 5. 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 5 or something else, the decision–theoretic approach does not solve the fundamental problem at hand. Indeed, 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 decision–theoretic 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.⁷ 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 5, the decision–theoretic approach generates predictions that are similar to our preferred method, *maximum entropy*, which is discussed in section 4.

With the maximum entropy method of the information theory literature (Jaynes, 1957a,b; 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).⁸ There are numerous philosophical justifications for the application of Occam’s razor; we refer the reader to Baker (2013).

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.⁹

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

⁷This objection can be removed by focusing on worst–case risk (minmax), but minmax takes us back to the midpoint prediction, which we have already discussed.

⁸With maximum entropy, the closest choice is unique up to trivial deviations.

⁹There 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.

maximum entropy on a monotonic transformation of \mathbf{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 (e.g. Cover and Thomas, 2012) and since the random variable at the center of the problem is the selection probability \mathbf{p} , it is more natural to deal with the distribution of \mathbf{p} than that of e.g. \mathbf{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. Finally, as noted above, for the choices made in section 5, the Dirichlet and maximum entropy approaches appear to generate similar predictions while the maximum entropy approach is considerably easier to implement.

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 ρ and construct corresponding identified sets in the form of sharp bounds; see section 3. We further derive bounds if 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 5 and a new point prediction method based on the maximum entropy concept from the information theory literature in section 4.¹⁰ Further, we demonstrate the virtues of the maximum entropy approach (and to a lesser extent the Dirichlet approach) in a number of examples in section 6. Finally, analogs of our 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

Consider a standard two player binary decision game with strategic substitutability with complete information and pure strategies, i.e.

$$\begin{cases} y_1 = \mathbb{1}\{\tau_1(\mathbf{x}, y_2) + e_1 \geq 0\}, \\ y_2 = \mathbb{1}\{\tau_2(\mathbf{x}, y_1) + e_2 \geq 0\}, \end{cases} \quad (4)$$

where \mathbf{x} is a vector of exogenous covariates, $\mathbf{e} = (e_1, e_2)$ are errors that are independent of \mathbf{x} , and $\tau_i(x, 1) \leq \tau_i(x, 0)$ for all x and $i = 1, 2$. The model in (4) is *incomplete* in that \mathbf{x} and \mathbf{e} do not necessarily determine a unique outcome $\mathbf{y} = (y_1, y_2)$ due to the possible presence of multiple Nash equilibria (see e.g. Bresnahan and Reiss, 1991a; Tamer, 2003).

The properties of the model in (4) have been studied extensively and conditions under which τ_j and the distribution of (\mathbf{e}, \mathbf{x}) are identified are well–understood (e.g. Tamer, 2003). From hereon, we therefore take the functions τ_1, τ_2 and the distribution of (\mathbf{e}, \mathbf{x}) as given.

¹⁰Others 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.

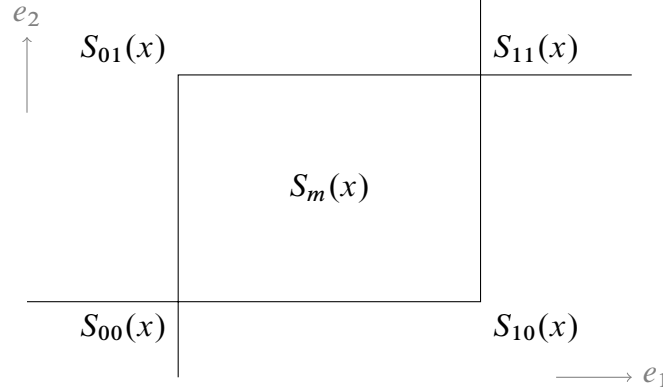


Figure 1: Regions of e -values corresponding to either a specific equilibrium and a multiplicity region.

We follow Tamer (2003) and split the (e_1, e_2) -space (\mathbb{R}^2) into five regions: four regions ($S_{00}(x)$, $S_{01}(x)$, $S_{10}(x)$, $S_{11}(x)$) in which a unique Nash equilibrium arises plus a *multiplicity region* $S_m(x)$, in which either $(1, 0)$ or $(0, 1)$ is the outcome. See figure 1.

Since we are taking the functions τ_j 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.

Recall from (1) that for some random variable $p = p(e, v, x)$ and all p, x , if $e \in S_m(x)$ then

$$\mathbb{P}\{y = (1, 0) \mid e = e, p = p, x = x\} = p, \quad (5)$$

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 .

The relationship in (5) is a characterization, not an assumption. In other words, there always exists such a function p and such a random variable v . Equation (5) says that in a given market, under multiplicity $(1, 0)$ is selected with probability p , where p can depend on payoff shifters and on (other) unobserved market heterogeneity v . In view of the (weak) monotonicity of p in v , v reflects factors (unobservable to the econometrician) that affect the probability of choosing one equilibrium over another over and above what is explained by the payoff shifters. However, (5) is not necessarily a description of players' behavior in the multiplicity region. In particular, it does *not* assume that players randomize in their choice of equilibrium: it merely says that from the econometrician's perspective this choice *can* be viewed as random (p is allowed to be binary).

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) and also nests the structures in Bjorn and Vuong (1984), where $p(e, v, x) = \mu_m(x)$, and Jia (2008), where $p(e, v, x) = 1$. But, unlike Grieco (2014), we allow for an unobservable v in p and p can be flat in e, v . Our specification is a *characterization* instead of an *assumption* and it is consistent with *any* equilibrium refinement.¹¹ The same can be said for Grieco (2014) but, as noted in the introduction, we want to allow for the possibility that v is fixed and u changes, i.e. that some of the unobserved payoff-irrelevant heterogeneity is fixed and some changes. 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).

Although our characterization is consistent with *any* equilibrium refinement, our characterization does limit the nature of the counterfactual experiments discussed in section 3. Indeed, in section 3 we hold v (and possibly e) fixed (in our leading cases), but it is possible that an underlying structural model of 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 (5) 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, (5) does not restrict the dependence of e and p .

In view of (5) we can represent the outcome y by

$$y = y(e, u, v, x) = \begin{cases} y, & e \in S_y(x), \\ (1, 0), & e \in S_m(x) \text{ and } u \leq p(e, v, x), \\ (0, 1), & e \in S_m(x) \text{ and } u > p(e, v, x), \end{cases} \quad (6)$$

where u is uniform and independent of e, v, x . So, the outcome is determined by four random elements: e, x are unobservable and observable payoff variables, which *can* also impact p in the multiplicity region,¹² and v captures extra unobserved market heterogeneity that affects p independent of e, x . Further, u is an unobservable variable, which, together with p , determines which equilibrium is played in the multiplicity region. Note that fixing the payoff variables and v does not determine the outcome because $y(e, u, v, x)$ is still random: only the probability of which outcome is selected under multiplicity is determined by e, v, x , not (necessarily) the outcome itself.

In view of (6), (5) *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 (5) imposes is by itself of limited help for the prediction problem studied in this paper.

¹¹Kalai and Kalai (2012) (and references therein) provides refinements for two by two games, albeit that in Kalai and 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.

¹²We do not impose exclusion restrictions here.

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) .¹³ 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 (6) 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 counterfactual of interest is $y^* = y(e, u^*, v^*, x^*)$.¹⁴ 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 v unchanged and keeping both e and 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 e, u, v in appendix B,¹⁵ and provide a brief summary thereof in section 3.3.

For each of the counterfactuals, we wish to construct a prediction of y^* . Since y^* is a categorical variable, we obtain the distribution of y^* given observables. We define $q_c(y^* | x^*, x, y)$ to be the conditional probability that $y^* = y^*$ given $x^* = x^*, x = x, y = y$, where the subscript c indicates which of e^*, u^*, v^* are fixed. For instance,

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

where q represents the case with no constraints and q_v, q_{ev} represent the cases where (part of) unobserved market heterogeneity is kept unchanged. Please note that we use both x and y to predict y^* , because the values of both x and y contain information about p .

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

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

¹³So, 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.

¹⁴Therefore, the ‘condition’ $e^* = e$ does not mean that we are ‘conditioning.’ It means that we keep 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 ξ_1, ξ_2 are random variables with a standard exponential distribution. If ξ_2 is the same random variable as ξ_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. ξ_1, ξ_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.

¹⁵There 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_{evv} ; and predictions that are omitted because of their similarity to ones that are discussed in the paper, namely q_u and q_{eu} .

The *regression prediction* $\tilde{\pi}_{y^*}(x^*)$ is simple, but it represents the case where nothing 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. For this purpose, we define

$$\mu_m(x) = \mathbb{E}\{\mathbf{p} \mid \mathbf{e} \in S_m(x), \mathbf{x} = x\}. \quad (7)$$

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 $\pi_m(x) > 0$.¹⁶ In fact, μ_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

$$\delta_y = \mathbb{1}\{y = (1, 0)\} - \mathbb{1}\{y = (0, 1)\}. \quad (8)$$

3.1 Case 1: $\mathbf{v}^* = \mathbf{v}$: We now study identification of $q_v(y^* \mid x^*, x, y)$. 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

$$\begin{aligned} \rho(x, x^*) &= \text{Cov}\{p(\mathbf{e}, \mathbf{v}, x), p(\mathbf{e}^*, \mathbf{v}, x^*) \mid \mathbf{e} \in S_m(x), \mathbf{e}^* \in S_m(x^*)\} \\ &= \text{Cov}\{k(\mathbf{v}, x), k(\mathbf{v}, x^*)\}, \end{aligned} \quad (9)$$

where

$$k(\mathbf{v}, x) = \mathbb{E}\{\mathbf{p} \mid \mathbf{e} \in S_m(x), \mathbf{v} = \mathbf{v}, \mathbf{x} = x\} = \mathbb{E}\{p(\mathbf{e}, \mathbf{v}, x) \mid \mathbf{e} \in S_m(x)\}. \quad (10)$$

The function ρ 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 it belongs to the multiplicity region or not: 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 δ_y was defined in (8). Then, $q_v(y^* \mid x^*, x, y) = \tilde{\pi}_{y^*}(x^*) + \tilde{\rho}_{yy^*}(x, x^*)$. \square

The results of theorem 1 are intuitive. 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. 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 $\pi_m(x^*)$

¹⁶If $\pi_m(x) = 0$ then the value of $\mu_m(x)$ is immaterial.

is large) then the correction term is large. If the probability that we are in the multiplicity region is large relative to the probability that $y = y$, 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

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

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. 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, both the maximum difference between the optimal and regression predictions and the maximum attainable length of the identified set of $q_v\{(1, 0) | x^*, x, (1, 0)\}$ equal

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

which 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, the maximum length of the identified set as $\pi_m(x), \pi_m(x^*) \rightarrow 1$ is $\min\{1, \mu_m(x^*) / \mu_m(x)\} - \mu_m(x^*)$. More generally, the maximum attainable length of the identified set varies with the ratio $\pi_{10}(x) / \pi_m(x)$ as depicted in figure 2. Much the same applies to other combinations of $y, y^* \in \{(1, 0), (0, 1)\}$.

3.2 Case 2: $\mathbf{e}^* = \mathbf{e}$ and $\mathbf{v}^* = \mathbf{v}$: We now turn to the case where both \mathbf{e} and \mathbf{v} are fixed, i.e. the probability \mathbf{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 \mathbf{e} is important for prediction. For instance, if $x^* = x$ then the S -regions are unchanged and $\mathbf{p}^* = \mathbf{p}$.¹⁷ If $x^* \neq x$ then the S -regions will be different and \mathbf{p}^* can be different from \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 $\pi_m(x), \pi_m(x^*)$ 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 $S_m(x)$ and $S_m(x^*)$ becomes relevant. We now introduce further notation to facilitate our analysis.

¹⁷By ‘ S -region’ we mean one of the sets $S(\cdot)$.

트위터

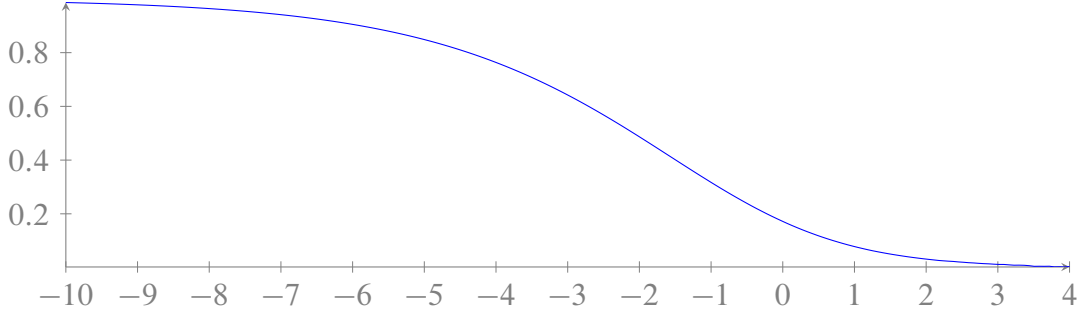


Figure 2: Maximum possible value of $\tilde{\rho}_{yy^*}(x, x^*)$ as a function of $\log\{\pi_y(x) / \pi_m(x)\}$ if $\mu_m(x^*) = \mu_m(x)$ and $y^* = y = (1, 0)$

Let $S_{my}(x, x^*) = S_m(x) \cap S_y(x^*)$, $\pi_{my}(x, x^*) = \mathbb{P}\{e \in S_{my}(x, x^*)\}$, and let $S_{mm}, S_{yy^*}, \pi_{mm}, \pi_{yy^*}$, etcetera, be analogously defined.¹⁸ We further define $S_r(x) = S_r^y(x) = \{S_y(x) \cup S_m(x)\}^c$ and $S_{r^*}(x) = S_r^{y^*}(x) = \{S_{y^*}(x) \cup S_m(x)\}^c$. In order to determine the bounds on q_{ev} , it matters where the S -regions intersect. Let

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

which is similar to (9) but uses $e = e^*$ in both p functions, whereas in (9) e^* was an independent copy of e . Indeed, $\rho(x, x^*)$ serves the same role as in theorem 1: to correct for the information contained in x, y about the value of p^* .

But we need another correction term to address the fact that x, y contain information about the $S(x^*)$ -region to which 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 $c_y(e, x) = \mathbb{1}\{e \in S_y(x)\}$, $c_m(x) = \mathbb{1}\{e \in S_m(x)\}$, and

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

Theorem 3. Let $\bar{a}_y(e, x) = c_y(e, x) + c_m(x)\mathbb{E}b_y(e, v, x)$ and

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

Then, in view of (5), $q_{ev}(y^* \mid x^*, x, y) = \tilde{\pi}_{y^*}(x^*) + \alpha_{yy^*}(x, x^*) + \tilde{\rho}_{yy^*}(x, x^*)$. \square

The intuition for the $\tilde{\rho}$ -term in theorem 3 is similar to that in theorem 1. For the α term in theorem 3, consider first the case that $y = y^* = (0, 0)$. Then the α_{yy^*} term simply increases the

¹⁸Please note that the orders of the input arguments and the subindices are relevant here.

트위티

probability that $\mathbf{y}^* = (0, 0)$ because we know that $\mathbf{y} = (0, 0)$ implies that $\mathbf{e}_1, \mathbf{e}_2$ are both small. The intuition for the other cases are merely more complicated versions of this argument.

Before we proceed, we introduce some notation. 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

$$\begin{aligned} q_{ev}\{(1, 0) \mid x^*, x, y\} + q_{ev}\{(0, 1) \mid x^*, x, y\} \\ = \tilde{\pi}_{10}(x^*) + \tilde{\pi}_{01}(x^*) + \frac{\text{Cov}\{a_y(\mathbf{e}, x), a_{10}(\mathbf{e}, x^*) + a_{01}(\mathbf{e}, x^*)\}}{\tilde{\pi}_y(x)}. \end{aligned} \quad (11)$$

The covariance on the right hand side in (11) 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 $\rho(\mathbf{e}, \mathbf{v}, x)$. Note however that $\rho(\mathbf{e}, \mathbf{v}, x^*)$ 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: i.e. the probability of monopoly conditional on no monopoly 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.

트위티

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

$$U_{yy^*}(x, x^*) = \min. \text{ of } \begin{cases} \phi_y(x) + \phi_{y^*}(x^*), \\ \phi_y(x) + \pi_{ym}(x, x^*), \\ \phi_{y^*}(x^*) + \pi_{my^*}(x, x^*), \\ \pi_{mm}(x, x^*) + \pi_{ym}(x, x^*) + \pi_{my^*}(x, x^*), \end{cases} \quad (12)$$

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 π -values and the values of $\phi_y(x)$ and $\phi_{y^*}(x^*)$, as is illustrated in figure 3. Depending on the values of $\phi_y(x)$ and $\phi_{y^*}(x^*)$, different bounds are binding.

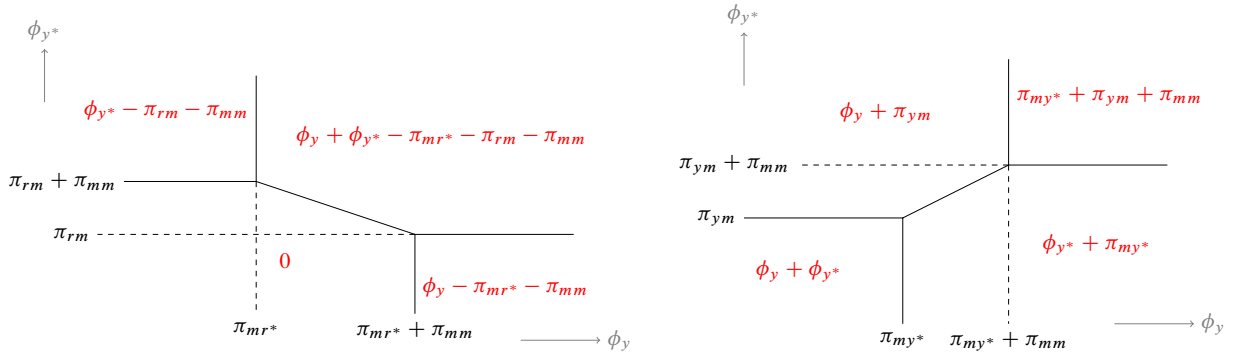


Figure 3: Regions of $(\phi_y(x), \phi_{y^*}(x^*))$ and the corresponding values of $L_{yy^*}(x, x^*)$, $U_{yy^*}(x, x^*)$.

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 fact that the formulas for the bounds in theorem 5 depend fundamentally on whether or not $x = x^*$ 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 $\rho(\mathbf{e}, \mathbf{v}, x)$ is continuous in x (which we do not assume) then $|\rho(\mathbf{e}, \mathbf{v}, x^*) - \rho(\mathbf{e}, \mathbf{v}, x)| / \|x - x^*\|$ can still be arbitrarily large. Since ρ need not be monotonic in \mathbf{e} , the covariance between \mathbf{p} and \mathbf{p}^* can be negative, unlike the variance of \mathbf{p} .

If the function ρ is flat in \mathbf{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 6.3 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

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

where $\bar{\rho}(\mathbf{e}, x) = \mathbb{E}\rho(\mathbf{e}, \mathbf{v}, x)$. 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 ρ is not necessarily monotonic in \mathbf{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.

$$\begin{cases} q(y^* | x^*, x, y) \leq q_v(y^* | x^*, x, y) \leq q_{uv}(y^* | x^*, x, y), \\ q_e(y^* | x^*, x, y) \leq q_{ev}(y^* | x^*, x, y) \leq q_{euv}(y^* | x^*, x, y), \end{cases} \quad (13)$$

which is intuitive. However, if $x \neq x^*$, as is shown in theorems 8 to 12 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} .¹⁹ Since the bounds coincide, so do the midpoint predictions. Thus, despite the fact that the q -functions can be ordered as indicated in (13), the midpoint predictions are the same in each case.²⁰ The maximum entropy solution proposed in section 4 does not have this unfortunate feature.

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,b). The entropy measures the amount of uncertainty that a probability distribution represents.²¹

In our case, this entails finding the joint density of e, p, x which maximizes the entropy: recall from (5) 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 as being 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 (7), is identified for all x , and hence it imposes another constraint on f . Thus, using all information available, we consider maximizing the *conditional entropy*²²

$$f^* = \underset{f}{\operatorname{argmin}} \iiint_0^1 f(p | e, x) \log f(p | e, x) dp f_e(e) f_x(x) de dx$$

$$\text{subject to } \begin{cases} \forall e, x : \int_0^1 f(p | e, x) dp = 1, \\ \forall x : \int_{S_m(x)} \int_0^1 pf(p | e, x) dp f_e(e) de = \mu_m(x) \pi_m(x). \end{cases} \quad (14)$$

In (14) we are minimizing minus the entropy to get rid of the minus sign in the entropy definition. The optimization problem in (14) can be solved using standard constrained optimization techniques (see e.g. Cover and Thomas, 2012). The result is provided in theorem 6 below. Let

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

Theorem 6. The solution to (14) is given by

$$f^*(p | e, x) = \begin{cases} \mathbb{1}(0 \leq p \leq 1), & e \notin S_m(x), \\ A\{p, \lambda_m(x)\} \mathbb{1}(0 \leq p \leq 1), & e \in S_m(x), \end{cases} \quad (15)$$

¹⁹We focus on the case $x \neq x^*$ because $q_{euv}(y^* | x, x, y) = \mathbb{1}(y^* = y)$.

²⁰The midpoint predictions are different if one compares fixing e and not fixing e , however.

²¹Admittedly, one could define alternative such measures; our use of entropy is motivated by it being the dominant choice in information theory.

²²The problem can be equivalently formulated by the unconditional entropy of the joint density of (e, p, x) .

Can we create 'bounds' by showing what a loss of e.g. one unit of entropy would do to the predictions?

트위터

where

$$\lambda_m(x) = \underset{\lambda}{\operatorname{argmin}}\{\mathcal{L}(\lambda) - \mu_m(x)\lambda\}, \quad A(p, \lambda) = \exp(p\lambda) / I(\lambda). \quad (16)$$

□

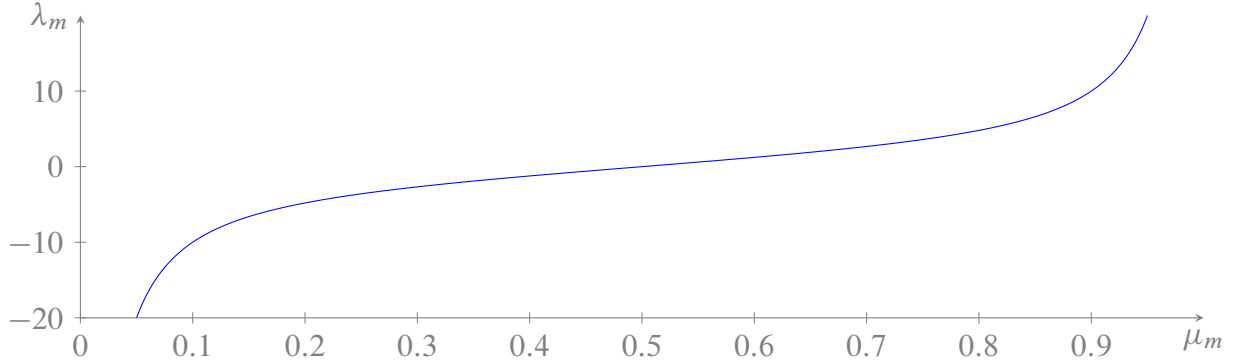


Figure 4: λ_m as a function of μ_m .

The relationship between λ_m and μ_m is depicted in figure 4. We can see from the first order condition of the minimization problem in (16) that the solution λ_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 well-behaved at zero.²³

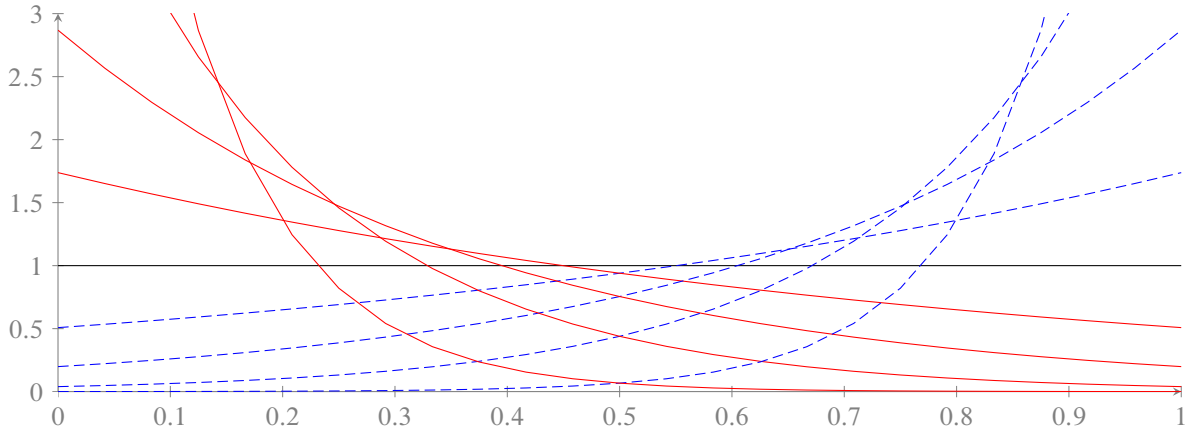


Figure 5: Conditional densities $f^*(p | e, x)$ for $e \in S_m(x)$ for values of $\mu_m = 0.1, \dots, 0.9$, with blue dashed curves corresponding to $\mu_m > 0.5$ and red solid curves to $\mu_m < 0.5$.

²³For instance, \mathcal{L}' is continuous and differentiable at zero.

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

Recall again 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 (10) 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

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

and

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

트위터

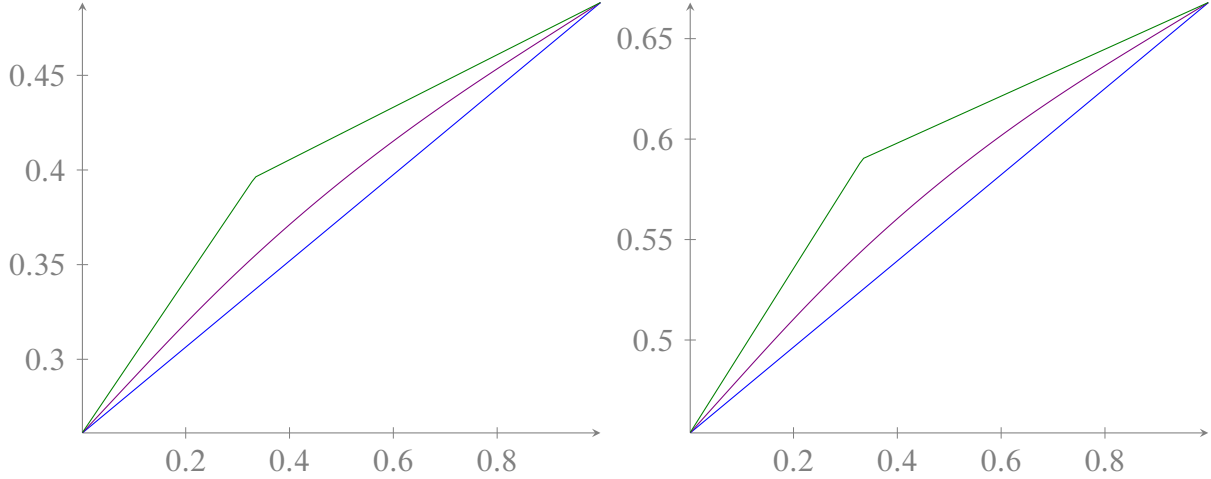


Figure 6: Maximum entropy solutions for $q.$, q_v , q_{uv} (left graph, bottom to top) and q_e , q_{ev} , q_{euv} (right graph, bottom to top) as a function of $\mu_m(x^*)$ if $y = y^* = (1, 0)$, $\mu_m(x) = 1 / 3$, e_1, e_2 independent $N(0, 1)$, shift in x value moves multiplicity region from $[-1, -1]^2$ to $[0, 2]^2$.

For one simple design, the maximum entropy predictions are depicted in figure 6. 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 6, we will compare the maximum entropy solutions with various alternatives using a number of stylized examples.

The maximum entropy solution is our preferred method to make a point decision. Its main alternative, the Dirichlet solution, is discussed in section 5. The Dirichlet solution is of intrinsic interest on its own but since its derivation is more technical, readers who are more interested in the comparison of methods than in the details of the Dirichlet solution can skip ahead to section 6.

5. Dirichlet approach

5.1 Overview: In this section we discuss an alternative possibility, for which Dirichlet processes will be used. The approach we discuss has a decision-theoretic foundation, but its implementation is considerably more complicated than maximum entropy and it does not do much to mitigate the inherent arbitrariness of the problem at hand. We will compare the results to the maximum entropy ones in section 6.

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.²⁴ Suppose further that \mathbf{G} is a draw from the prior. Then, our prediction based on the draw \mathbf{G} will be

$$T_{\mathbf{G}^{-1}}(x, x^*, y, y^*), \quad (17)$$

which is a random object. The approach considered in this section is based on the probability distribution of (17). We show in section 5.2 that the mean of (17) 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 ρ in the distribution of observables is the regression function μ_m , which we assume to be known (throughout this paper).

We now obtain the likelihood function (in terms of the parameter ρ). Note from (5) that for $e \in S_m(x)$,

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

which implies

$$\frac{1}{\pi_m(x)} \int_{S_m(x)} \mathbb{P}\{\mathbf{y} = (1, 0) \mid \mathbf{e} = e, \mathbf{x} = x; \rho\} f_e(e) de = \mu_m(x). \quad (18)$$

Equation (18) shows that the likelihood function depends on ρ only through μ_m , which is known. Therefore, since μ_m is known and fixed, the likelihood is also known and it does not depend on ρ other than via μ_m , which is a different way of saying (as we have before) that μ_m is the only identifiable object about behavior in the multiplicity region.

5.2 Average risk optimality: Our decision problem is to make the ‘best’ guess for $T_\rho(x, x^*, y, y^*)$ based on $\mathbf{x} = x, \mathbf{y} = y$ and $\mathbf{x}^* = x^*, \mathbf{y}^* = y^*$. Denote our decision by $d(x, x^*, y, y^*)$. For a given loss function ℓ , the risk (for any combination x, x^*) is then

$$\begin{aligned} R(d, \rho, x, x^*) &= \sum_{y, y^*} \ell\{d(x, x^*, y, y^*), T_\rho(x, x^*, y, y^*)\} \mathbb{P}(\mathbf{y} = y \mid \mathbf{x} = x; \rho) \mathbb{P}(\mathbf{y}^* = y^* \mid \mathbf{x}^* = x^*; \rho) \\ &= \sum_{y, y^*} \ell\{d(x, x^*, y, y^*), T_\rho(x, x^*, y, y^*)\} L(y \mid x) L(y^* \mid x^*), \end{aligned} \quad (19)$$

where the likelihood function L is known and independent of ρ (other than via μ_m), as was established in (18). However, since the prediction of interest depends on ρ , the risk in (19) depends on ρ , also. Therefore, we consider the average risk, i.e. averaging out over ρ , for which we need a prior for ρ .

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

$$\mathcal{F}_x = \left\{ \rho^{-1}(\cdot, \cdot, x) : \frac{1}{\pi_m(x)} \int_{S_m(x)} \int_0^1 \rho(e, v, x) f_e(e) de dv = \mu_m(x) \right\}. \quad (20)$$

²⁴ ρ^{-1} is the (generalized) inverse of the quantile function ρ with respect to the argument v . For a quantile function Q , its generalized inverse G^{-1} is defined by $Q^{-1}(p) = \sup\{p : Q(v) \leq p\}$, 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.

Therefore, once we have priors ω_x, ω_{x^*} on $\mathcal{F}_x, \mathcal{F}_{x^*}$, we can induce a prior ω_{xx^*} on \mathcal{F}_{xx^*} . Since p^{-1} is a (conditional) distribution function, Dirichlet processes provide a natural way of choosing ω_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 5.3.

Once we have ω_{xx^*} , the average risk will be given by

$$\begin{aligned} \bar{R}(d) &= \iiint_{\mathcal{F}_{xx^*}} R(d, p, x, x^*) d\omega_{xx^*}(p^{-1}) f_x(x) f_x(x^*) dx dx^* \\ &\geq \iint \sum_{y, y^*} \min_s \int_{\mathcal{F}_{xx^*}} \ell\{s, T_p(x, x^*, y, y^*)\} d\omega_{xx^*}(p^{-1}) L(y | x) L(y^* | x^*) f_x(x) f_x(x^*) dx dx^*, \end{aligned}$$

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 \bar{R} is given by

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

Thus, when the loss function is quadratic, $d^*(x, x^*, y, y^*)$ is simply a mean as we explained below (17).

5.3 Dirichlet prior: We now describe a way of drawing probability distributions from ω_{xx^*} , for which it suffices that we can draw probability distributions from ω_x on support \mathcal{F}_x . To keep things simple, we will restrict \mathcal{F}_x by considering elements p^{-1} that do not depend on e for all $e \in S_m(x)$. With maximum entropy, the fact that the p 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

$$\mathcal{F}_x = \left\{ p^{-1}(\cdot, x) : \int_0^1 p(v, x) dv = \mu_m(x) \right\}, \quad (21)$$

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 ω_x , and a hyperparameter $\zeta > 0$. H and ζ 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 $\beta_1^*, \beta_2^*, \dots$

Paul: why is this intro in here twice?

to the mass points τ_1, τ_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 μ_m in lieu of $\mu_m(x)$ for the remainder of this section.

Let the β_j ’s and β_j^* ’s be defined as above. Suppose that the τ_j ’s are i.i.d. $U[0, 1]$ and (unconditionally) independent of the β_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 τ_j ’s to ensure that

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

for which we will draw τ_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

$$\forall j : W_j = \beta_j \tau_j + (1 - \beta_j) W_{j+1}. \quad (22)$$

Further, the Bayes rule says that

$$f_{\tau_j | \beta_j, W_j}(\tau | \beta, \mu) = \frac{f_{W_j | \tau_j, \beta_j}(\mu | \tau, \beta)}{f_{W_j | \beta_j}(\mu | \beta)}. \quad (23)$$

Here, by (22) and the fact that the distribution of W_j does not depend on j , it follows that

$$\mathbb{P}(W_j \leq w | \tau_j = \tau_j, \beta_j = \beta_j) = \mathbb{P}\{\beta_j \tau_j + (1 - \beta_j) \tau_j W_{j+1} \leq w\} = F^\circ\left(\frac{w - \beta_j \tau_j}{1 - \beta_j}\right), \quad (24)$$

where F° is the distribution function of W_1 . We discuss further down how to compute F° efficiently.

Finally, note from (23) and (24) that $F_{\tau_j | \beta_j, W_j}(\tau | \beta, \mu) = \tilde{F}^\circ(\tau; \beta, \mu)$, where

$$\tilde{F}^\circ(t; \tilde{\beta}, \tilde{\mu}) = \frac{F^\circ\left(\frac{\tilde{\mu}}{1 - \tilde{\beta}}\right) - F^\circ\left(\frac{\tilde{\mu} - \tilde{\beta}t}{1 - \tilde{\beta}}\right)}{F^\circ\left(\frac{\tilde{\mu}}{1 - \tilde{\beta}}\right) - F^\circ\left(\frac{\tilde{\mu} - \tilde{\beta}}{1 - \tilde{\beta}}\right)}, \quad \max\left\{0, \frac{\tilde{\mu} - (1 - \tilde{\beta})}{\tilde{\beta}}\right\} \leq t \leq \min\left(1, \frac{\tilde{\mu}}{\tilde{\beta}}\right). \quad (25)$$

This motivates the following procedure.

Procedure 1. Do the following:

1. Draw $\beta_1, \beta_2, \dots \sim U(0, 1)$;

트위티

2. Draw τ_1 from $\tilde{F}^\circ(\cdot; \beta_1, \mu_m)$;
3. Draw τ_2 from $\tilde{F}^\circ\{\cdot; \beta_2, (\mu_m - \beta_1\tau_1) / (1 - \beta_1)\}$;
4. Draw τ_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 \mathbf{G} be the probability distribution with mass points $\beta_1^*, \beta_2^*, \beta_3^*, \dots$ at $\tau_1, \tau_2, \tau_3, \dots$, where β_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 μ_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 (25) 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 μ_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 τ^* is a draw from a $U(0, 1)$.

It now remains to be shown how to compute F° . From (22) 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° numerically.

6. Comparison of methods

In this section we compare the maximum entropy approach with other possibilities. Section 6.1 shows that the midpoint prediction method is inconsistent in the sense that there generally is no single function p that produces the midpoint predictions in all cases. Section 6.2 uses two simple designs to compare the sharp bounds on q_v to regression predictions, midpoint predictions, and maximum entropy predictions. Section 6.3 does the same for q_{ev} . Then, section 6.4 compares and contrasts all methods to the Dirichlet–process–based idea of section 5.

6.1 Inconsistency of the midpoint method: As promised, we now show that there generally exists no single function p that is consistent with the midpoint prediction for all values of x^* . In other words, we can construct examples in which any function p 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 p° such that the predictions based on p° are the midpoint predictions. Then, for any \bar{x} in the support of x , the predictions based on p° must satisfy

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

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

$$\mathbb{E}\nabla\{p^\circ(e, v, x) | v\} = 0,$$

from which we conclude that p° 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° is to be consistent with the midpoint predictions we must have

$$q_{ev}^\circ(y^* | x^*, x, y) = \frac{\mathbb{E}\{p^\circ(v, x)p^\circ(v, x^*)\}}{\tilde{\mu}} = \frac{1}{2}, \quad (26)$$

$$q_{ev}^\circ(y^* | x, x, y) = \frac{\mathbb{E}p^{\circ 2}(v, x)}{\tilde{\mu}} = \frac{3}{4}, \quad (27)$$

$$q_{ev}^\circ(y^* | x^*, x^*, y) = \frac{\mathbb{E}p^{\circ 2}(v, x^*)}{\tilde{\mu}} = \frac{3}{4}. \quad (28)$$

Equation (26) 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° is monotonic in v . But that conclusion contradicts either (27) or (28). Therefore, there exists no function p° 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 6.

6.2 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 6.3 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

트위터

$$\begin{aligned} \rho^*(x, x^*) &= \mathbb{V}\{h^*(v, x) \mid x = x\} \\ &= \int_0^1 p^2 A\{p, \lambda_m(x)\} dp - \left\{ \int_0^1 p A\{p, \lambda_m(x)\} dp \right\}^2 = \mathcal{L}''\{\lambda_m(x)\}, \end{aligned} \quad (29)$$

where A is defined in (16). Example 2 is based on (29).

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 7 depicts prediction probabilities as a function of $\tilde{\mu}$ in each case.

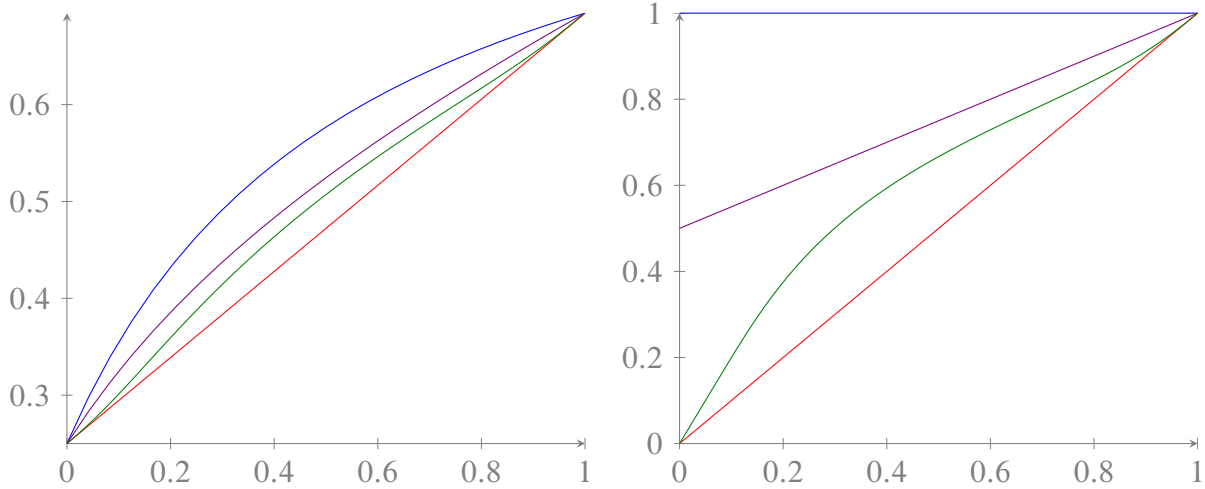


Figure 7: Infeasible predictions as a function of $\mu_m(x) = \mu_m(x^*)$ if $y = y^* = (1, 0)$ using regression (red; also smallest), largest (blue), midpoint (purple), and maximum entropy (green) methods. In both panels $\pi_m(x) = \pi_m(x^*)$ and $\pi_y(x) = \pi_y(x^*)$. In the left panel, $\pi_m(x) = 4/9$ and $\pi_y(x) = 1/4$, whereas in the right panel $\pi_m(x) = 1$ and $\pi_y(x) = 0$.

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 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

$$\rho(e, v, x) = \mathbb{1}(v > 1 - \tilde{\mu}), \quad e \in S_m(x), \quad (30)$$

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

$$q_v\{(1, 0) \mid x^*, x, (1, 0)\} = \frac{1}{\tilde{\mu}} \int_0^1 h(v, x)h(v, x^*) dv. \quad (31)$$

For most functions p , (31) will be close to zero if $\tilde{\mu}$ is close to zero. However, for p in (30), 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 p in (30); see also example 5. The maximum entropy method, in contrast, does not have this problem. \square

6.3 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 ρ in theorem 3 using p^* is given by

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

We continue to focus on the case in which $\mu_m(x) = \mu_m(x^*)$ and hence $h^*(v, x) = h^*(v, x^*)$ for all v , as explained in section 6.2. Consequently, the formula for $\rho^*(x, x^*)$ is again given by (29), i.e. $\rho^*(x, x^*) = \mathcal{L}''\{\lambda_m(x^*)\}$. The function α_y in theorem 3 using ρ^* is provided in (47) 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

$$[\pi_y(x)\{1 - \pi_y(x)\} + \mu_m^2(x)\pi_m(x)\{1 - \pi_m(x)\} - 2\pi_y(x)\pi_m(x)\mu_m(x)] / \tilde{\pi}_y(x). \quad (32)$$

Example 3 below is based on (29) and (32).

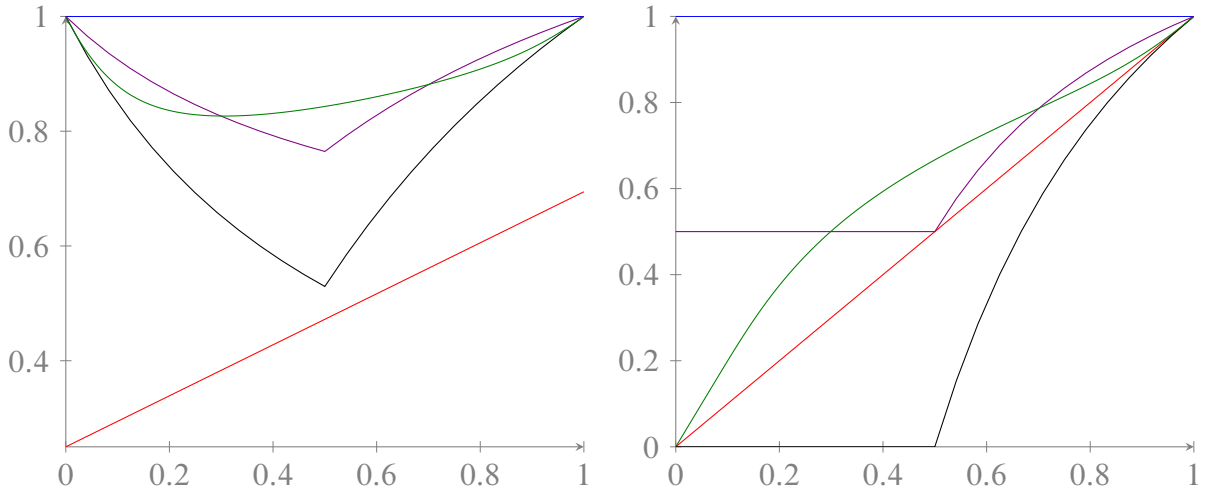


Figure 8: Infeasible predictions as a function of $\mu_m(x) = \mu_m(x^*)$ if x and x^* correspond to the same S -regions with $x \neq x^*$, and $y = y^* = (1, 0)$. The predictions are produced using the regression (red), smallest (black), midpoint (purple), largest (blue), and maximum entropy (green) 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 $\pi_y(x) = 0$.

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 8 and 9 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 $q_{ev}(y^* \mid x, x^*, y)$.

First consider the right panel in figure 8. 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

트위터

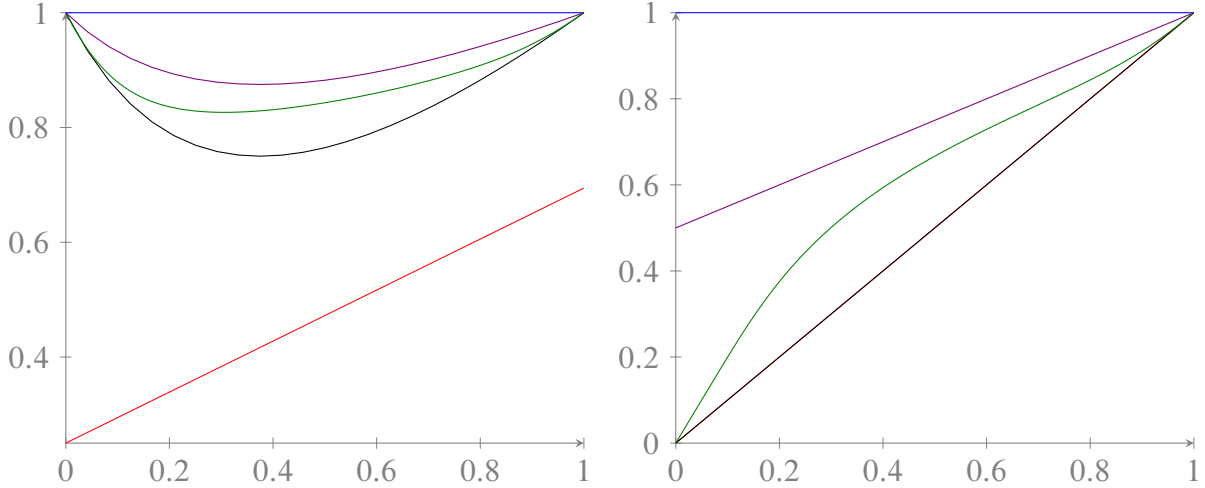


Figure 9: Infeasible predictions as a function of $\mu_m(x)$ if $x = x^*$ and $y = y^* = (1, 0)$, using regression (red), smallest (black), midpoint (purple), largest (blue), and maximum entropy (green) methods. In the left panel, $\pi_m(x) = 4 / 9$ and $\pi_y(x) = 1 / 4$, whereas in the right panel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

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 8 and 9 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

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

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 $p(e, v, x)$ and $p(e, v, x^*)$ 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 5, we restricted p to be flat in e . However, since the parameter space \mathcal{F}_x in (20) 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 .

6.4 Dirichlet–based predictions, midpoints, and maximum entropy: We now consider the average risk optimal predictions using the prior we proposed in section 5. Instead of choosing a particular loss function, we consider the distribution of the prediction when p is drawn by the algorithm

described in section 5.3. 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 7 and 8, 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

$$\frac{\mathbb{E}\{k(\mathbf{v}, x)k(\mathbf{v}, x^*)\}}{\mu_m(x)}. \quad (33)$$

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

$$\frac{\mathbb{E}\{p(\mathbf{e}, \mathbf{v}, x)p(\mathbf{e}, \mathbf{v}, x^*)\}}{\mu_m(x)}. \quad (34)$$

Now, consider the restricted parameter space defined in (21). Then, in our current design, (33) and (34) coincide and equal

$$T_p(x, x^*, y, y^*) = \frac{\mathbb{E}\{p(\mathbf{v}, x)p(\mathbf{v}, x^*)\}}{\mu_m(x)}. \quad (35)$$

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 (35) is an upper bound of both (33) and (34).

Figure 10 depicts the distribution of (35) when the function p is a random draw from the Dirichlet-like-process $\mathcal{D}_{\mu_m}^*$ described in section 5.3. Recall from section 5 that this entails pretending that each probability distribution is a possible distribution of \mathbf{p} for given μ_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.

For each value of μ_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 (35) for a different function p .

The results of our experiments are depicted in figure 10. The left panel in figure 10 is identical to the right panel in figure 7 and the right panel in figure 10 is identical to the right panel in figure 8, 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 μ_m , the dashed line the median, and the bottom and top of the grey 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 (35), 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 10 is for essentially all μ_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.

Paul:restricted b/c of independence of e ?

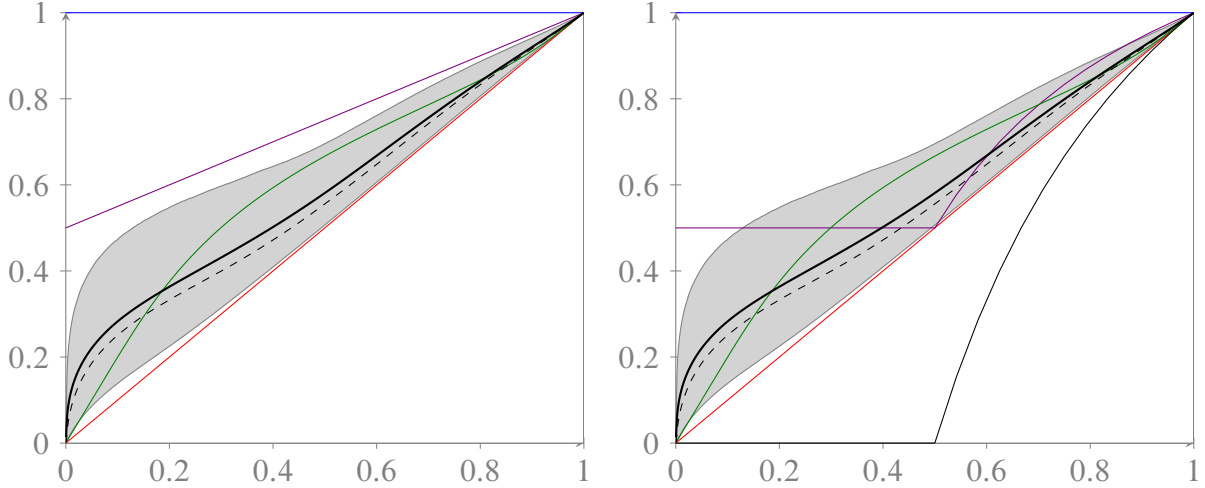


Figure 10: Infeasible predictions as a function of $\mu_m(x) = \mu_m(x^*)$ if $y = y^* = (1, 0)$ using regression (red), largest (blue), midpoint (purple), and maximum entropy (green) methods. The grey regions depict the area between 10% and 90% quantiles of the Dirichlet process–based experiment described in section 5.3 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.

Further, we chose a uniform prior and pseudocount hyperparameter equal to one for the Dirichlet process for convenience.²⁵ 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 (35) 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,

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

where $\hat{h}(\mathbf{v}, x) = \mathbb{E}p(\mathbf{e}, \mathbf{v}, x)$, 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 \hat{h} cannot exceed one,

²⁵The 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 μ_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.

트위터

it must be true that $\mu_m(x^*) \geq \mu_m(x)$ and that $\hat{k}(v, x^*) = 1$ whenever $\hat{k}(v, x) > 0$. Since \hat{k} is nondecreasing in v , it must be true that $\hat{k}(v, x) = \mathbb{1}\{v \geq 1 - S_m(x)\}$ and that $\hat{k}(v, x^*) = \hat{k}(v, x)$ 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 $\hat{k}(v, x) = \mu_m(x)$ for all values of v . \square

7. 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 four 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, 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,

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

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

$$\text{Cov}[\mathbb{1}\{y(\mathbf{e}^*, \mathbf{u}^*, \mathbf{v}^*, x^*) = y^*\}, \mathbb{1}\{y(\mathbf{e}, \mathbf{u}, \mathbf{v}^*, x) = y\} | \mathbf{e} \in S_m(x), \mathbf{e}^* \in S_m(x^*)]. \quad (36)$$

The covariance in (36) equals zero unless both y and y^* belong to $\{(1, 0), (0, 1)\}$. For $y = y^* = (1, 0)$ the covariance in (36) equals (by the law of iterated expectations)

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

The remaining three cases follow analogously. \square

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

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

트위터

Proof of Theorem 3. We have

$$\begin{aligned}
& \mathbb{P}\{y(\mathbf{u}^*, \mathbf{e}, \mathbf{v}, x^*) = y^*, y(\mathbf{u}, \mathbf{e}, \mathbf{v}, x) = y\} \\
&= \mathbb{E}\{c_y(\mathbf{e}, x)c_{y^*}(\mathbf{e}, x^*)\} + \mathbb{E}\{c_y(\mathbf{e}, x)c_m(\mathbf{e}, x^*)b_{y^*}(\mathbf{e}, \mathbf{v}, x^*)\} \\
&+ \mathbb{E}\{c_m(\mathbf{e}, x)c_{y^*}(\mathbf{e}, x^*)b_y(\mathbf{e}, \mathbf{v}, x)\} + \mathbb{E}\{c_m(\mathbf{e}, x)c_m(\mathbf{e}, x^*)b_y(\mathbf{e}, \mathbf{v}, x)b_{y^*}(\mathbf{e}, \mathbf{v}, x^*)\} \\
&= \pi_{mm}\text{Cov}\{b_y(\mathbf{e}, \mathbf{v}, x), b_{y^*}(\mathbf{e}, \mathbf{v}, x^*) \mid \mathbf{e} \in S_{mm}(x, x^*)\} \\
&\quad + \text{Cov}\{a_y(\mathbf{e}, x), a_{y^*}(\mathbf{e}, x^*)\} + \mathbb{E}a_y(\mathbf{e}, x)\mathbb{E}a_{y^*}(\mathbf{e}, x^*). \quad (37)
\end{aligned}$$

The stated result follows if one divides the right hand side in (37) 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 π_B, μ_C , etcetera, be similarly defined. We then solve

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

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 μ_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

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

$$\begin{aligned}
& + \pi_{ym}(x, x^*) \mathbb{E}\{\bar{b}_{y^*}(e, x^*) \mid e \in S_{ym}(x, x^*)\} \\
& + \pi_{mm}(x, x^*) \mathbb{E}\{b_y(e, v, x) b_{y^*}(e, v, x^*) \mid e \in S_{mm}(x, x^*)\}. \quad (38)
\end{aligned}$$

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

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

and let \bar{z} 's with different subscript and superscript combinations be defined analogously. Then the right hand side in (38) reduces to

$$\begin{aligned}
& \pi_{yy^*}(x, x^*) + \pi_{my^*}(x, x^*) \bar{z}_{my^*}^{xy}(x, x^*) + \pi_{ym}(x, x^*) \bar{z}_{ym}^{x^*y^*}(x, x^*) \\
& + \pi_{mm}(x, x^*) \mathbb{E}\{b_y(e, x) b_{y^*}(e, x^*) \mid e \in S_{mm}(x, x^*)\}. \quad (39)
\end{aligned}$$

Therefore, we seek to minimize/maximize (39) subject to the (mean) restrictions, namely (dropping the (x, x^*) arguments)

$$\begin{cases} \pi_{ym} \bar{z}_{ym}^{x^*y^*} + \pi_{mm} \bar{z}_{mm}^{x^*y^*} + \pi_{rm} \bar{z}_{rm}^{x^*y^*} = \phi_{y^*}(x^*), \\ \pi_{my^*} \bar{z}_{my^*}^{xy} + \pi_{mm} \bar{z}_{mm}^{xy} + \pi_{mr^*} \bar{z}_{mr^*}^{xy} = \phi_y(x). \end{cases} \quad (40)$$

Since $b_y(e, x)$ must belong to $[0, 1]$, (39) 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, x, 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 (39) subject to (40) is illustrated in figure 11: 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 11 so that $\mathbb{P}(y = y, y^* = y^* \mid x = x, x^* = x^*) = a + b + c + c$. Then, note that a, b, c, d correspond to the four terms in (39), (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 (40) can likewise be expressed in terms of the probabilities illustrated at the top of figure 11. So, using the shorthand $\phi = \phi_y(x)$ and $\phi^* = \phi_{y^*}(x^*)$, minimizing (39) subject to (40) can be formulated as

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

트위터

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 11. Then, maximizing (39) subject to (40) 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

$$\begin{cases} f_e(e)f_x(x)\{1 + \log f^*(p | e, x)\} - \nu_1(e, x) - \nu_2(x)p f_e(e) \mathbb{1}\{e \in S_m(x)\} = 0, \\ \int f^*(p | e, x) dp = 1, \\ \int_{S_m(x)} \int_0^1 p f^*(p | e, x) dp f_e(e) de = \mu_m(x)\pi_m(x), \end{cases}$$

almost everywhere, where ν_1, ν_2 are Lagrangean parameters. Thus, taking $\lambda(x) = \nu_2(x) / f_x(x)$, 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 (15).

Further, the second condition in (16) 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 (15) 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

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 8. $q_e(y^* | x^*, x, y) = \tilde{\pi}_{y^*}(x^*) + \alpha_{yy^*}(x, x^*)$. 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

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

which coincides with (38), except that b in the last term of (38) is replaced with \bar{b} here. Since the bounds in theorem 5 obtain for functions ρ that are flat in v , as noted in the proof of theorem 5, the bounds must coincide. \square

Theorem 9.

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

where $\rho_u(x, x^*) = \mathbb{E}\{\min\{p(e, \mathbf{v}, x), p(e^*, \mathbf{v}^*, x^*)\} \mid e \in S_m(x), e^* \in S_m(x^*)\} - \mu_m(x)\mu_m(x^*)$. 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 $\mathbb{E}\{\min(\mathbf{p}, \mathbf{p}^*) \mid \cdot\} \leq \min\{\mathbb{E}(\mathbf{p} \mid \cdot), \mathbb{E}(\mathbf{p}^* \mid \cdot)\}$, 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 $\mathbb{E}\{\min(\mathbf{p}, \mathbf{p}^*) \mid \cdot\} \geq \mathbb{E}(\mathbf{p}\mathbf{p}^* \mid \cdot)$, where \mathbf{p}, \mathbf{p}^* are independent given e, e^* . \square

Theorem 10.

$$q_{uv}(y^* | x^*, x, y) = \tilde{\pi}_{y^*}(x^*) + \delta_y \delta_{y^*} \frac{\pi_m(x)\pi_m(x^*)}{\tilde{\pi}_y(x)} \rho_{uv}(x, x^*), \quad (41)$$

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

트위터

Proof. Define

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

Note that

$$\begin{aligned} & \mathbb{P}\{y(\mathbf{e}^*, \mathbf{u}, \mathbf{v}, x^*) = y^*, y(\mathbf{e}, \mathbf{u}, \mathbf{v}, x) = y\} \\ &= \tilde{\pi}_y(x)\tilde{\pi}_{y^*}(x, x^*) + \\ & \quad \pi_m(x)\pi_m(x^*)\text{Cov}\{h_y(\mathbf{e}, \mathbf{u}, \mathbf{v}, x), h_{y^*}(\mathbf{e}^*, \mathbf{u}, \mathbf{v}, x^*) \mid \mathbf{e} \in S_m(x), \mathbf{e}^* \in S_m(x^*)\}. \end{aligned} \quad (42)$$

Tedious but simple algebra shows that the covariance in (42) equals $\delta_y\delta_{y^*}\rho_{uv}(x, x^*)$, which yields (41).

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 $\mathbb{E}\{\min(\mathbf{p}, \mathbf{p}^*) \mid \cdot\} \leq \min\{\mathbb{E}(\mathbf{p} \mid \cdot), \mathbb{E}(\mathbf{p}^* \mid \cdot)\}$, with equality for $p(e, v, x) = \mathbb{1}\{v \geq 1 - \mu_m(x)\}$. For the lower bound,

$$\mathbb{E}\{\min(\mathbf{p}, \mathbf{p}^*) \mid \mathbf{e}, \mathbf{e}^*\} \geq \mathbb{E}(\mathbf{p}\mathbf{p}^* \mid \mathbf{e}, \mathbf{e}^*) \geq \mathbb{E}(\mathbf{p} \mid \mathbf{e})\mathbb{E}(\mathbf{p}^* \mid \mathbf{e}^*) \quad (43)$$

where the inequalities hold with equality for $p(e, v, x) = \mathbb{1}\{e \in S_m(x)\}$. Take the expectations in (43) 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

$$g_y(e, v, x) = c_y(e, x) + c_m(e, x)b_y(e, v, x).$$

Theorem 11.

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

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 $q_{eu}(y^* \mid x^*, x, y)$ are identical to those of $q_{ev}(y^* \mid x^*, x, y)$ 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, $\tilde{\pi}_y(x)q_{eu}(y^* \mid x^*, x, y)$ equals

$$\begin{aligned} & \pi_{yy^*}(x, x^*) + \pi_{ym}(x, x^*)\mathbb{E}\{p(\mathbf{e}, \mathbf{v}^*, x^*) \mid \mathbf{e} \in S_{ym}(x, x^*)\} \\ & \quad + \pi_{my^*}(x, x^*)\mathbb{E}\{p(\mathbf{e}, \mathbf{v}, x) \mid \mathbf{e} \in S_{my^*}(x, x^*)\} \\ & \quad + \pi_{mm}(x, x^*)\mathbb{E}[\min\{p(\mathbf{e}, \mathbf{v}, x), p(\mathbf{e}, \mathbf{v}^*, x^*)\} \mid \mathbf{e} \in S_{mm}(x, x^*)]. \end{aligned} \quad (44)$$

If $x = x^*$, the second and third terms in (44) are equal to zero, so we only need to consider the last term. The upper bound follows from the fact that

$$\mathbb{E}\{\min\{\mathbf{p}, \mathbf{p}^*\} \mid \cdot\} \leq \min\{\mathbb{E}(\mathbf{p} \mid \cdot), \mathbb{E}(\mathbf{p}^* \mid \cdot)\}. \quad (45)$$

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 (45) 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 (44) is bounded above by $\pi_{yy^*}(x, x^*) + U_{yy^*}(x, x^*)$, where U_{yy^*} is as defined in (12). Sharpness follows from the fact that (44) is bounded below by (38). 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 (38), except that b_y is replaced with \bar{b}_y . \square

Theorem 12.

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

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

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

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 $\tilde{\pi}_y(x)q_{euv}(y^* \mid x^*, x, y)$ equals

$$\begin{aligned} &\pi_{yy^*}(x, x^*) + \pi_{ym}(x, x^*)\mathbb{E}\{p(\mathbf{e}, \mathbf{v}, x^*) \mid \mathbf{e} \in S_{ym}(x, x^*)\} \\ &\quad + \pi_{my^*}(x, x^*)\mathbb{E}\{p(\mathbf{e}, \mathbf{v}, x) \mid \mathbf{e} \in S_{my^*}(x, x^*)\} \\ &\quad + \pi_{mm}(x, x^*)\mathbb{E}[\min\{p(\mathbf{e}, \mathbf{v}, x), p(\mathbf{e}, \mathbf{v}, x^*)\} \mid \mathbf{e} \in S_{mm}(x, x^*)]. \end{aligned} \quad (46)$$

Because $\mathbb{E}\{\min(\mathbf{p}, \mathbf{p}^*) \mid \cdot\} \leq \min\{\mathbb{E}\{\mathbf{p} \mid \cdot\}, \mathbb{E}\{\mathbf{p}^* \mid \cdot\}\}$, $\mathbb{E}\{p(\mathbf{e}, \mathbf{v}, x) \mid \mathbf{e} \in S_m(x)\} = \mu_m(x)$, and $\mathbb{E}\{p(\mathbf{e}, \mathbf{v}, x^*) \mid \mathbf{e} \in S_m(x^*)\} = \mu_m(x^*)$, (46) is bounded above by $\pi_{yy^*}(x, x^*) + U_{yy^*}(x, x^*)$, where U_{yy^*} is as defined in (12). For sharpness, note that (46) is bounded below by (38).

Now the lower bound. Since (46) is no less than (38), the lower bound for (38) is a lower bound here, also. It remains to be shown that it can be attained. Recall from the proof of theorem 5 that (38) 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, (38) equals (46). \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.

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

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

$$\begin{aligned} \bar{a}_y^*(e, x) &= c_y(e, x) + c_m(e, x)\mathbb{E}\delta_y^*(e, \mathbf{v}, x) \\ &= \begin{cases} c_y(e, x) + c_m(e, x)\mu_m(x), & y = (1, 0), \\ c_y(e, x) + c_m(e, x)\{1 - \mu_m(x)\}, & y = (0, 1), \\ c_y(e, x), & y \in \{(0, 0), (1, 1)\}, \end{cases} \end{aligned} \quad (47)$$

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

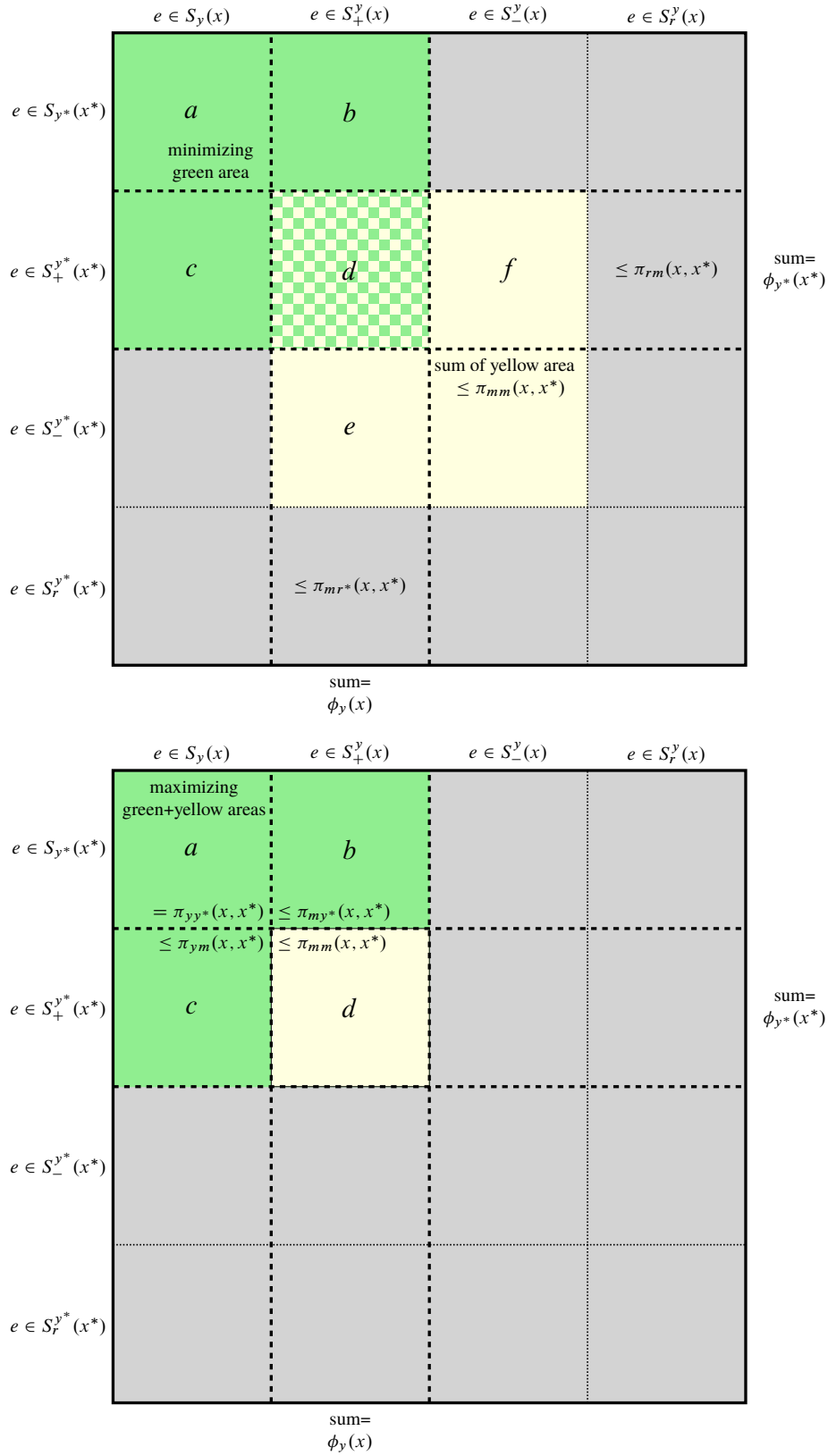


Figure 11: Finding the lower bound (top) and the upper bound (bottom)

References

- AGUIRREGABIRIA, V., AND P. MIRA (2007): “Sequential estimation of dynamic discrete games,” *Econometrica*, 75(1), 1–53.
- ARADILLAS-LOPEZ, A. (2011): “Nonparametric probability bounds for Nash equilibrium actions in a simultaneous discrete game,” *Quantitative Economics*, 2(2), 135–171.
- ARADILLAS-LÓPEZ, A., AND E. TAMER (2008): “The identification power of equilibrium in simple games,” *Journal of Business & Economic Statistics*, 26(3), 261–283.
- ATHEY, S., AND G. W. IMBENS (2007): “Discrete choice models with multiple unobserved choice characteristics*,” *International Economic Review*, 48(4), 1159–1192.
- AUMANN, R. J. (1961): “Borel structures for function spaces,” *Illinois Journal of Mathematics*, 5(4), 614–630.
- BAJARI, P., C. L. BENKARD, AND J. LEVIN (2007): “Estimating dynamic models of imperfect competition,” *Econometrica*, 75(5), 1331–1370.
- BAJARI, P., H. HONG, AND S. P. RYAN (2010): “Identification and estimation of a discrete game of complete information,” *Econometrica*, 78(5), 1529–1568.
- BAKER, A. (2013): “Simplicity,” in *Stanford Encyclopedia of Philosophy*, ed. by E. N. Zalta. Fall 2013 edn.
- BJORN, P. A., AND Q. H. VUONG (1984): “Simultaneous equations models for dummy endogenous variables: a game theoretic formulation with an application to labor force participation,” Discussion paper, CalTech.
- BRESNAHAN, T. F., AND P. C. REISS (1991a): “Empirical models of discrete games,” *Journal of Econometrics*, 48(1), 57–81.
- (1991b): “Entry and competition in concentrated markets,” *Journal of Political Economy*, pp. 977–1009.
- BRIESCH, R. A., P. K. CHINTAGUNTA, AND R. L. MATZKIN (2012): “Nonparametric discrete choice models with unobserved heterogeneity,” *Journal of Business & Economic Statistics*.
- BULOW, J. I., J. D. GEANAKOPOLOS, AND P. D. KLEMPERER (1985): “Multimarket oligopoly: Strategic substitutes and complements,” *Journal of Political economy*, 93(3), 488–511.
- CASS, D., AND K. SHELL (1983): “Do sunspots matter?,” *Journal of Political Economy*, 91(2), 193–227.
- CHEN, X., AND D. POUZO (2012): “Estimation of nonparametric conditional moment models with possibly nonsmooth generalized residuals,” *Econometrica*, 80(1), 277–321.
- CILIBERTO, F., AND E. TAMER (2009): “Market structure and multiple equilibria in airline markets,” *Econometrica*, 77(6), 1791–1828.

- COVER, T. M., AND J. A. THOMAS (2012): *Elements of information theory*. John Wiley & Sons.
- GALICHON, A., AND M. HENRY (2011): “Set identification in models with multiple equilibria,” *The Review of Economic Studies*, p. rdr008.
- GOLAN, A., G. G. JUDGE, AND D. MILLER (1996): *Maximum entropy econometrics: Robust estimation with limited data*. Wiley New York.
- GRIECO, P. L. (2014): “Discrete games with flexible information structures: An application to local grocery markets,” *Rand Journal of Economics*, 45(2), 303–340.
- HECKMAN, J. J. (2005): “The scientific model of causality,” *Sociological methodology*, 35(1), 1–97.
- JAYNES, E. (1957a): “Information theory and statistical mechanics,” *Physical Review*, 106, 620–630.
- (1957b): “Information theory and statistical mechanics II,” *Physical Review*, 108, 171–190.
- JIA, P. (2008): “What happens when Wal-Mart comes to town: An empirical analysis of the discount retailing industry,” *Econometrica*, pp. 1263–1316.
- KALAI, A., AND E. KALAI (2012): “Cooperation in strategic games revisited,” *Quarterly Journal of Economics*, p. qjs074.
- KASHAEV, N. (2015): “Testing for Nash behavior in entry games with complete information,” Discussion paper, Penn State.
- KASHAEV, N., AND B. SALCEDO (2015): “Identification of solution concepts for discrete semi-parametric games with complete information,” Discussion paper, Penn State.
- KASY, M. (2011): “Identification in triangular systems using control functions,” *Econometric Theory*, 27(03), 663–671.
- KLINE, B. (2015): “Identification of complete information games,” *Journal of Econometrics*, 189(1), 117–131.
- KLINE, B., AND E. TAMER (2012): “Bounds for best response functions in binary games,” *Journal of Econometrics*, 166(1), 92–105.
- KOOREMAN, P. (1994): “Estimation of econometric models of some discrete games,” *Journal of Applied Econometrics*, 9(3), 255–268.
- LIU, N., Q. VUONG, AND H. XU (2013): “Rationalization and identification of discrete games with correlated types,” Discussion paper, University of Texas.
- MAGNOLFI, L., AND C. RONCORONI (2016): “Estimation of discrete games with weak assumptions on information,” Discussion paper, Yale.
- MANSKI, C. (2015): “Interpreting point predictions,” Discussion paper, Northwestern University.
- PAKES, A., M. OSTROVSKY, AND S. BERRY (2007): “Simple estimators for the parameters of discrete dynamic games (with entry/exit examples),” *The Rand Journal of Economics*, 38(2), 373–399.

- REGUANT, M. (2016): “Bounding equilibria in counterfactual analysis,” Discussion paper, Northwestern University.
- SCHENNACH, S. M. (2014): “Entropic latent variable integration via simulation,” *Econometrica*, 82(1), 345–385.
- SEIM, K. (2006): “An empirical model of firm entry with endogenous product–type choices,” *Rand Journal of Economics*, 37(3), 619–640.
- SOETEVEENT, A. R., AND P. KOOREMAN (2007): “A discrete-choice model with social interactions: with an application to high school teen behavior,” *Journal of Applied Econometrics*, 22(3), 599–624.
- SONG, K. (2014): “Point decisions for interval–identified parameters,” *Econometric Theory*, 30(02), 334–356.
- TAMER, E. (2003): “Incomplete simultaneous discrete response model with multiple equilibria,” *Review of Economic Studies*, 70(1), 147–165.
- TEH, Y. W. (2010): “Dirichlet process,” in *Encyclopedia of machine learning*, pp. 280–287. Springer.
- XU, H. (2014): “Estimation of discrete games with correlated types,” *The Econometrics Journal*, 17(3), 241–270.