

# Information and admissible sets

Jeff Rowley

1<sup>st</sup> September, 2014

## Abstract

«Abstract here»

**Acknowledgements.** *I acknowledge the R and L<sup>A</sup>T<sub>E</sub>X communities, and the wealth of knowledge that they have made freely available to all. I thank Andrew Chesher and Toru Kitagawa for their supervision and support. I further thank Adam Rosen for helpful discussion. I gratefully acknowledge financial support from the Economic and Social Research Council (ESRC).*

I explore the effect of incorporating information for a non-parametric binary choice model. The model permits endogenous variation in a scalar random variable that is due to non-random selection, and it is the average causal effect of this endogenous variable on the outcome variable that is of interest. The model embeds an exclusion restriction and an independence restriction that together define an instrumental variable but is silent as to the relationship between the endogenous variable and the instrumental variable. I restrict the relationship between the outcome variable and the endogenous variable up to a non-parametric threshold crossing function. The model is credible (Manski, 2013) in that it embeds restrictions that impose weaker constraints on assumed behaviour, but does not identify the average causal effect of the endogenous variable on the outcome variable.<sup>1</sup> Rather, the model partially identifies the average causal effect of the endogenous variable on the outcome variable.

I define information to be those additional characteristics of economic agents that are observable with the caveat that these characteristics are exogenous and are relevant to the latent structure. It is convenient to think of such characteristics as being predetermined and immutable; characteristics that result from choices that are made jointly with the outcome variable are excluded by the definition. Accordingly, exogenous variables and instrumental variables are each regarded as information, and I distinguish between these classes of information. I study how the admissible set of values for the average causal effect of the endogenous variable on the outcome variable changes as each class of information is incorporated into the model separately.

It is useful to distinguish between classes of information since each class enters the latent structure in a different way. Exogenous variables are permitted to enter the structural equation for the outcome variable and to determine the endogenous variable. As such, exogenous variables can be seen to enrich both individual response and individual selection, respectively. An important consequence is that the causal effect of the endogenous variable on the outcome variable depends upon

---

<sup>1</sup>Assumptions that cannot be tested using data. The model does embed some non-trivial non-verifiable restrictions that might be relaxed.

the values of exogenous variables when individual response is enriched. In contrast, instrumental variables are excluded from the structural equation for the outcome variable by definition and so only enrich individual selection. Given this, the effect of incorporating information is different depending upon the class of information that is being incorporated into the model.

Incorporating information of either class is generally sensible for a number of reasons. Firstly, incorporating information is known to be efficient; variation that is attributable to an observable variable is instead attributable to unobservable heterogeneity when that variable is omitted. Secondly, the effect of incorporating information for partially identifying models is not well-documented. In identifying models, if an omitted variable is a determinant of the endogenous variable then its effect on the outcome variable is instead attributable to the endogenous variable, and is a bias. A contribution that I make is in showing that incorporating information negates this bias, and is equivalent to a shift in the location of the admissible set. A further reason to particularly favour incorporating exogenous variables is that the average causal effect of the endogenous variable on the outcome variable in identifiable sub-populations can be recovered. I name this structural characteristic the conditional average causal effect of the endogenous variable on the outcome variable, and index it by the conditioning value.<sup>2</sup> Understanding the effect of an intervention in sub-populations can be interesting if the intervention can be targeted or if the intervention is to be applied elsewhere in a population that differs according to its observable characteristics.

A relevant question is how to relate conditional causal effects to (unconditional) causal effects. More precisely, how does the average causal effect of the endogenous variable on the outcome variable relate to its conditional counterparts? I show that the average causal effect of the endogenous variable on the outcome variable can be expressed as a Minkowski summation of its conditional counterparts when the non-parametric binary choice model is augmented. I derive sharp bounds on the conditional average causal effect by applying random set theory.

I demonstrate application of the non-parametric binary choice model, elucidating the practical difficulties that arise when estimating set identifying models (focusing on those issues that arise from incorporating information). As in Chesher and Rosen (2013), I estimate the average causal effect of additional children on a mother's employment using US census data. I extend Chesher and Rosen (2013) in a number of ways. Firstly, I report statistical uncertainty in the estimate of the average causal effect of additional children on a mother's employment using a method that is outlined in Chernozhukov et al. (2013). Secondly, I enrich the support of the instrumental variable and explore the effect that this has on the admissible set of values for the average causal effect of additional children on a mother's employment, and on its accompanying confidence region. Thirdly, I enrich individual response by permitting the structural equation for employment to depend upon predetermined and immutable characteristics of mothers. I discuss the complication of calculating statistical uncertainty when exogenous variables are permitted to enter the structural equation for employment. With respect to the second and third extensions, it is necessary that I augment the model by embedding additional restrictions. In fact, Chesher and Rosen (2013) describe the augmented non-parametric binary choice model that I assume but simplify this model for application (by excluding exogenous variables from the structural equation for the outcome variable). I discuss how the augmented model relates to the simplified model in each case and the credibility of the additional restrictions that are embedded in the augmented model.

## Related research

Other notable non-parametric binary choice models are described in Balke and Pearl (1997) and Shaikh and Vytlacil (2011), and general non-parametric models of choice are described in Chesher (2005), Kitagawa (2009) and Chesher (2010).

Balke and Pearl (1997) assumes a triangular model (the model embeds a structural equation for the outcome variable and a structural equation for the endogenous variable; see Strotz and Wold

---

<sup>2</sup>The conditioning value is specifically the value of the exogenous variables. Heckman and Vytlacil (2005) defines a parameter  $ATE(x)$  that is equivalent to the conditional average causal effect of the endogenous variable on the outcome variable at the conditioning value  $x$ . Khan and Tamer (2010) and Abrevaya et al. (2013) instead refer to this parameter as the conditional average treatment effect and abbreviate this to  $CATE(x)$ .

(1960) for a detailed discussion of triangular models) that relaxes separability of unobservable heterogeneity in the structural equation for the outcome variable. The cost is that the model is no longer silent as to the relationship between the endogenous variable and the instrumental variable. The model does not permit exogenous variables to enter the structural equation for the outcome variable. I discuss the credibility of separability of unobservable heterogeneity in the main text. Shaikh and Vytlacil (2011) assumes a triangular model but maintains separability of unobservable heterogeneity in the structural equation for the outcome variable. The model permits exogenous variables to enter the structural equation for the outcome variable.

Chesher (2005) and Kitagawa (2009) describe non-parametric models that permit continuous variation in the outcome variable. Chesher (2005) assumes a triangular model that relaxes separability of unobservable heterogeneity in the structural equation for the outcome variable. The model permits exogenous variables to enter the structural equation for the outcome variable, although local invariance of the structural equation for the outcome variable to variation in the exogenous variables is embedded. The model is uninformative when there is binary variation in the endogenous variable but is informative when there is discrete variation. Kitagawa (2009) extends Balke and Pearl (1997) to permit discrete and continuous variation in the outcome variable, and studies commonly invoked restrictions on covariation of the instrumental variable and unobservable heterogeneity.

Chesher (2010) describes an ordered choice model that permits discrete variation in the outcome variable. Chesher (2010) assumes a single equation model that relaxes separability of unobservable heterogeneity in the structural equation for the outcome variable, although monotonicity of the structural equation for the outcome variable in unobservable heterogeneity is embedded. The model permits binary or discrete variation in the endogenous variable.

## Notation

I study a probability space  $(\Omega, \Sigma, \mathbb{P})$ . I define random variables on this probability space. I write random variables as upper case Latin letters, and I write realisations (or specific values) of random variables as lower case Latin letters. I write the support of  $A$  as  $\mathcal{R}_A$ . I write the counterfactual value of  $A$  when  $B$  has a causal interpretation and is externally fixed as  $A(b)$ . I write the average causal effect of  $B$  on  $A$  as  $ACE(B \rightarrow A)$ , and the conditional average causal effect of  $B$  on  $A$  given  $C$  as  $ACE(B \rightarrow A|c)$ .

I refer to  $Y$  as the outcome variable, to  $D$  as the endogenous variable, to  $X$  as the endogenous variable, to  $Z$  as the instrumental variable, and to  $U$  as unobservable heterogeneity. Despite the use of *the*, I permit  $(X, Z)$  to be vectors. I write the structural equation for the outcome variable as  $h$ , and the structural equation for the endogenous variable as  $g$ .

I write the dimension of a set  $S$  as  $\dim(S)$ , and its length as  $\|S\|$ . I write the expectation operator as  $\mathbb{E}$ , and the indicator function as  $\mathbb{1}$ . I write  $A$  is independent of  $B$  as  $A \perp\!\!\!\perp B$ . To distinguish between population and sample quantities, I subscript sample quantities by  $n$ .

I introduce further terminology and notation in Figure 1 through Figure 4. This specifically relates to models and structures, and is consistent with the approach that is formally laid out in Hurwicz (1950) and in Koopmans and Reiersøl (1950).

## Application

I estimate the average causal effect of additional children on a mother's employment using United States census data. The data are obtainable from Angrist (2014), and are described in Angrist and Evans (1998). To summarise, the dataset consists of 254,654 households that were recorded as part of the 1980 United States census. The dataset specifically contains observations of married households with at least two children under the age of 18 years and where the mother is aged between 21 years and 35 years. For clarity, I translate each variable in the data into the mathematical notation that I employ.

$$Y \equiv \mathbb{1}[\text{Mother is employed in 1979}]$$

$$D \equiv \mathbb{1}[\text{Three or more children in the household}]$$

As  $(X, Z)$  are continually redefined in the main text, I do not define these variables as I do  $(Y, D)$ . Instead, I note those variables that at some point or other form part of the definition of  $(X, Z)$ .  $X$  is a function of mother's race or ethnicity (shortened to race).<sup>3</sup>  $Z$  is a function of whether the oldest two children in the household share the same gender (shortened to child gender), and whether the second pregnancy was a multiple pregnancy.

I refer to the application throughout the main text so as to illustrate how technical conditions on variables and on the relationship between variables restrict the behaviour of economic agents, in this case mothers. For brevity, I simply refer to race when discussing  $X$  in the context of the application, and child gender when discussing  $Z$  in the context of the application.

## 1 A non-parametric model of binary choice

**Axiom.** *Economic agents are utility maximising, selecting between alternatives in a choice set according to the utility that they attach to that choice. Utility is perfectly observable by economic agents and is determined by a well-defined utility function for each choice. Each agent is permitted to value each choice differently.*

I introduce the non-parametric binary choice model that is described in Chesher and Rosen (2013) (hereafter, the single equation model). The single equation model constitutes the set of structures that are consistent with Restriction M1 through Restriction M6.

**M1. Discrete support.**  $(Y, D, X, Z)$  are observable and have discrete supports (with at least two points of support). Further,  $(Y, D)$  have binary supports and are normalised such that

(a)  $\mathcal{R}_Y = \{0, 1\}$  and

(b)  $\mathcal{R}_D = \{0, 1\}$ ,

respectively.

Restriction M1 is a verifiable restriction.  $(Y, D, X, Z)$  are observable and so it is trivial to verify that each variable satisfies its support restriction.  $(Y, D)$  are normalised to be consistent with the application, but any other supports  $\{y_0, y_1\}$  and  $\{d_0, d_1\}$  can be generated by an affine transformation of  $h$  and of  $g$ .

**M2. Scalar  $U$ .**  $U$  is an unobservable scalar such that  $\mathcal{R}_U$  is an open subset of  $\mathbb{R}$  with strictly positive Lebesgue measure.

Restriction M2 is a non-verifiable restriction. The dimension of  $U$  is a normalisation rather than a restriction since  $D$  has binary support.

**M3. Joint independence.**  $U \perp (X, Z)$ .

Restriction M3 is a non-verifiable restriction. The restriction nests the restrictions  $U \perp Z|X$  (conditional independence) and  $U \perp X$ , and nests the restrictions  $U \perp X$  and  $U \perp Z$  (marginal independence).<sup>4</sup> In the context of the application, the restriction implies that variables such as opportunity are independent of child gender conditional on race, and are independent of race.

**M4. Exclusion.**  $Y = h(D, X, U)$ .

Restriction M4 is a non-verifiable restriction. The restriction excludes  $Z$  from  $h$  and so excludes  $Z$  from having a causal effect on  $Y$ . The restriction is equivalent to an order condition. In the context of the application, the restriction implies that child gender does not have a causal effect on the employment of mothers.

**M5. Monotonicity.**  $h$  is a non-parametric threshold crossing function that is separable in  $U$ .  $h$  is normalised to be increasing in  $U$ , and  $U$  is normalised to be distributed uniformly on the unit interval.

<sup>3</sup>Angrist and Evans (1998) also treats mother's age, and mother's age at the time of her first birth (shortened to birthing age) as exogenous variables.

<sup>4</sup>The restriction also implies the restrictions  $U \perp X|Z$  and  $U \perp Z$ .

Restriction M5 is a non-verifiable restriction. The restriction implies that individual response is monotonic. In the context of the application, the restriction implies that the causal effect of additional children on a mother's employment is positive for all mothers or is negative for all mothers. The restriction permits the threshold to be a non-parametric function of  $(D, X)$  and implies that the distribution of  $U$  can be relatively unrestricted beyond Restriction M2.

**M6. Relevance.** There exist values  $(z, z')$  such that  $\mathbb{P}(d|z) \neq \mathbb{P}(d|z')$  for all  $d \in \mathcal{R}_D$  and for some  $(z, z') \in \mathcal{R}_Z^2$ .

Restriction M6 is a verifiable restriction. The restriction states that  $Z$  covaries with  $D$ . A simple interpretation is that  $Z$  causes  $D$ , but the restriction itself is weaker than this in that it permits  $Z$  to be correlated with a variable that causes  $D$ .<sup>5</sup> The restriction is equivalent to a rank condition. In the context of the application, the restriction implies that the probability of having three or more children varies with child gender. For example, if the probability of having three or more children is greater when the oldest two children in the household share the same gender.

The single equation model partially identifies  $ACE(D \rightarrow Y)$ . The single equation model also partially identifies  $ACE(D \rightarrow Y|x)$  for all  $x \in \mathcal{R}_X$ . Restriction M1 through Restriction M6 can be written more compactly as Restriction M1' through Restriction M6'.

**M1'.**  $Y = \mathbb{1}[p(D, X) < U]$ .

**M2'.**  $U|(X, Z) \sim \text{unif}(0, 1)$ .

**M3'.**  $\mathbb{P}(d|z) \neq \mathbb{P}(d|z')$  for all  $d \in \mathcal{R}_D$  and for some  $(z, z') \in \mathcal{R}_Z^2$ .

**M4'.**  $\mathcal{R}_D = \{0, 1\}$ .

**M5'.**  $\mathcal{R}_X = \{x_1, \dots, x_K\}$  and  $K < \infty$ .

**M6'.**  $\mathcal{R}_Z = \{z_1, \dots, z_L\}$  and  $L < \infty$ .

## 2 Credibility in economic modelling

I define credibility. I discuss the conditions under which a model is more credible than another. I discuss opposition to the assumption and use of partially identifying models.

Credibility is a statement of the validity and the plausibility of the restrictions that a model embeds, and is a desirable property. The need to discuss both validity and plausibility arises because restrictions can be either verifiable or non-verifiable. The distinction between verifiable and non-verifiable restrictions is that verifiable restrictions are testable using data while non-verifiable restrictions cannot be tested even if data is collected for the population. As verifiable restrictions can be rejected or not rejected on the basis of observed behaviour, it makes sense to talk about such restrictions as being valid or invalid. In contrast, the validity of non-verifiable restrictions is indeterminable. Whether to accept a set of non-verifiable restrictions as an accurate representation of how economic agents behave is subjective and depends upon how plausible the restrictions seem. Restrictions that are founded in economic theory, or that impose weaker constraints on assumed behaviour are more plausible (a view that is consistent with Occam's razor, a widely accepted principle of parsimony).

I regard a model as incredible if the verifiable restrictions that it embeds are invalid. I regard a model as more credible relative to another if the verifiable restrictions that it embeds are valid and if the sum of the non-verifiable restrictions that it embeds are more plausible. Manski (2013) adopts an equivalent stance, formalised as The Law of Decreasing Credibility.

Models that embed restrictions that impose weaker constraints on assumed behaviour are typically not uniformly identifying. Instead, such models are typically partially identifying. More commonly,

(a) a more restrictive model is assumed that identifies a feature of interest; or,

<sup>5</sup>This point is the subject of Figure 5, which studies causality.

- (b) identification of a different feature is sought and a model that embeds restrictions that impose weak constraints on assumed behaviour is assumed.

I suggest that (a) and (b) are motivated by two concerns. Namely, that characterising the admissible set of structures or the admissible set of values for a structural characteristic of interest can be complex and computationally difficult, and that partially identifying models do not produce unique conclusions. Although tractability is a legitimate concern, there is an inherent and widespread misunderstanding that models that do not produce unique conclusions are inferior regardless of the restrictions that they embed. Conclusions that are produced by more credible models should always be preferred, even if these conclusions display ambiguity.

I caution against both (a) and (b). In (a), a more restrictive model is assumed specifically for the purpose of achieving identification. Koopmans and Reiersøl (1950) remarks that a model should be constructed purely from prior knowledge of the studied behaviour, and to do otherwise violates scientific honesty. In (b), the feature that is identified is often less valuable than the original feature of interest. Nonetheless, it is promising that (b) should implicitly recognise the importance of credibility.

## 2.1 The credibility of the single equation model

I discuss the credibility of the single equation model, generally and in the context of the application. I focus on the non-verifiable restrictions that the model embeds since it is these restrictions that are of principal interest when selecting from competing models.

First, the single equation model embeds the restriction that  $h$  is a non-parametric function (Restriction M5). In general, non-parametric restrictions are plausible since they permit the output of a function to depend arbitrarily on the value of its arguments. Non-parametric functions are flexible and are able to capture arbitrary variation that could otherwise only be captured using high-order polynomial functions or indicator functions. In particular, non-parametric functions are well-suited to capturing interaction between the arguments of a function. For example, if the difference in the employment of mothers between the counterfactual environments of two children in the household versus three or more children in the household varies systematically with race. Non-parametric functions are also well-suited to settings in which an argument is a categorical (and discrete) variable with no natural ordering. For example, race is a categorical variable with no natural ordering.

Second, the single equation model embeds the normalisation that  $U$  is distributed uniformly on the unit interval (Restriction M5). Note that, as a normalisation and not as a restriction, the normalisation imposes no constraints on the distribution of  $U$ . In general, restrictions that impose constraints on the distribution of  $U$  are implausible. This follows from the definition of  $U$  as a projection of unobservable determinants of utility onto an ordered set: there is no reason to suppose that economic agents should be distributed on this set according to some well-behaved distribution. Further, the normalisation is a normalisation and not a restriction because the single equation model embeds the restriction that  $h$  is a non-parametric function. The restriction that  $h$  is a non-parametric function implies not only that  $h$  depends arbitrarily on the value of its argument, but that it permits  $U$  to be distributed non-parametrically.

Third, the single equation model embeds the restriction that  $h$  excludes  $Z$  (Restriction M4). The restriction is an essential element of the single equation model if identification of causality is sought (and if  $Z$  is an instrumental variable). Further, the restriction is an essential element of any model of causality and not just an essential element of the single equation model. This point is the subject of Figure 5, which studies causality. In the context of the application, the restriction (together with Restriction M3) implies that child gender is conditionally independent of the employment of mothers. If a mother chooses to participate in the labour market only when her children share the same gender say, then this restriction is violated. It is plausible that child gender should not affect the decision of a mother to participate in the labour market and so the restriction is plausible.

Fourth, the single equation model is silent as to the relationship between  $D$  and  $Z$ . More importantly, it is silent as to the determination of  $D$ . That is, the single equation model embeds the restriction that the codomain of  $g$  is  $\{0, 1\}$  but otherwise is silent as to the arguments and

functional form of  $g$ . The lack of constraints on  $g$  is important since the single equation model then permits any endogenous relationship between  $Y$  and  $D$ , and does not restrict the relationship between  $D$  and  $Z$  to be causal. In the context of the application, whether there are three or more children in the household is endogenous. If mothers that incur large costs from employment also incur small costs from having children say, then whether there are three or more children in the household is correlated with a mother's employment. The single equation model does not restrict how strong this negative correlation should be, nor does it limit the reasons for or direction of endogenous variation.

Fifth, the single equation model embeds the restriction that  $U \perp\!\!\!\perp (X, Z)$  (Restriction M3). Joint independence is a strong restriction on the joint distribution of  $(D, X, U)$ . Joint independence is stronger than both conditional independence and marginal independence, and restricts the full distribution unlike mean independence or quantile independence restrictions. In general, joint independence is not a plausible restriction. In the context of the application, the restriction implies that variables such as opportunity are independent of child gender conditional on race, and are independent of race. It is plausible that opportunity is independent of child gender conditional on race since child gender is randomly assigned. However, there is no reason to suppose that opportunity is independent of race. Minority individuals need not have access to the same opportunities as otherwise equivalent white individuals. It is plausible that the distribution of opportunities for white individuals is more negatively skewed than the distribution of opportunities for minority individuals.

Sixth, the single equation model embeds the restriction that  $h$  is separable in  $U$  (Restriction M5). In general, this is an implausible restriction since it implies that individual response is monotonic. In the context of the application, the restriction implies that the causal effect of additional children on a mother's employment is positive for all mothers or is negative for all mothers. There is no reason to suppose that individual response is monotonic. A mother that uses a paid-for childcare service might find the cost of such services prohibitive for additional children and so exit employment. In contrast, a mother that relies on family members for childcare might find need to return to employment to increase household income.

As a final remark regarding the credibility of the single equation model, Restriction M6 is a weaker restriction than the monotonicity restriction that is embedded in the model that is described in Imbens and Angrist (1994) (hereafter, the late model), which identifies the local average treatment effect. Restriction M6 permits the existence of compliers and defiers. In other words, it permits individual selection to be both positive and negative. In fact, the monotonicity restriction that is embedded in the late model nests Restriction M6. In the context of the application, it is plausible that whether there are three or more children in the household depends upon the gender of the oldest two children. If preferences are convex then it is plausible that whether there are three or more children in the household is more likely when the oldest two children share the same gender, but there may be other considerations which determine perceived child quality.

The single equation model is reasonably credible in general, and in the context of the application. Nonetheless, some of the restrictions that the single equation model embeds are unsatisfactory in general, and are implausible in the context of the application. Joint independence and monotonicity are two such restrictions.

Conditional independence is a more plausible restriction than joint independence. Recall that joint independence nests conditional independence and marginal independence, and that it is marginal independence that is implausible in the context of the application. Further, it is plausible that partial identification of the average causal effect of the endogenous variable on the outcome variable can be maintained under conditional independence. In other words, it is plausible that joint independence is not an essential element of the model. I propose relaxing joint independence in favour of conditional independence as an extension.

Non-separability of  $h$  in  $U$  enriches individual response and is a more plausible restriction than monotonicity. To clarify, if  $h$  is non-separable in  $U$  then Restriction M1' is replaced by

$$Y = \mathbb{1}[q(D, X, U) < 0]$$

for some non-parametric function  $q$ . In general, non-separability of  $h$  in  $U$  permits the causal effect of  $(D, X)$  to vary with  $U$ . At the extreme, non-separability of  $h$  in  $U$  permits the causal

effect of  $(D, X)$  to be positive for some values of  $U$  and negative for others. In the context of the application, it is plausible that additional children lead some mothers to enter the labour market and some mothers to exit the labour market, and that this is dependent upon the opportunity that a mother faces. This behaviour is consistent with non-separability of  $h$  in  $U$ . I propose relaxing monotonicity in favour of non-separability as an extension.

Finally, the single equation model embeds the restriction that the codomain of  $g$  is  $\{0, 1\}$  but otherwise is silent as to the arguments and functional form of  $g$ . Suppose that this restriction is replaced with the restriction that  $g$  is a non-parametric function that varies with  $(U, X, Z)$ , and is non-separable in  $U$ . The resulting triangular model that embeds this restriction might have greater informational content than the single equation model with only a small loss in credibility. I propose strengthening the lack of constraints on  $g$  in favour of a non-parametric restriction as an extension.

## 2.2 Falsifiability

I show that the single equation model is falsifiable. I derive an instrumental inequality (Pearl, 1995) that is a sufficient condition for the single equation model to be observationally restrictive.

It is important that a model is falsifiable. If a model is not falsifiable then it cannot be rejected for any probability distribution of observable random variables that is consistent with the observable supports of these random variables. In other words, a model that is not falsifiable is always valid. The single equation model is falsifiable, and can be rejected if the following instrumental inequality is violated.

$$\max_{d \in \mathcal{R}_D} \max_{x \in \mathcal{R}_X} \sum_{y \in \mathcal{R}_Y} \max_{z \in \mathcal{R}_Z} \mathbb{P}(y, d|x, z) \leq 1 \quad (2.1)$$

*Proof.* I write  $\mathbb{P}(y, d|x, z)$  as

$$\int_{u \in \mathcal{R}_U} \mathbb{P}(y, d, u|x, z) dU, \quad (2.2)$$

which is valid by the Law of Total Probability. Further, I write the integrand as the product decomposition

$$\mathbb{P}(y|d, x, z, u) \mathbb{P}(d|x, z, u) \mathbb{P}(u|x, z), \quad (2.3)$$

which is valid by Bayes' Theorem. I postulate that  $\mathbb{P}(y, d|x, z)$  is generated by the single equation model. I write (2.3) as the product decomposition

$$\mathbb{P}(y|d, x, u) \mathbb{P}(d|x, z, u) \mathbb{P}(u), \quad (2.4)$$

which is valid by the restrictions that are embedded in the single equation model. Specifically, I use Restriction M3 and Restriction M4. I substitute (2.4) into (2.2) in place of the integrand.

$$\mathbb{P}(y, d|x, z) = \int_{u \in \mathcal{R}_U} \mathbb{P}(y|d, x, u) \mathbb{P}(d|x, z, u) \mathbb{P}(u) dU \quad (2.5)$$

I note that (2.5) holds for any  $z$ , and so holds for  $z(y, d, x)$  that I define as  $\operatorname{argmax}_{z \in \mathcal{R}_Z} \mathbb{P}(y, d|x, z)$ .

$$\mathbb{P}(y, d|x, z(y, d, x)) = \int_{u \in \mathcal{R}_U} \mathbb{P}(y|d, x, u) \mathbb{P}(d|x, z(y, d, x), u) \mathbb{P}(u) dU \quad (2.6)$$

I sum both sides of (2.6) over  $\mathcal{R}_Y$

$$\sum_{y \in \mathcal{R}_Y} \mathbb{P}(y, d|x, z(y, d, x)) = \sum_{y \in \mathcal{R}_Y} \int_{u \in \mathcal{R}_U} \mathbb{P}(y|d, x, u) \mathbb{P}(d|x, z(y, d, x), u) \mathbb{P}(u) dU$$

and write

$$\sum_{y \in \mathcal{R}_Y} \mathbb{P}(y, d|x, z(y, d, x)) = \sum_{y \in \mathcal{R}_Y} \int_{u \in \mathcal{R}_U} \mathbb{P}(y|d, x, u) \mathbb{P}(d|x, z(y, d, x), u) \mathbb{P}(u) dU.$$



I note that  $\mathbb{P}(d|x, z(y, d, x), u)$  is a probability, and so is bounded from above by unity.

$$\sum_{y \in \mathcal{R}_Y} \mathbb{P}(y, d|x, z(y, d, x)) \leq \sum_{y \in \mathcal{R}_Y} \int_{u \in \mathcal{R}_U} \mathbb{P}(y|d, x, u) \mathbb{P}(u) dU \quad (2.7)$$

I note that the right-hand side of (2.7) is an expectation, and so evaluates to a well-defined probability.

$$\sum_{y \in \mathcal{R}_Y} \mathbb{P}(y, d|x, z(y, d, x)) \leq \sum_{y \in \mathcal{R}_Y} \mathbb{P}(y|d, x) \quad (2.8)$$

I write

$$\sum_{y \in \mathcal{R}_Y} \mathbb{P}(y, d|x, z(y, d, x)) \leq 1,$$

which is valid by the Law of Total Probability. I write

$$\sum_{y \in \mathcal{R}_Y} \max_{z \in \mathcal{R}_Z} \mathbb{P}(y, d|x, z) \leq 1, \quad (2.9)$$

which is valid by the definition of  $z(y, d, x)$ . I note that (2.9) is constant for any  $(d, x)$ , and so holds for those values that maximise the left-hand side of the inequality. This completes the proof.  $\square$

I note that the instrumental inequality and its proof are adapted from Pearl (1995). I extend Pearl (1995) in that I permit the existence of an exogenous variable.

### 3 Identification

I introduce random set theory (Molchanov, 2005) as a tool for identification analysis. I discuss random set theory in the context of the single equation model. I exploit joint independence (Restriction M3) and monotonicity (Restriction M5), and derive sharp bounds on the distribution of  $U$ .

Artstein's Inequality (Artstein, 1983) is an important theorem of random set theory that is useful for deriving bounds on latent probability distributions. The usefulness of Artstein's Inequality is its dual representation as a capacity functional and as a containment functional, and that it defines a sharp set. Together these properties determine that identification analysis that is conducted using Artstein's Inequality defines the identified set of a functional of a latent probability distribution, as opposed to a proper superset (of the identified set). For a selection  $\xi$  from a random closed set  $\Xi$ ,

$$\mathbb{P}(\xi \in T|\mathcal{I}) \leq \mathbb{P}(\Xi \cap T \neq \emptyset|\mathcal{I}) \quad (3.1)$$

$$\mathbb{P}(\xi \in T|\mathcal{I}) \geq \mathbb{P}(\Xi \subseteq T|\mathcal{I}) \quad (3.2)$$

for all test sets  $T$  in the class of compact sets  $\mathcal{T}$ . Here, I write a conditional form (conditional on an arbitrary information set  $\mathcal{I}$ ) of Artstein's Inequality since the identification analysis that is conducted exploits the independence relations that the single equation model embeds. (3.1) is the capacity functional representation of Artstein's Inequality, and (3.2) is its containment functional representation.

In the context of the single equation model,  $(Y(0), Y(1))$  is a random closed set that has a well-defined probability distribution. I follow Chesher and Rosen (2013) in defining a level set  $\mathcal{U}_h$ .

$$\mathcal{U}_h(y, d, x) \equiv \{u : y = h(d, x, u)\}$$

The usefulness of  $\mathcal{U}_h$  to the identification analysis is its synonymity with  $(Y(0), Y(1))$ . This point is the subject of Figure (12). I translate Artstein's Inequality as

$$\mathbb{P}(U \in T|x, z) \leq \mathbb{P}(\mathcal{U}_h(Y, D, x) \cap T \neq \emptyset|x, z) \quad (3.3)$$

$$\mathbb{P}(U \in T|x, z) \geq \mathbb{P}(\mathcal{U}_h(Y, D, x) \subseteq T|x, z) \quad (3.4)$$

for all  $T$  in the class of all compact sets on  $[0, 1]$ , and for some  $(x, z) \in \mathcal{R}_X \times \mathcal{R}_X$ .<sup>6</sup> Further, I exploit joint independence and monotonicity, and write the left-hand side of (3.3) and of (3.4) as  $\|T\|$ .

The class of all compact sets on  $[0, 1]$  is a large class, and it is infeasible to compute Artstein's Inequality for all test sets in this class. Chesher et al. (2013) shows that it is sufficient to compute Artstein's Inequality for a smaller class of test sets, and names this class the class of core determining sets. For the single equation model, the class of core determining sets is the collection of

$$[0, p(d, x)] \text{ and } (p(d, x), 1]$$

over  $\mathcal{R}_D \times \mathcal{R}_X$ .<sup>7</sup> The left-hand side of (3.3) and of (3.4) is then either  $p(d, x)$  or  $1 - p(d, x)$ , depending upon the core determining set that is tested. Further, the class of core determining sets can be reduced to the collection of

$$[0, p(d, x)]$$

over  $\mathcal{R}_D \times \mathcal{R}_X$ . This refinement is valid since  $[0, p(d, x)]$  is the complement of  $(p(d, x), 1]$ , and so the capacity functional representation of Artstein's Inequality for  $[0, p(d, x)]$  is equivalent to the containment functional representation function of Artstein's Inequality for  $(p(d, x), 1]$ .

### 3.1 Identification analysis

I describe the identified set of the collection of  $p(d, x)$  over  $\mathcal{R}_D \times \mathcal{R}_X$ . I employ random set theory as a tool for the identification analysis. I distinguish between the identified set of a functional of a latent probability distribution and the admissible set of values for that functional.

Assume that there is a particular ordering of the collection of  $p(d, x)$  over  $\mathcal{R}_D \times \mathcal{R}_X$ . First, I define correspondences (set-valued functions)  $\mathcal{A}_p$  and  $\mathcal{B}_p$  as follows.

$$\begin{aligned}\mathcal{A}_p(\eta; d, x) &\equiv \{a : p(a, \eta) \leq p(d, x)\} \\ \mathcal{B}_p(\eta; d, x) &\equiv \{b : p(b, \eta) \geq p(d, x)\}\end{aligned}$$

I write  $\eta$  so as to emphasise the distinction between  $\eta$  as a conditioning value, and  $x$  as a determinant of a test set. That is, for a given test set  $[0, p(d, x)]$ , there is a conditional form of Artstein's Inequality for each  $\eta \in \mathcal{R}_X$ . Second, I describe the identified set of  $p(d, x)$  as

$$\sup_{\eta \in \mathcal{R}_X} \sup_{z \in \mathcal{R}_Z} \mathbb{P}(Y = 0, D \in \mathcal{A}_p(\eta; d, x) | \eta, z) \leq p(d, x) \leq 1 - \sup_{\eta \in \mathcal{R}_X} \sup_{z \in \mathcal{R}_Z} \mathbb{P}(Y = 1, D \in \mathcal{B}_p(\eta; d, x) | \eta, z), \quad (3.5)$$

which is valid by Artstein's Inequality. Note that alternative forms of Artstein's Inequality that rely only on marginal independence relations that are implied by joint independence also describe sets of  $p(d, x)$ , but that these sets constitute proper or improper supersets of (3.5). Further, (3.5) is defined for each possible ordering of the collection of  $p(d, x)$  over  $\mathcal{R}_D \times \mathcal{R}_X$ . If (3.5) is empty for a particular ordering of the collection of  $p(d, x)$  over  $\mathcal{R}_D \times \mathcal{R}_D$ , then that ordering can be rejected.

Generally, the sampling process does not identify  $\mathbb{P}$ . Instead, the sampling process identifies  $\mathbb{P}_n$  that is representative of a proper subset of the population, and that is often distinct from  $\mathbb{P}$ . To estimate (3.5), it is natural to replace  $\mathbb{P}$  with  $\mathbb{P}_n$ , which is valid by the assumption that the sampling process is informative of  $\mathbb{P}$  as  $n \rightarrow \infty$ . This assumption is known as the analogue principle (Manski, 1988). It is unnatural to regard the sample analogue of (3.5) as an identified set since identification is a concept of the population. Instead, I refer to the sample analogue of (3.5) as an admissible set to emphasise the distinction.

<sup>6</sup>The restriction on the class of test sets is valid by joint independence and monotonicity

<sup>7</sup>Notice that  $(p(d, x), 1]$  is an open set, but that Artstein's Inequality is defined for compact test sets. Strictly speaking, the class of core determining sets should include  $\text{cl}(p(d, x), 1]$  rather than  $(p(d, x), 1]$ . Since  $U$  is distributed continuously, there is zero mass at  $p(d, x)$  and so the measure of  $\text{cl}(p(d, x), 1]$  is equal to the measure of  $(p(d, x), 1]$ .

### 3.2 Dimensionality

I calculate the number of inequalities that describe the identified set, and how this number changes as  $\|\mathcal{R}_X\|$  increases. I show that there is a curse of dimensionality.

Tractability is a legitimate concern: regardless of whether incorporating information is advantageous in the context of the single equation model, the computational cost of incorporating information may be prohibitive. Let

$$\|\mathcal{R}_X\| = K \text{ and } \|\mathcal{R}_Z\| = L,$$

and fix a particular ordering of the collection of  $p(d, x)$  over  $\mathcal{R}_D \times \mathcal{R}_X$ . There are  $2K$  parameters in the collection of  $p(d, x)$  over  $\mathcal{R}_D \times \mathcal{R}_X$ , and there are  $2K \times L$  inequality relations for each parameter. In sum, the number of inequalities for each ordering of the collection of  $p(d, x)$  over  $\mathcal{R}_D \times \mathcal{R}_X$  is

$$4K^2 \times L.$$

The total number of orderings of the collection of  $p(d, x)$  over  $\mathcal{R}_D \times \mathcal{R}_X$  is then

$$\sum_{j=1}^{2K} j! \binom{2K}{j},$$

where  $\binom{n}{m}$  counts the number of ways that a set of length  $n$  can be partitioned into  $m$  non-empty subsets (the Stirling Number of the Second Kind). Each  $\binom{n}{m}$  is multiplied by  $m!$ , which accounts for the possible ordering over the  $m$  non-empty subsets. Combining the number of inequalities for a particular ordering of the collection of  $p(d, x)$  over  $\mathcal{R}_D \times \mathcal{R}_X$  with the number of possible orderings yields

$$4K^2 \times \sum_{j=1}^{2K} j! \binom{2K}{j} \times L$$

as the number of inequalities that describe the identified set. There is a clear curse of dimensionality. It is sensible to automate the task of calculating the inequalities that describe the identified set.

### 3.3 Average causal effects

I state the definition of  $ACE(D \rightarrow Y)$  when  $D$  is a binary variable. I discuss the relationship between  $ACE(D \rightarrow Y)$  and its conditional counterparts. I show that  $ACE(D \rightarrow Y)$  is expressible as a function of the parameters of the single equation model.

First,  $ACE(D \rightarrow Y)$  is defined as

$$\mathbb{E}[Y(1)] - \mathbb{E}[Y(0)],$$

which is commonly referred to as the average treatment effect in economics. For the single equation model,  $ACE(D \rightarrow Y)$  is equal to

$$\mathbb{E}[p(0, X)] - \mathbb{E}[p(1, X)]. \tag{3.6}$$

*Proof.* I focus on  $\mathbb{E}[Y(d)]$  since  $ACE(D \rightarrow Y)$  is an additive function of this expectation. To prove (3.6), I must show that

$$\mathbb{E}[Y(d)] = 1 - \mathbb{E}[p(d, X)].$$

First, I note that  $\mathbb{E}[Y(d)]$  is equivalent to  $\mathbb{P}(Y(d) = 1)$  since  $Y(d)$  is a binary variable. Substituting for  $Y(d)$  yields

$$\mathbb{P}(Y(d) = 1) = \mathbb{P}(p(d, X) < U).$$

Then,

$$\mathbb{P}(p(d, X) < U) = \sum_{x \in \mathcal{R}_X} \mathbb{P}(p(d, x) < U|x) \mathbb{P}(x), \quad (3.7)$$

which is valid by the Law of Iterated Expectations. The restrictions that the single equation model embeds determine that

$$\mathbb{P}(p(d, x) < U|x) = \mathbb{P}(p(d, x) < U),$$

and that

$$\mathbb{P}(p(d, x) < U) = 1 - p(d, x).$$

I rewrite (3.7) as

$$\mathbb{P}(p(d, X) < U) = 1 - \sum_{x \in \mathcal{R}_X} p(d, x) \mathbb{P}(x), \quad (3.8)$$

which is an expectation. This completes the proof of (3.6).  $\square$

Second,  $p(d, x)$  is a parameter of the single equation model. Although  $p$  is a non-parametric function it can be summarised by the collection of  $p(d, x)$  over  $\mathcal{R}_D \times \mathcal{R}_X$ . Further, knowledge of  $p$  (and of the distribution of  $X$ ) is sufficient to determine the distribution of  $Y(d)$ .  $ACE(D \rightarrow Y)$  is a function of the parameters of the single equation model, which is valid by (3.8).

Third,  $ACE(D \rightarrow Y|x)$  is defined as

$$\mathbb{E}[Y(1)|x] - \mathbb{E}[Y(0)|x],$$

which for the single equation model is equivalent to  $p(0, x) - p(1, x)$  (stated without proof). The relationship between  $ACE(D \rightarrow Y)$  and its conditional counterparts is then

$$ACE(D \rightarrow Y) = \sum_{x \in \mathcal{R}_X} ACE(D \rightarrow Y|x) \mathbb{P}(x), \quad (3.9)$$

which is established in the proof of (3.6). I exploit this relationship to describe the identified set of  $ACE(D \rightarrow Y)$ .

## 4 Incorporating information

I explore the effect of incorporating information for the single equation model. I study a special case of the single equation model that embeds the restriction that  $h$  excludes  $X$ . That is, I study a special case of the single equation model that embeds the restriction that

$$Y = \mathbb{1}[r(D) < U]$$

for a non-parametric function  $r$ . I name this special case of the single equation model the simple single equation model, and note that this model is assumed in Chesher and Rosen (2013) for the purpose of estimating the average causal effect of additional children on a mother's employment. I study how the admissible set of values of  $ACE(D \rightarrow Y)$ , equal to  $r(0) - r(1)$ , changes as

(a.) I enrich  $\mathcal{R}_Z$ ; and,

(b.) I permit  $X$  to enter  $h$ .

The model that I assume for (b) is then the single equation model, rather than the simple single equation model.

Chesher and Rosen (2013) does not report statistical uncertainty in the estimate of the average causal effect of additional children on a mother's employment. I extend Chesher and Rosen (2013) by constructing a confidence region for the estimate of the average causal effect of additional children on a mother's employment, and use this as a baseline for comparison in studying how the admissible set of values of the average causal effect of additional children on a mother's employment changes as information is incorporated into the simple single equation model. I use a method that is outlined in Chernozhukov et al. (2013) to compute statistical uncertainty. Although valid by Bonferroni's inequality, the statistical uncertainty that I report is conservative. That is, for some statistical size  $\alpha$ , the confidence region covers the average causal effect of additional children on a mother's employment in at least  $1 - \alpha$  samples. I discuss the inferential problem in Figure 6.

In the context of the application, I consider two definitions for  $Z$  in the framework of the simple single equation model. This approach is consistent with Chesher and Rosen (2013). In the first instance, I define  $Z$  as an indicator for the event that the oldest two children in the household share the same gender. I refer to this instance as Experiment 1. In the second instance, I define  $Z$  as an indicator for the event that the second birth is a multiple birth. I refer to this instance as Experiment 6.

$$\begin{aligned} \text{Experiment 1} \quad Z &= \begin{cases} 1 & \text{Male-Male} \cup \text{Female-Female}, \\ 0 & \text{Male-Female} \cup \text{Female-Male}. \end{cases} \\ \text{Experiment 6} \quad Z &= \begin{cases} 1 & \text{Multiple birth}, \\ 0 & \text{Single birth}. \end{cases} \end{aligned}$$

Figure 9 and Figure 10 show the admissible set of values of  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$  for Experiment 1 and Experiment 6, respectively.

In Experiment 1, the admissible set of values of  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$  is large and disconnected. The gender of the oldest two children is not a good predictor of whether there are three or more children in the household.

$$\mathbb{P}_n(D = 1|Z = 1) = 0.414 \text{ versus } \mathbb{P}_n(D = 1|Z = 0) = 0.346$$

There is little statistical uncertainty in the admissible set; the intersection of the admissible set and the confidence region for the admissible set is large relative to the confidence region for the admissible set. There is little statistical uncertainty in the admissible set as

- (a) the event that the oldest two children in the household share the same gender is approximately as frequent as the event that the oldest two children in the household do not share the same gender; and,
- (b) the gender of the oldest two children is not a good predictor of whether there are three or more children in the household.

Specifically, (a) means that each of the inequalities that describe the admissible set is a functional of a large number of data, and so there should be little statistical uncertainty in each of these inequalities. Further, (b) means that there is little difference between those inequalities that are binding constraints, and those that are not binding constraints.

In Experiment 6, the admissible set of values of  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$  collapses to a line, which covers a medium interval. The admissible set collapses due to the nature of the relationship between whether there are three or more children in the household and the event that the second birth is a multiple birth: clearly, if there is a multiple second birth then there must be three or more children in the household. The event that there is a multiple second birth is a special instrumental variable as it excludes some events in the probability space.

$$\mathbb{P}_n(D = 1|Z = 1) = 1 \text{ versus } \mathbb{P}_n(D = 1|Z = 0) = 0.375$$

An implication of its special status is that  $\mathbb{E}_n[Y(1)]$  is identified. Essentially, if there is a multiple second birth then there is no selection and so  $\mathbb{E}_n[Y(1)]$  is identified from the sub-population of mothers that experience a multiple second birth. There is considerable statistical uncertainty in the

admissible set; the intersection of the admissible set and the confidence region for the admissible set is small relative to the confidence region for the admissible set. In particular, there is considerable uncertainty in the  $\mathbb{E}_n[Y(1)]$ -direction of the admissible set of values of  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$  as

- (a) the event that there is a multiple second birth is rare; and,
- (b) the event that there is a multiple second birth is a good predictor of whether there are three or more children in the household.

In fact, whether the instrumental variable is a good predictor of the endogenous variable is of first-order importance only for estimation, and not for inference. The difference between those inequalities that are binding constraints, and those that are not binding constraints only becomes important if the statistical uncertainty in those inequalities that are binding constraints is considerable. In other words, if the instrumental variable is a rare event. In sum, the ideal instrumental variable is a good predictor of the endogenous variable and is approximately as frequent as not.

Figure 8 shows the admissible set of values of  $ACE_n(D \rightarrow Y)$  for Experiment 1 and for Experiment 6 (and all other experiments that I conduct). In Experiment 1, the admissible set is large and disconnected. The event that the oldest two children in the household share the same gender is insufficient to sign  $ACE_n(D \rightarrow Y)$ . Further, either side of zero the admissible set is not informative beyond what anecdotal evidence might suggest, and certainly cannot inform policy. The only information that Experiment 1 conveys with respect to  $ACE_n(D \rightarrow Y)$  is that there is a statistically significant effect (since the confidence region for the admissible set does not cover zero). In Experiment 6, the admissible set is not disconnected. The event that there is a multiple second birth is sufficient to sign  $ACE_n(D \rightarrow Y)$  as negative. Although Experiment 6 is informative with respect to the sign of  $ACE_n(D \rightarrow Y)$ , it conveys little other information to inform policy.

## 4.1 Enriching the support of the instrumental variable

I explore the effect of enriching  $\mathcal{R}_Z$  for the simple single equation model in the context of the application. Specifically, I disaggregate the event that the oldest two children in the household share the same gender into the possible permutations of gender for two children. Further, I combine the event that there is a multiple second birth with the possible permutations of gender for two children. I show that incorporating informational is advantageous.

In all, I consider eight experiments. Each experiment is characterised by a different definition of  $Z$ . The results of these experiments are the subject of Table 1 through Table 4, and the tables also give the definition of  $Z$  in each experiment. I make the following observations.

- (a) Disaggregating the event that the oldest two children in the household share the same gender into either the event that the oldest two children in the household are male, or into the event that the oldest two children in the household are female has less informational content.

In Experiment 2 I disaggregate the event that the oldest two children in the household share the same gender into the event that the oldest two children are male. In Experiment 3 I disaggregate the event that the oldest two children in the household share the same gender into the event that the oldest two children are female. In both experiments, the admissible set of values of  $ACE_n(D \rightarrow Y)$  is larger than the admissible set for Experiment 1. Further, there is greater statistical uncertainty in the admissible set for Experiment 2 and for Experiment 3.

Either disaggregation maintains  $Z$  as a binary variable and, if preferences are convex, aggregates an event that the oldest two children in the household share the same gender with the event that the oldest two children do not share the same gender. For example, in Experiment 2 the disaggregation aggregates the event that the oldest two children in the household are female with the event that the oldest two children do not share the same gender. Importantly, this dilutes the identifying power of  $Z$  (since the event that the oldest two children in the household share the same gender is associated with both values of  $Z$ ). Further, each value of  $Z$  is no longer as frequent as not. The imbalance in the sub-populations leads to greater statistical uncertainty.

- (b) Disaggregating the event that the oldest two children in the household share the same gender into the possible permutations of gender for two children has more informational content.

In Experiment 4 I disaggregate the event that the oldest two children in the household share the same gender into the event that the oldest two children are male, the event that the oldest two children are female, and the event that the oldest two children do not share the same gender. In Experiment 5 I disaggregate the event that the oldest two children in the household share the same gender into the possible permutations of gender for two children. In both experiments, the admissible set of values of  $ACE_n(D \rightarrow Y)$  is smaller than the admissible set for Experiment 1. There is greater statistical uncertainty in the admissible set for Experiment 4 and for Experiment 5, but the net effect of greater accuracy and less precision is to make the admissible set and the confidence region of the admissible set smaller in both cases. The admissible set of values of  $ACE_n(D \rightarrow Y)$  is smaller than the admissible set for Experiment 1.

Most importantly, either disaggregation increases  $\|\mathcal{R}_Z\|$ . The advantage of disaggregating the event that the oldest two children in the household is that more variation is introduced that might be exploited. To illustrate, notice that

$$\mathbb{P}_n(A|\text{Male-Male} \cup \text{Female-Female}) \approx \frac{1}{2}\mathbb{P}_n(A|\text{Male-Male}) + \frac{1}{2}\mathbb{P}_n(A|\text{Female-Female})$$

for some event  $A$ . If

$$\mathbb{P}_n(A|\text{Male-Male}) \neq \mathbb{P}_n(A|\text{Female-Female})$$

(rather, the probabilities are not sufficiently close) then there is advantage to conditioning solely on one of these events. The reason is that one of the probabilities will constitute a more binding constraint. If preferences are convex but also exhibit gender bias then such an effect might be observed. For example, if

$$\text{Male-Male-Unknown} > \text{Female-Female-Unknown}$$

in expectation then the probability that there are three or more children in the household is different depending upon whether the oldest two children in the household are male, or are female. Such an effect might lead to an imbalance in the conditional probabilities. The disaggregation does not dilute the identifying power of  $Z$ .

- (c) Combining the event that there is a multiple second birth with the possible permutations of gender for two children has more informational content.

In Experiment 8 I combine the event that there is a multiple second birth with the possible permutations of gender for two children. There is greater statistical uncertainty in the admissible set for Experiment 8, but the net effect of greater accuracy and less precision is to make the admissible set and the confidence region of the admissible set smaller than the admissible set for Experiment 6.

Whether there are three or more children is dependent upon the quality margin (if preferences are convex then whether the oldest two children in the household share the same gender) and the quantity margin (whether there is a multiple second birth), and so combining information on both margins yields a more complete picture of the relationship. Some caution must be taken in combining information on both margins; namely, that the quantity margin has primacy over the quality margin. For example, if the oldest child in the household is male and there is a multiple second birth say, then whether the second oldest child is male or female is not relevant to whether there are three or more children in the household. A multiple second birth implies that there are three or more children in the household, regardless. The advantage of incorporating information on the quality margin is that quality is relevant when there is not a multiple second birth.

Incorporating information is advantageous, but (a) emphasises that simple redefinition of  $Z$  without enriching  $\mathcal{R}_Z$  is insufficient for a model to have more informational context. To enrich  $\mathcal{R}_Z$ , it is essential that  $\|\mathcal{R}_Z\|$  increases such that there is greater variation to exploit. Although the gains from enriching  $\mathcal{R}_Z$  in conducting estimation and inference are not substantial, the cost of incorporating information is low since the number of inequalities increases linearly in  $\mathcal{R}_Z$ . If the task of calculating the inequalities that describe the admissible set is automated then there is

little reason not to incorporate information. Further, the availability of information on so many events implies that the single equation model is overidentified, and so its credibility is testable (to an extent).

## 4.2 Enriching individual response

I explore the effect of permitting  $X$  to enter  $h$  for the simple single equation model. I discuss the implications for individual behaviour of assuming the simple single equation model versus assuming the single equation model, and discuss the misspecification issue that arises when individual behaviour is dependent upon  $X$  but the single equation model is assumed.

The principal difference between the simple single equation model and the single equation model is the treatment of  $X$ . The simple single equation model embeds the restriction that individual behaviour is not dependent upon  $X$ .<sup>8</sup> In the context of the application, the restriction implies that race does not affect a mother's employment. Cultural norms or discriminatory policies that reduce the incentive to seek employment for mothers of a particular race say, independently of ability, are excluded by the simple single equation model. Supposing that the restriction that  $U \perp X$  is maintained (a reasonable restriction given the definition of  $X$ ) then the interpretation of  $X$  for the simple single equation model is

- (a) as an exogenous variable that is irrelevant to the determination of any component of the structural model; or,
- (b) as an instrumental variable if  $X$  is figured to be a determinant of  $D$ , or is a good predictor of  $D$ .

In either case, if individual behaviour is dependent upon  $X$  but the simple single equation model is assumed, the restriction that  $U \perp X$  is no longer valid. Note that

$$Y = \mathbb{1}[p(D, X) < U]$$

can always be written as

$$Y = \mathbb{1}[\zeta(D, X) < \lambda(D, V)] \quad (4.1)$$

if the restrictions that the single equation model embeds are maintained. I define a function  $\pi$  such that

$$F_{v|d}(v < \pi(d)|d) = r(d),$$

where  $F_v$  is the distribution function of some as yet undefined variable.<sup>9</sup> Adding and subtracting  $\pi(D)$  inside the indicator of (4.1) yields

$$Y = \mathbb{1}[\pi(D) < \lambda(D, V) + \pi(D) - \zeta(D, X)],$$

which upon defining

$$v \equiv \lambda(D, V) + \pi(D) - \zeta(D, X)]$$

demonstrates why the restriction that  $U \perp X$  is no longer valid ( $v$  is a function of  $X$ ). The interpretation of the bias term  $\zeta(D, X)$  is then as omitted variable bias in the case of (a), but more worryingly as a violation of the exclusion restriction if  $X$  is re-branded as an instrumental variable as in the case of (b).

<sup>8</sup>In fact, individual behaviour is permitted to depend upon  $X$  provided that this dependency is in the aggregate. In other words, if only the aggregate of  $X$  matters for individual behaviour, and so is constant across all economic agents.

<sup>9</sup>I depart from the notation that is employed throughout the rest of the paper in order to emphasise the variable to which this distribution corresponds.



### 4.2.1 Identification analysis

I describe the identified set of the collection of  $r(d)$  over  $\mathcal{R}_D$  for the simple single equation model. I employ random set theory as a tool for the identification analysis, and proceed as for the identification analysis of the single equation model.

Assume that there is a particular ordering of the collection of  $r(d)$  over  $\mathcal{R}_D$ . First, I define correspondences  $\mathcal{A}_r$  and  $\mathcal{B}_r$  as follows.

$$\begin{aligned}\mathcal{A}_r(d) &\equiv \{a : r(a) \leq r(d)\} \\ \mathcal{B}_r(d) &\equiv \{b : r(b) \geq r(d)\}\end{aligned}$$

Second, I describe the identified set of  $r(d)$  as

$$\sup_{z \in \mathcal{R}_Z} \mathbb{P}(Y = 0, D \in \mathcal{A}_r(d)|z) \leq r(d) \leq 1 - \sup_{z \in \mathcal{R}_Z} \mathbb{P}(Y = 1, D \in \mathcal{B}_r(d)|z), \quad (4.2)$$

which is valid by Artstein's Inequality.

### 4.2.2 Model misspecification and bias

I discuss the bias that arises when individual behaviour is dependent upon  $X$  but the simple single equation model is assumed. I show that the simple single equation model does not identify a superset of the identified set if the restrictions that the single equation model embeds are maintained (model misspecification).

If the simple single equation model is misspecified then it does not (partially) identify  $\mathbb{E}[Y(d)]$ . Rather, the simple single equation model (partially) identifies the functional  $r(d)$  that encompasses a bias. A natural question is then how does the identified set of  $r(d)$  compare to the identified set of  $\mathbb{E}[Y(d)]$ , which is identified by the single equation model? I focus on the lower set of inequalities that describe the identified set in each case since the analysis is sufficient to address the question. I expand the left-hand inequality of (4.2) to

$$\sup_{z \in \mathcal{R}_Z} \sum_{x \in \mathcal{R}_X} \mathbb{P}(Y = 0, D \in \mathcal{A}_r(d)|x, z) \mathbb{P}(x|z), \quad (4.3)$$

which is valid by the Law of Total Expectation. Contrast this with the lower bound that is identified by the single equation model.

$$\sum_{x \in \mathcal{R}_X} \mathbb{P}(x) \sup_{\eta \in \mathcal{R}_X} \sup_{z \in \mathcal{R}_Z} \mathbb{P}(Y = 0, D \in \mathcal{A}_p(\eta; d, x)|\eta, z) \quad (4.4)$$

I study the differences between (4.3) and (4.4), emphasising each difference that I study in turn to facilitate comparison. I study the partial (*ceteris paribus*) effect of each difference.

- (a) First, the supremum operator in (4.4) guarantees that the conditional probability  $\mathbb{P}(\cdot|\eta, z)$  is at least as great as in (4.3).

$$\sum_{x \in \mathcal{R}_X} \mathbb{P}(x) \sup_{\eta \in \mathcal{R}_X} \sup_{z \in \mathcal{R}_Z} \mathbb{P}(Y = 0, D \in \mathcal{A}_p(\eta; d, x)|\eta, z)$$

This difference determines that the bias from model misspecification is a downwards bias.

- (b) Second, the location of the supremum operator in (4.4) guarantees that the conditional probability  $\mathbb{P}(\cdot|\eta, z)$  is at least as great as in (4.3).

$$\sum_{x \in \mathcal{R}_X} \mathbb{P}(x) \sup_{\eta \in \mathcal{R}_X} \sup_{z \in \mathcal{R}_Z} \mathbb{P}(Y = 0, D \in \mathcal{A}_p(\eta; d, x)|\eta, z)$$

It is somewhat misleading to study this difference independently of the weighting function but it is nonetheless enlightening to study the partial effect of moving the supremum operator inside the sum. If the difference between the probability weights is ignored then moving the supremum operator inside the sum determines that the bias from model misspecification is a downwards bias.

(c) Third, the probability weights in (4.4) are unconditional probabilities.

$$\sum_{x \in \mathcal{R}_X} \mathbb{P}(x) \sup_{\eta \in \mathcal{R}_X} \sup_{z \in \mathcal{R}_Z} \mathbb{P}(Y = 0, D \in \mathcal{A}_p(\eta; d, x) | \eta, z)$$

Together with moving the supremum operator inside the sum, the difference in probability weights determines that the bias from model misspecification is ambiguous.

(d) Fourth,  $\mathcal{A}_r$  is not necessarily equal to  $\mathcal{A}_p$ .

$$\sum_{x \in \mathcal{R}_X} \mathbb{P}(x) \sup_{\eta \in \mathcal{R}_X} \sup_{z \in \mathcal{R}_Z} \mathbb{P}(Y = 0, D \in \mathcal{A}_p(\eta; d, x) | \eta, z)$$

Notice that the set  $\mathcal{A}_r$  is constant across the summation, but that  $\mathcal{A}_p$  can be different in each term of the summation. Further, it is unclear whether  $\mathcal{A}_r$  and  $\mathcal{A}_p$  are equal in any term of the summation. This difference determines that the bias from model misspecification is ambiguous.

Overall, the differences between (4.3) and (4.4) determine that the bias from model misspecification is ambiguous; the bias from (c) and (d) can outweigh the bias from (a) and (b). Extending this finding to the upper bound that is identified by the single equation model, the effect of permitting  $X$  to enter  $h$  is to alter the length of the identified set, to shift the location of the identified set, or a combination of both. With respect to location shift, the finding is consistent with the behaviour of identifying models.

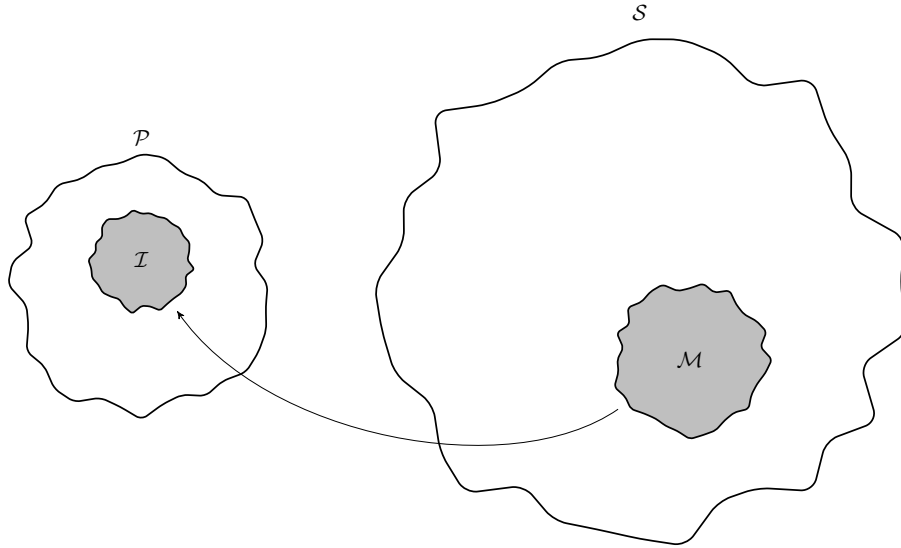
#### 4.2.3 Hispanic ethnicity

- Add in  $X$  to  $h$ .
- Context of application.
- Define  $X$  as event that a mother is Hispanic;  $Z$  is multi2nd.
- Of 75 possible orderings there is only one possible ordering that is admitted (72 are not internally consistent - lower bound greater than upper bound on at least one of the parameters; and 2 of the 3 non-excluded orderings do not respect their ordering over the parameters).
- Ordering that is preserved is

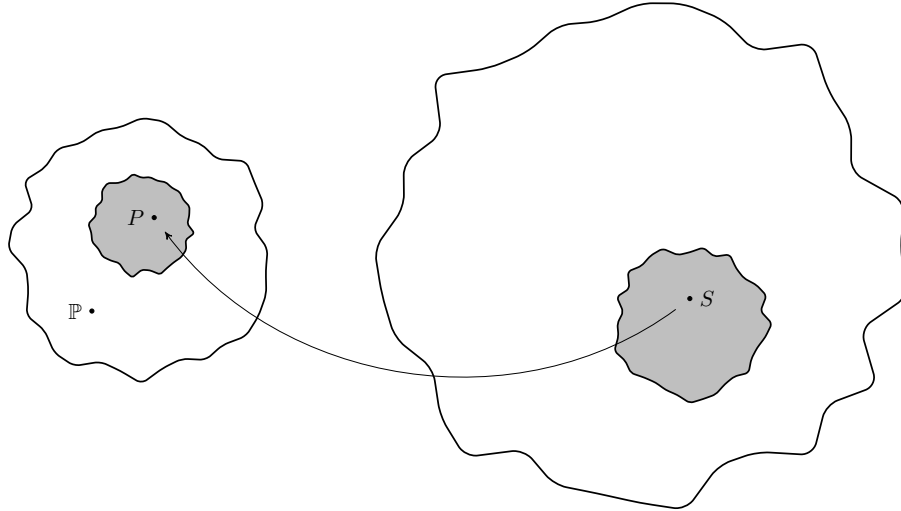
$$p(0, 0) = p(0, 1) < p(1, 0) < p(1, 1),$$

or that the conditional average causal effect is greater for Hispanics than for non-Hispanics.

- $[0.272, 0.470] < 0.521 < 0.555$  - again multi2nd point identifies  $p(1, x)$ .
- Probability of Hispanic is 0.074.
- $ACE_n(D \rightarrow Y|1) \in [-0.250, -0.052]$  and  $ACE_n(D \rightarrow Y|0) \in [-0.283, -0.085]$ .
- Minkowski sum of the two sets weighted by the probability of Hispanic gives  $ACE(D \rightarrow Y) \in [-0.252, -0.054]$ .
- Gives a narrower set than the original admissible set without  $X$ .

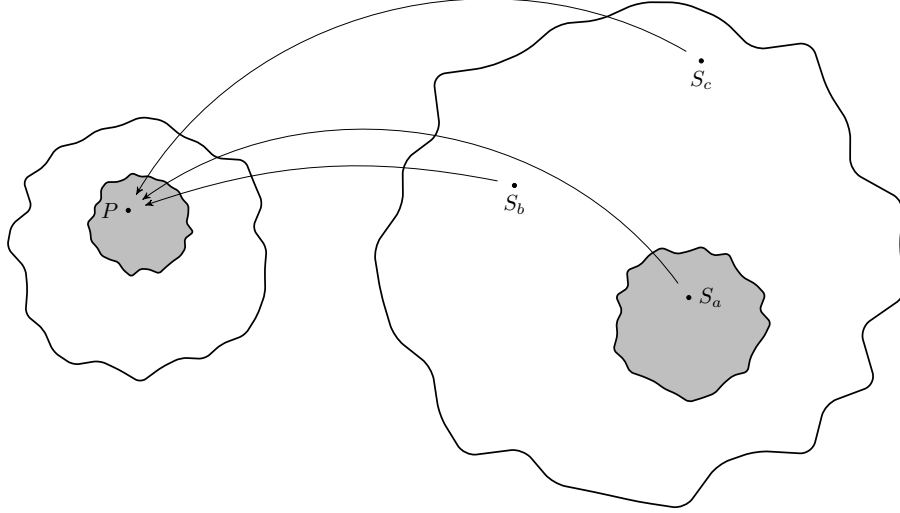


(a) A model  $\mathcal{M}$  is a set of structures that forms a proper subset of the class of all structures  $\mathcal{S}$ . Each structure in  $\mathcal{M}$  generates a probability distribution in the class of all probability distributions (of observable variables)  $\mathcal{P}$ . Then the image  $\mathcal{I}$  is the set of all probability distributions that are generated by structures in  $\mathcal{M}$ .

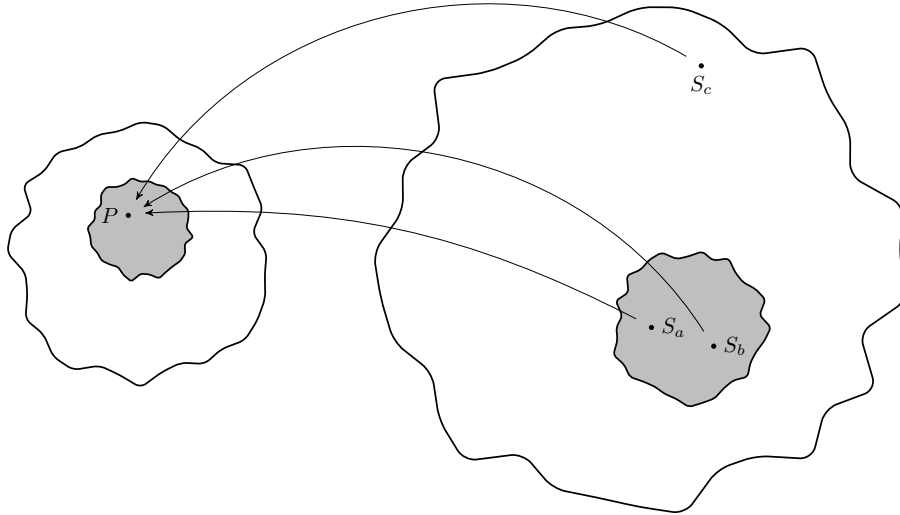


(b) A structure  $S$  is incompatible with data if it generates a probability distribution (of observable variables)  $P$  that is distinct from a realised probability distribution  $\mathbb{P}$ . If all structures in  $\mathcal{M}$  are incompatible with data then  $\mathcal{M}$  is said to be observationally restrictive, and is falsified. This condition is equivalent to  $\mathbb{P} \in \mathcal{P} \setminus \mathcal{I}$ .

Figure 1: Structures, models, probability distributions (of observable variables), and falsifiability.

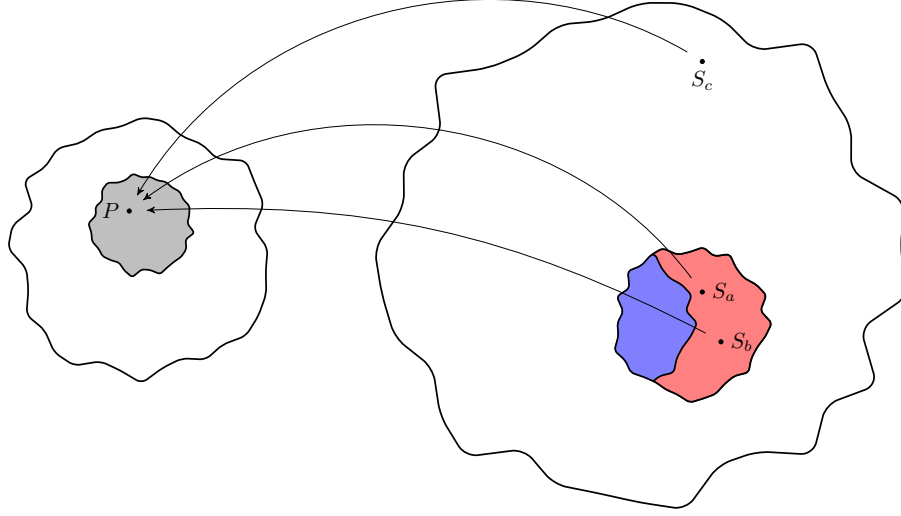


(a) A model  $\mathcal{M}$  is said to identify a structure  $S$  if the probability distribution (of observable variables)  $P$  that is generated by  $S$  is distinct from those generated by other structures in  $\mathcal{M}$ . The structures  $S_a$ ,  $S_b$  and  $S_c$  are said to be observationally equivalent as they all generate  $P$  but  $S_b$  and  $S_c$  are not admitted by  $\mathcal{M}$ . As  $S_a$  is the only structure that is admitted by  $\mathcal{M}$  and that generates  $P$ ,  $S_a$  is identified by  $\mathcal{M}$ . For completeness,  $\mathcal{M}$  is said to be uniformly identifying if it identifies each structure that it admits.

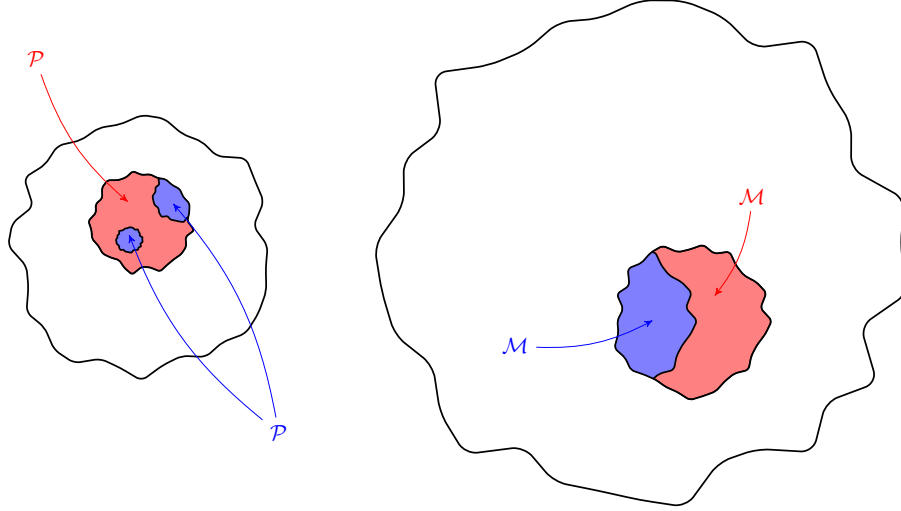


(b) As  $S_a$  and  $S_b$  are observationally equivalent and are both admitted by  $\mathcal{M}$  then  $\mathcal{M}$  does not identify either  $S_a$  or  $S_b$ . Nonetheless, as  $\mathcal{M}$  restricts the set of observationally equivalent structures that generate  $P$  to  $S_a$  and  $S_b$  then  $\mathcal{M}$  partially identifies  $S_a$  (and  $S_b$  to within  $\{S_a, S_b\}$ ).

Figure 2: Identification and non-identification of a structure, and partial identification of a structure.



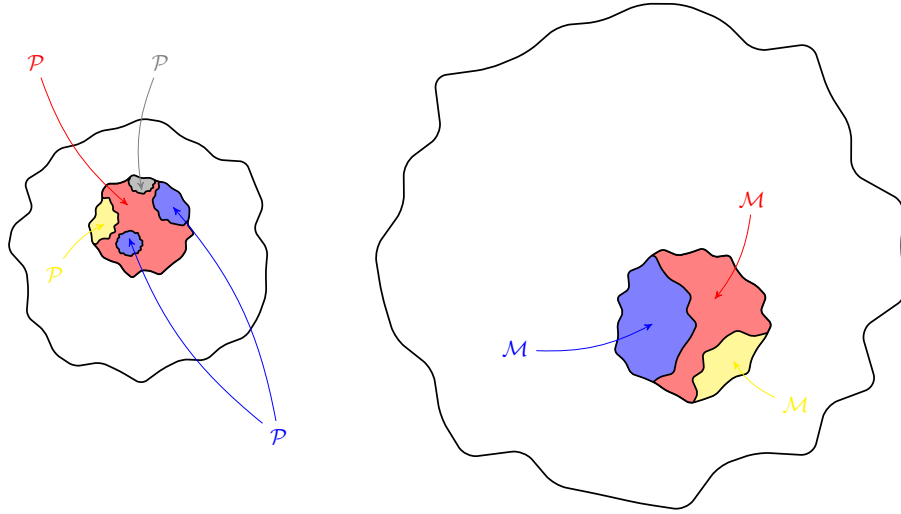
(a) A structural characteristic  $\chi$  is a function of a structure  $S$ . A model  $\mathcal{M}$  can be partitioned such that structures in a partition deliver the same value for  $\chi$ . Structures in the red partition  $\mathcal{M}$  deliver the value  $a$  for  $\chi$ , and structures in the blue partition  $\mathcal{M}$  deliver the value  $b$  for  $\chi$ . If  $\chi$  is constant across all observationally equivalent structures that  $\mathcal{M}$  admits then  $\mathcal{M}$  is said to identify  $\chi$ . As  $\chi(S_a)$  is equal to  $\chi(S_b)$  (is equal to  $a$ )  $\mathcal{M}$  identifies  $\chi$ .



(b) If  $\mathcal{M}$  identifies  $\chi$  for all structures in  $\mathcal{M}$  then  $\mathcal{M}$  is said to uniformly identify  $\chi$ . The class of all probability distributions (of observable variables) is partitioned into the red partition  $\mathcal{P}$  and into the blue partition  $\mathcal{P}$ . Probability distributions in  $\mathcal{P}$  are generated by (potentially many) structures in  $\mathcal{M}$ , and probability distributions in  $\mathcal{P}$  are generated by (potentially many) structures in  $\mathcal{M}$ . It is important that the number of partitions in  $\mathcal{M}$  and in  $\mathcal{P}$  are equal, although that number can be countably infinite. In the context of Figure 3b  $\mathcal{M}$  uniformly identifies  $\chi$  since observationally equivalent structures that  $\mathcal{M}$  admits are in the same colour of  $\mathcal{M}$ . More conveniently, whether  $\mathcal{M}$  uniformly identifies  $\chi$  can be determined by the existence of an identifying correspondence  $G$ , a functional.  $\mathcal{P}$  is a probability distribution in  $\mathcal{P}$ , and  $\mathcal{P}$  is a probability distribution in  $\mathcal{P}$ . Then  $\mathcal{M}$  uniformly identifies  $\chi$  if the value of  $G(\mathcal{P})$  is  $a$  and if the value of  $G(\mathcal{P})$  is  $b$ , holding for any such  $\mathcal{P}$  and  $\mathcal{P}$ . Notice that if  $\mathcal{M}$  uniformly identifies all  $\chi$  then  $\mathcal{M}$  also uniformly identifies structures.

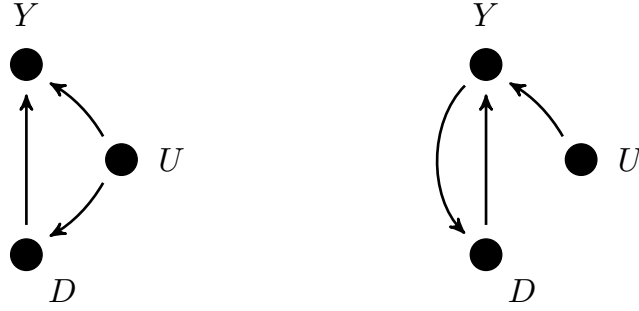
Figure 3: The identification of structural characteristics, and identifying correspondences.

(a) A structural characteristic  $\chi$  is a function of a structure  $S$ . A model  $\mathcal{M}$  can be partitioned such that structures in a partition deliver the same value for  $\chi$ . Structures in the red partition  $\mathcal{M}$  deliver the value  $a$  for  $\chi$ , structures in the blue partition  $\mathcal{M}$  deliver the value  $b$  for  $\chi$ , and structures in the yellow partition  $\mathcal{M}$  deliver the value  $c$  for  $\chi$ . The class of all probability distributions (of observable variables)  $\mathcal{P}$  is partitioned into the red partition  $\mathcal{P}$ , into the blue partition  $\mathcal{P}$ , into the yellow partition  $\mathcal{P}$  and into the grey partition  $\mathcal{P}$ . Probability distributions in a colour of  $\mathcal{P}$  are generated by (potentially many) structures in the same colour of  $\mathcal{M}$ ; the exception is probability distributions in  $\mathcal{P}$  which are generated by (potentially many) structures in  $\mathcal{M}$  and in  $\mathcal{M}$ .  $P$  is a probability distribution in  $\mathcal{P}$  with probability distributions defined similarly for each colour in  $\mathcal{P}$ .

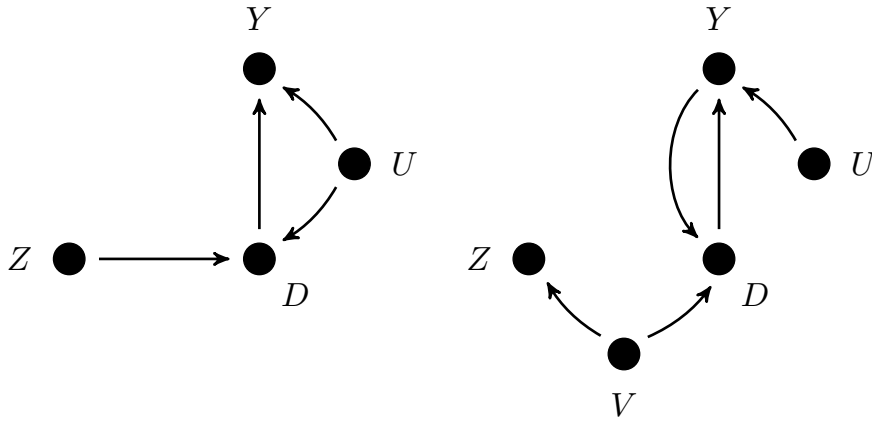


(b) That probability distributions in  $\mathcal{P}$  are generated by structures in  $\mathcal{M}$  and in  $\mathcal{M}$  creates a complication; the value of  $\chi$  is not constant across observationally equivalent structures that  $\mathcal{M}$  admits and that generate a probability distribution in  $\mathcal{P}$ . So  $\mathcal{M}$  does not uniformly identify  $\chi$ . Consideration of the identifying correspondence  $G$  determines that this corresponds to there being structures in  $\mathcal{M}$  for which  $G$  does not deliver the value of  $\chi$  when applied to the probability distributions that these structures generate. Nonetheless, if  $\mathcal{M}$  restricts the set of values of  $\chi$  for any probability distribution in  $\mathcal{P}$  then  $\mathcal{M}$  does have some non-trivial identifying power for  $\chi$ . Then  $\mathcal{M}$  is said to uniformly partially identify  $\chi$  if  $\mathcal{M}$  and  $\mathcal{P}$  can each be partitioned into countably many disjoint subsets and that a probability distribution in a partition of  $\mathcal{P}$  is not generated by a structure in at least one partition of  $\mathcal{M}$ , holding for any such partition of  $\mathcal{P}$ . In the context of Figure 4  $\mathcal{M}$  identifies  $\chi$  up to  $\{a, c\}$ ,  $\mathcal{M}$  identifies  $\chi$  uniquely to  $b$ , and  $\mathcal{M}$  identifies  $\chi$  up to  $\{a, c\}$ . Each partition of  $\mathcal{P}$  includes probability distributions that are generated by structures in at least one partition of  $\mathcal{M}$ . Equivalently, if  $G$  is permitted to be a multivalued functional (or one-to-many) then  $\mathcal{M}$  uniformly partially identifies  $\chi$  if  $G$  exists and if  $G(P)$  contains the set of values of  $\chi$  that are delivered by structures that generate  $P$ , holding for all such  $P$ . A caveat must be applied here;  $G$  cannot be trivial in the sense that it is constant across all such  $P$ . Clearly this definition of  $G$  does not exclude the possibility that there is multiplicity of identifying correspondences that satisfy this property. Sharpness is a desirable property in such circumstances; a functional  $G$  that can be shown to deliver smaller sets according to some well-defined distance measure across all possible  $P$  (and that satisfies the properties above) should be preferred to any alternative identifying correspondence.

Figure 4: Partial identification of a structural characteristic.



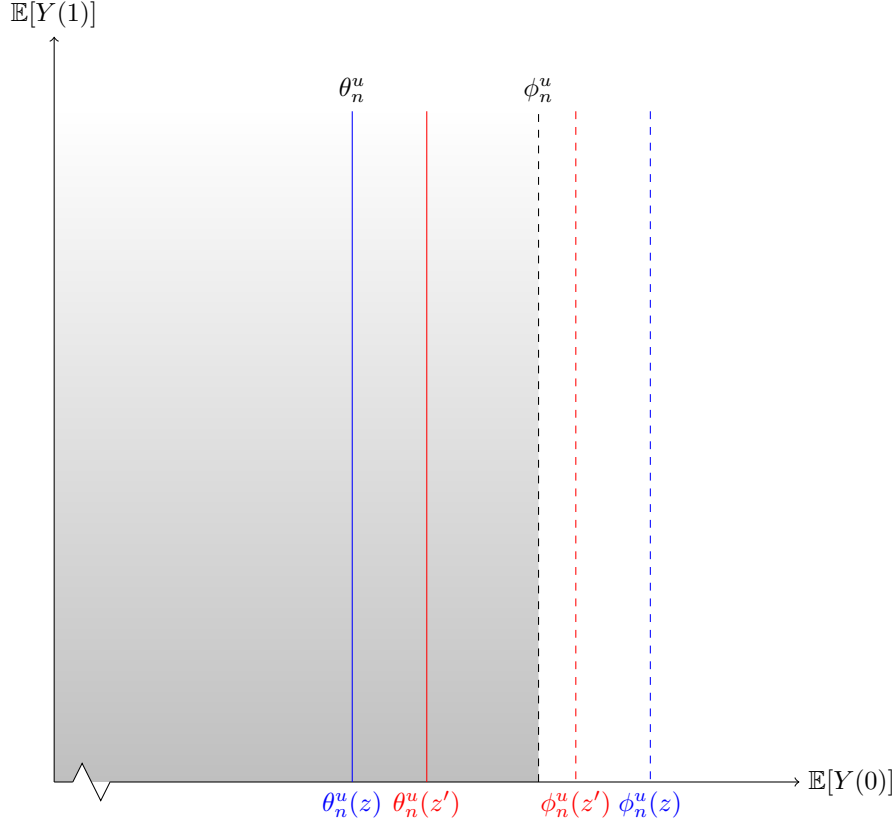
(a) The left- and right-hand panels show directed acyclic graphs that are able to represent the same probability distribution of  $(Y, D)$ . In the left-hand panel,  $U$  causes  $Y$  and  $D$ . The left-hand panel is a representation of selection in that particular values of  $D$  are strongly associated with particular values of  $Y$  independently of the causal effect of  $D$ . In the right-hand panel,  $U$  causes  $Y$  but has an indirect effect on  $D$  through  $Y$ . The right-hand panel is a representation of simultaneity in that  $Y$  is both a cause of and an effect of  $D$ . The equivalence between the two directed acyclic graphs is that  $Y$  can always be written as  $\xi(D, U)$  for  $\xi : \mathcal{R}_D \times \mathcal{R}_U \rightarrow \mathcal{R}_Y$  (provided that the right-hand panel converges to an equilibrium).



(b) To recover the causal effect of  $D$  on  $Y$ , it is necessary that there exists an external and measurable factor that causes variation in  $D$ . This external and measurable factor is known as an instrumental variable. In the left-hand panel,  $Z$  causes  $D$ . It is convenient to think of  $Z$  as a switch that forces  $D$  to take particular values. The difference between the value of  $Y$  when  $Z$  is **on** versus when  $Z$  is **off** is the causal effect of  $D$  on  $Y$ . In the right-hand panel, it is  $V$  that causes  $D$  ( $V$  is unobservable). As  $V$  causes  $D$  and  $Z$ , it may be sufficient to look at  $Z$  to measure exogenous variation in  $D$  (although it is not always). As such, the causal effect of  $D$  on  $Y$  is recoverable using variation in  $Z$ . This is an important point about the nature of an instrumental variable; namely, that the relationship between  $D$  and  $Z$  need not be causal.

Figure 5: A note on causality.

(a) Suppose that  $\theta_n^u(z)$  and  $\theta_n^u(z')$  are estimates of upper bounds on  $\mathbb{E}[Y(0)]$ . Similarly, suppose that  $\phi_n^u(z)$  and  $\phi_n^u(z')$  are one-sided  $1 - \alpha$  confidence regions for  $\theta_n^u(z)$  and  $\theta_n^u(z')$ .  $\phi_n^u(z') > \phi_n^u(z)$  if there is greater variation in the estimate of  $\theta_n^u(z)$  (if there are fewer observations of  $z$  than  $z'$ , say).

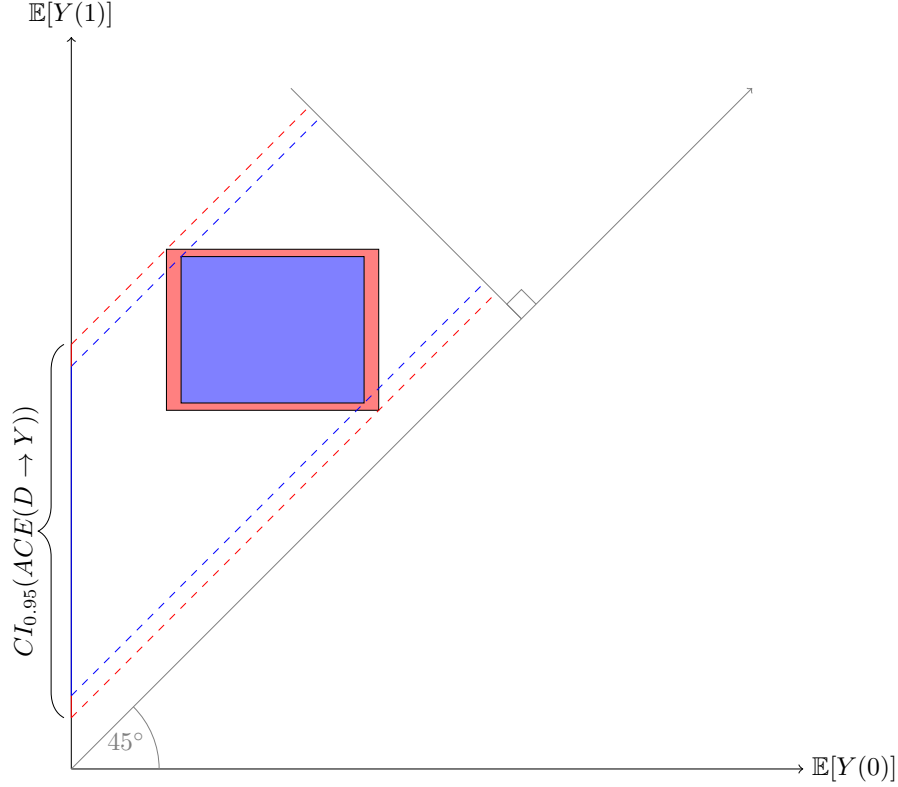


(b) As  $\theta_n^u(z)$  and  $\theta_n^u(z')$  are both binding constraints on  $\mathbb{E}[Y(0)]$ , the minimum of these upper bounds must be binding. The minimum upper bound is written as  $\theta_n^u$  and is  $\theta_n^u(z)$ . The inferential problem is to determine  $\phi_n^u$ , which satisfies  $\mathbb{P}(\mathbb{E}[Y(0)] > \phi_n^u) = \alpha$ . A naïve approach would be to look only at the one-sided  $1 - \alpha$  confidence region for  $\theta_n^u(z)$ , which is  $\phi_n^u(z)$ . This approach ignores variation in  $\theta_n^u(z')$ . An equally naïve approach would be to simply add aggregate variation to  $\theta_n^u(z)$ . In other words, to add the weighted average of  $\phi_n^u(z) - \theta_n^u(z)$  and  $\phi_n^u(z') - \theta_n^u(z')$  to  $\theta_n^u(z)$ . This approach is standard but fails in this case because it does not account for the fact that  $\theta_n^u(z) < \theta_n^u(z')$ . Inference must account for the fact that upward variation in  $\theta_n^u(z')$  does not matter so long as  $\theta_n^u(z) < \theta_n^u(z')$ ; equivalently, that  $\theta_n^u(z')$  is a one-sided  $1 - \gamma$  confidence region for  $\theta_n^u(z)$  for  $\gamma > \alpha$ . Chernozhukov et al. (2013) solves the inferential problem by adjusting the critical value that is associated with the one-sided  $1 - \alpha$ -confidence region. That is, Chernozhukov et al. (2013) adjusts  $k$  such that  $k$  that solves  $\mathbb{P}(\mathbb{E}[Y(0)] > \theta_n^u(z) + k\sigma) = \alpha$ . The solution for  $k$  yields  $\phi_n^u$  with the one-sided  $1 - \alpha$  confidence region for  $\theta_n^u$  given by the grey area. The distribution of  $\theta_n^u$  over repeated samples is non-standard in this case and the bootstrap is not necessarily consistent (Bugni, 2010).

Figure 6: A note on the inferential problem.



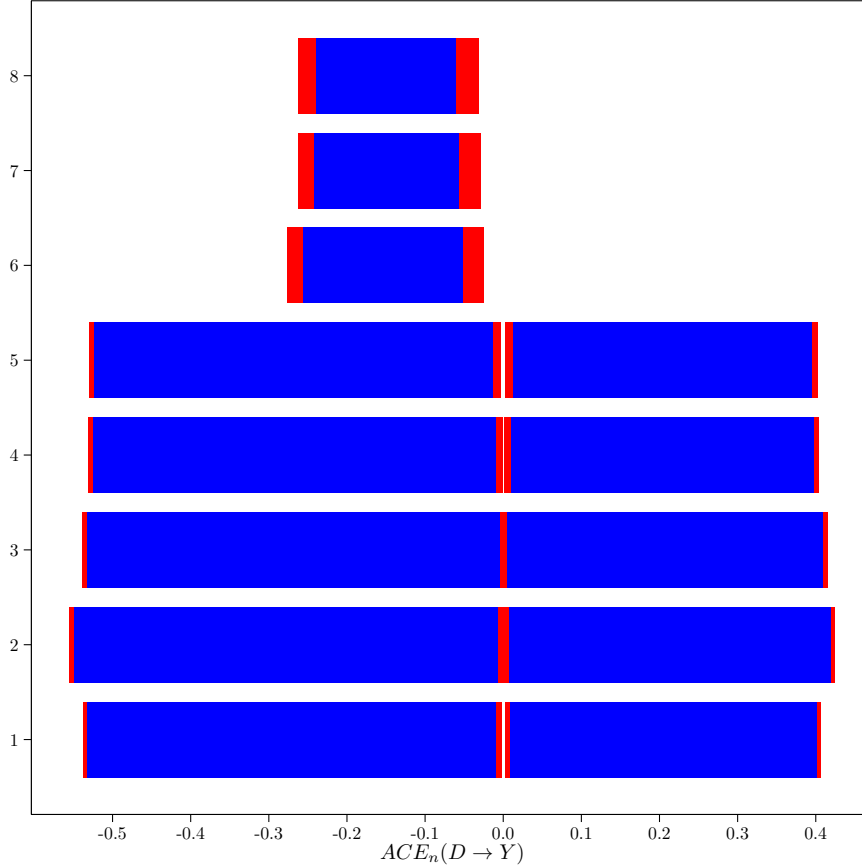
(a) Suppose that the admissible set of values of  $(\mathbb{E}[Y(0)], \mathbb{E}[Y(1)])$  is given by the blue rectangle, and that the  $1 - \alpha$ -confidence region for this set is the union of the blue rectangle and the red polygon.



(b) It is possible to recover  $ACE(D \rightarrow Y)$  from the plot. First, note that  $ACE(D \rightarrow Y)$  is increasing in the  $y$ -direction and is decreasing in the  $x$ -direction. Second, note that  $ACE(D \rightarrow Y)$  is constant along any line with unit gradient. Third, note that the value of  $ACE(D \rightarrow Y)$  along any line with unit gradient is dependent upon the value of the intercept of this line. Fourth, note that a projection from the normal of a line with unit gradient is a line that has unit gradient. For example, the blue dashed line is a projection from the normal of the  $45^\circ$  line; notice that the blue dashed line is parallel to the  $45^\circ$  line. Fifth, note that any projection from the normal of the  $45^\circ$  line that passes through the blue rectangle is an admissible value of  $(\mathbb{E}[Y(0)], \mathbb{E}[Y(1)])$ ; equivalently, that any projection from the normal of the  $45^\circ$  line that passes through the union of the blue rectangle and the red polygon is in the  $1 - \alpha$  confidence region of  $(\mathbb{E}[Y(0)], \mathbb{E}[Y(1)])$ . Together these five facts suggest that the admissible set of values of  $ACE(D \rightarrow Y)$  and its  $1 - \alpha$  confidence region can be recovered from the projection from the normal of the  $45^\circ$  line onto the  $y$ -axis. This method gives a geometric interpretation to  $ACE(D \rightarrow Y)$ .

Figure 7: A note on recovering the average causal effect.

(a) The plot shows how the admissible set of values of  $ACE_n(D \rightarrow Y)$  changes as  $\mathcal{R}_Z$  is varied. Each value on the  $y$ -axis corresponds to a different experiment. In each experiment,  $\mathcal{R}_Z$  is varied and the admissible set of values of  $ACE_n(D \rightarrow Y)$  is calculated. See Table 3 and Table 4 for the definition of  $\mathcal{R}_Z$  in each experiment (the column headed No. states the experiment number in each table, and the column headed  $\mathcal{R}_Z$  states the events that form the points of support for  $Z$ ). The plot is a graphical representation of Table 3 and Table 4. Blue regions represent the admissible set of values of  $ACE_n(D \rightarrow Y)$ , and the union of blue and red regions represent 0.950 confidence regions for  $ACE_n(D \rightarrow Y)$ .



(b) In experiment 1 through experiment 5, information relating to the incidence of a multiple second birth is ignored. The admissible sets of values of  $ACE_n(D \rightarrow Y)$  in these experiments are large and disconnected, which suggests that child gender is weakly associated with the number of children in a household. In other words, that child gender is a weak instrumental variable. Child gender is, by itself, insufficient to determine the sign of the average causal effect of the number of children in a household on a mother's employment. Experiment 2 and experiment 3 are unable to rule out a null effect at the 0.950 significance level (the central red regions overlap at zero). In experiment 6 through experiment 7, information relating to the incidence of a multiple second birth is incorporated. The sign of the average causal effect of the number of children in a household on a mother's employment is negative, and significant. The length of the admissible set of values of  $ACE_n(D \rightarrow Y)$  is decreasing as more information is incorporated, as is the 0.950 confidence region for  $ACE_n(D \rightarrow Y)$ .

Figure 8: Admissible sets of values of  $ACE_n(D \rightarrow Y)$  and  $\mathcal{R}_Z$ .

(a) The plot shows the admissible set of values of  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$  when  $\mathcal{R}_Z$  is formed of the events Male-Male  $\cup$  Female-Female and Male-Female  $\cup$  Female-Male. This is experiment 1 in Table 1. Blue regions represent the admissible set of values of  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$ , and the union of blue and red regions represent 0.950 confidence regions for  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$ . The confidence regions that are shown in the plot are different to those that are reported in Table 1. The distinction arises since Table 1 reports 0.950 confidence regions for  $\mathbb{E}_n[Y(0)]$  and for  $\mathbb{E}_n[Y(1)]$  whereas the plot shows the 0.950 confidence region for  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$ .

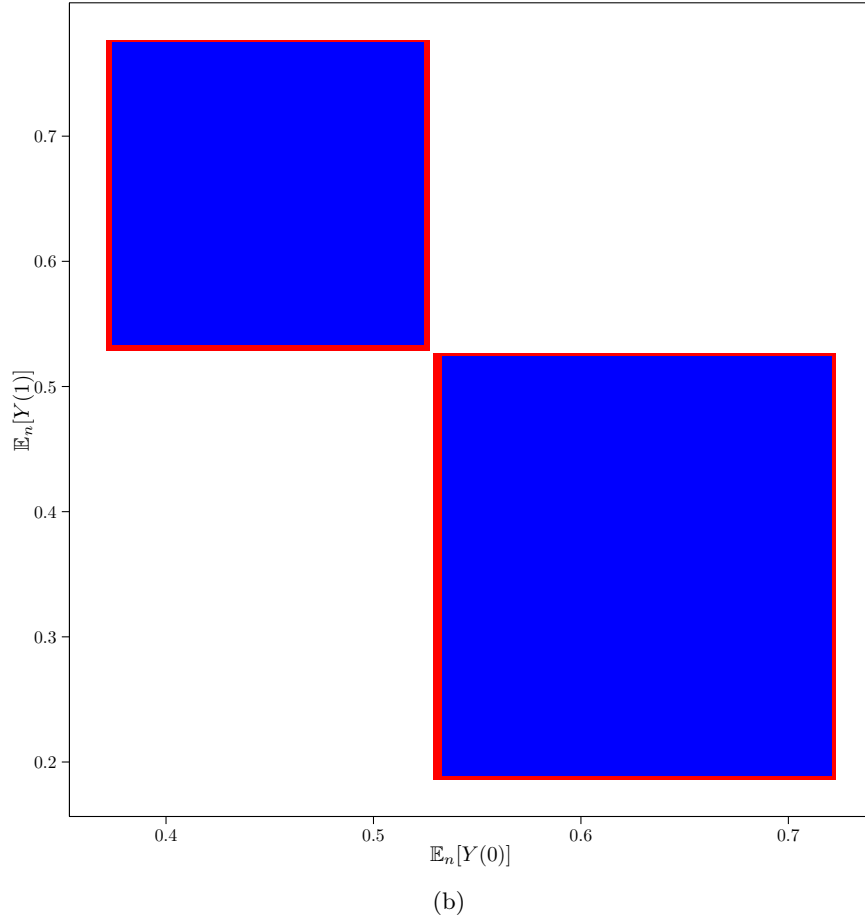


Figure 9

(a) The plot shows the admissible set of values of  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$  when  $\mathcal{R}_Z$  is formed of the events Multiple birth, Single birth  $\cap$  (Male-Male  $\cup$  Female-Female) and Single birth  $\cap$  (Male-Female  $\cup$  Male-Female). This is experiment 6 in Table 2. Blue regions represent the admissible set of values of  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$ , and the union of blue and red regions represent 0.950 confidence regions for  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$ . The confidence regions that are shown in the plot are different to those that are reported in Table 2. The distinction arises since Table 2 reports 0.950 confidence regions for  $\mathbb{E}_n[Y(0)]$  and for  $\mathbb{E}_n[Y(1)]$  whereas the plot shows the 0.950 confidence region for  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$ .

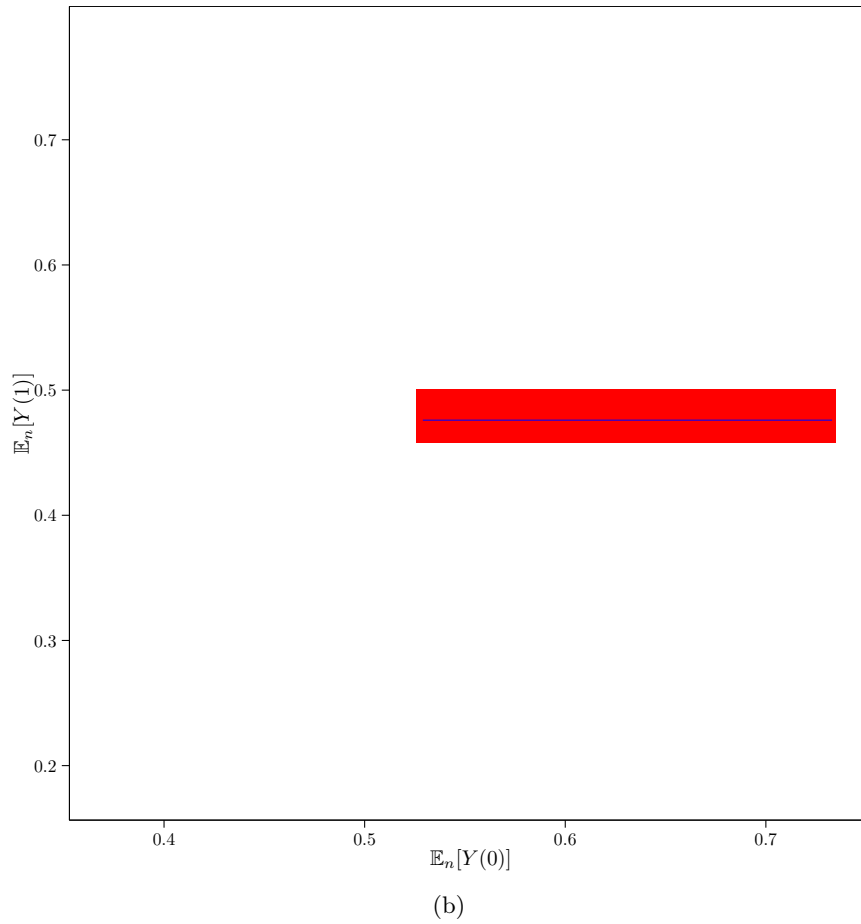


Figure 10

(a) The plot shows the admissible set of values of  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$  when  $\mathcal{R}_Z$  is formed of the events Multiple birth, Single birth  $\cap$  Male-Male, Single birth  $\cap$  Female-Female, Single birth  $\cap$  Male-Female and Single birth  $\cap$  Female-Male. This is experiment 8 in Table 2. Blue regions represent the admissible set of values of  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$ , and the union of blue and red regions represent 0.950 confidence regions for  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$ . The confidence regions that are shown in the plot are different to those that are reported in Table 2. The distinction arises since Table 2 reports 0.950 confidence regions for  $\mathbb{E}_n[Y(0)]$  and for  $\mathbb{E}_n[Y(1)]$  whereas the plot shows the 0.950 confidence region for  $(\mathbb{E}_n[Y(0)], \mathbb{E}_n[Y(1)])$ .

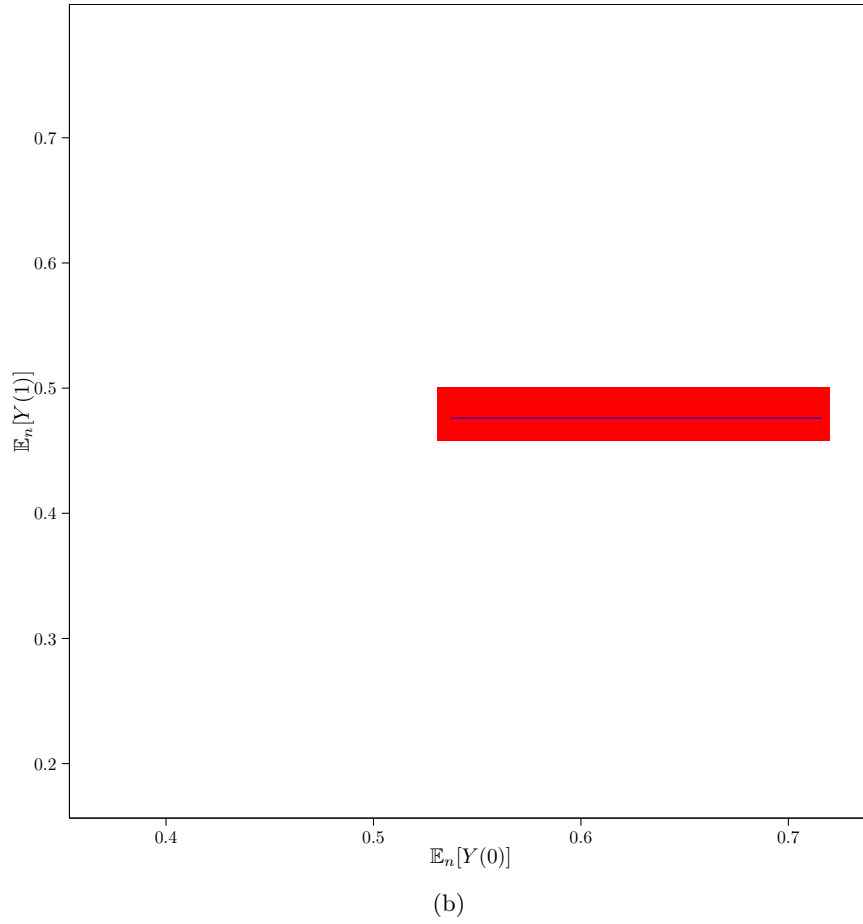
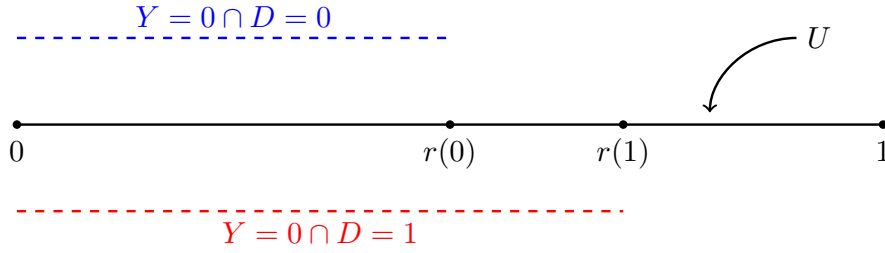


Figure 11

(a) The goal of the identification analysis is to determine the location of  $r(0)$  and of  $r(1)$ . The analysis is confounded by the fact that it is not known whether  $r(0) > r(1)$ ,  $r(0) < r(1)$  or  $r(0) = r(1)$ . The order of the collection of  $r(d)$  is crucial since it is the order that determines the value of  $\mathcal{A}_r$  and of  $\mathcal{B}_r$ .



(b) Suppose that  $r(0) < r(1)$ . What can be determined about the location of  $p(0)$  from observable variables alone? Any economic agent that is characterised by  $Y = 0 \cap D = 0$  must also be characterised by  $u$  somewhere in the blue dashed interval. Similarly, any economic agent that is characterised by  $Y = 0 \cap D = 1$  must also be characterised by  $u$  somewhere in the red dashed interval. The location of  $r(0)$  is confounded since it is not known what proportion of  $Y = 0 \cap D = 1$  are located in  $[0, r(0)]$  (the left partition of the red dashed interval), and what proportion of  $Y = 0 \cap D = 1$  are located in  $[r(0), r(1)]$  (the right partition of the red dashed interval). A lower bound on the location of  $r(0)$  is then that all economic agents that are characterised by  $Y = 0 \cap D = 1$  are located in  $[r(0), r(1)]$ . The measure of economic agents that are located in  $[0, r(0)]$  in this case is then  $\mathbb{P}(Y = 0, D = 0)$ . An upper bound on the location of  $r(0)$  is then that all economic agents that are characterised by  $Y = 0 \cap D = 1$  are located in  $[0, r(0)]$ . The measure of economic agents that are located in  $[0, r(0)]$  in this case is then  $\mathbb{P}(Y = 0, D = 0 \cup 1)$ . In this case,  $\mathcal{A}_r(0) = \{0\}$  and  $\mathcal{B}_r(0) = \{0, 1\}$ .

(c) Note that economic agents that are located in  $[0, r(0)]$  are characterised by  $Y(0) = 0$  and  $Y(1) = 0$ . In contrast, economic agents that are located in  $(r(0), r(1)]$  are characterised by  $Y(0) = 1$  and  $Y(1) = 0$ . The problem of determining the measure of economic agents that are characterised by each pair  $(Y(0), Y(1))$  is then equivalent to locating  $r(0)$  and  $r(1)$ . This point validates consideration of level sets of  $U$ .

Figure 12: A note on partial identification.

Table 1:  $\mathcal{R}_Z$  and admissible sets of values of  $\mathbb{E}_n[Y(d)]$ .

No.	$\mathcal{R}_Z$	Bound			
		$\mathbb{E}_n^-[Y(0)]$	$\mathbb{E}_n^-[Y(1)]$	$\mathbb{E}_n^+[Y(0)]$	$\mathbb{E}_n^+[Y(1)]$
1	Male-Male $\cup$ Female-Female Male-Female $\cup$ Female-Male	[0.533, 0.721] [0.530, 0.723]	[0.189, 0.524] [0.187, 0.527]	[0.374, 0.524] [0.371, 0.527]	[0.533, 0.775] [0.530, 0.777]
2	Male-Male Male-Female $\cup$ Female-Male $\cup$ Female-Female	[0.530, 0.731] [0.528, 0.733]	[0.182, 0.523] [0.179, 0.527]	[0.359, 0.523] [0.357, 0.527]	[0.530, 0.778] [0.528, 0.781]
3	Female-Female Male-Male $\cup$ Male-Female $\cup$ Female-Male	[0.529, 0.729] [0.527, 0.731]	[0.196, 0.524] [0.193, 0.529]	[0.362, 0.524] [0.360, 0.529]	[0.529, 0.771] [0.527, 0.774]
4	Male-Male Female-Female Male-Female $\cup$ Female-Male	[0.533, 0.721] [0.530, 0.723]	[0.196, 0.523] [0.192, 0.528]	[0.374, 0.523] [0.371, 0.528]	[0.533, 0.771] [0.530, 0.775]
5	Male-Male Female-Female Male-Female Female-Male	[0.536, 0.718] [0.531, 0.722]	[0.196, 0.523] [0.193, 0.528]	[0.376, 0.523] [0.372, 0.528]	[0.536, 0.771] [0.531, 0.775]

(a) The column headed No. states the experiment number. The column headed  $\mathcal{R}_Z$  states the events that form the points of support of  $Z$ . For example, in experiment 1  $Z$  is a random variable that takes the value  $z$  when the event Male-Male  $\cup$  Female-Female occurs and the value  $z'$  when this does not occur. The event *Male - Female* is the event that the oldest two children in a household are male and female, respectively. The columns headed Bound are the admissible set of values of  $\mathbb{E}_n[Y(D)]$ . A superscript - is written when the bounds are conditional on  $\mathbb{E}_n[Y(0)] \geq \mathbb{E}_n[Y(1)]$ , and a superscript + is written when the bounds are conditional on  $\mathbb{E}_n[Y(0)] \leq \mathbb{E}_n[Y(1)]$ . 0.950 confidence regions for  $\mathbb{E}_n[Y(D)]$  are in blue and are constructed from one-sided 0.975 confidence regions for the lower bound and the upper bound.

Table 2:  $\mathcal{R}_Z$  and admissible sets of values of  $\mathbb{E}_n[Y(d)]$ .

No.	$\mathcal{R}_Z$	Bound	
		$\mathbb{E}_n^-[Y(0)]$	$\mathbb{E}_n^-[Y(1)]$
6	Multiple birth	[0.529, 0.733]	0.476
	Single birth	[0.526, 0.734]	[0.460, 0.498]
7	Multiple birth	[0.534, 0.718]	0.476
	Single birth $\cap$ (Male-Male $\cup$ Female-Female)	[0.530, 0.721]	[0.460, 0.498]
	Single birth $\cap$ (Male-Female $\cup$ Female-Male)		
8	Multiple birth	[0.537, 0.716]	0.476
	Single birth $\cap$ Male-Male	[0.532, 0.720]	[0.460, 0.498]
	Single birth $\cap$ Female-Female		
	Single birth $\cap$ Male-Female		
	Single birth $\cap$ Female-Male		

(a) The column headed No. states the experiment number. The column headed  $\mathcal{R}_Z$  states the events that form the points of support of  $Z$ . For example, in experiment 1  $Z$  is a random variable that takes the value  $z$  when the event Male-Male  $\cup$  Female-Female occurs and the value  $z'$  when this does not occur. The event *Male - Female* is the event that the oldest two children in a household are male and female, respectively. The columns headed Bound are the admissible set of values of  $\mathbb{E}_n[Y(D)]$ . A superscript  $-$  is written when the bounds are conditional on  $\mathbb{E}_n[Y(0)] \geq \mathbb{E}_n[Y(1)]$ , and a superscript  $+$  is written when the bounds are conditional on  $\mathbb{E}_n[Y(0)] \leq \mathbb{E}_n[Y(1)]$ . 0.950 confidence regions for  $\mathbb{E}_n[Y(D)]$  are in blue and are constructed from one-sided 0.975 confidence regions for the lower bound and the upper bound.



Table 3:  $\mathcal{R}_Z$  and admissible sets of values of  $ACE_n(D \rightarrow Y)$ .

No.	$\mathcal{R}_Z$	Bound	
		$ACE_n^-(D \rightarrow Y)$	$ACE_n^+(D \rightarrow Y)$
1	Male-Male $\cup$ Female-Female	$[-0.532, -0.009]$	$[0.009, 0.401]$
	Male-Female $\cup$ Female-Male	$[-0.537, -0.002]$	$[0.002, 0.406]$
2	Male-Male	$[-0.549, -0.007]$	$[0.007, 0.419]$
	Male-Female $\cup$ Female-Male $\cup$ Female-Female	$[-0.555, 0.001]$	$[-0.001, 0.425]$
3	Female-Female	$[-0.533, -0.005]$	$[0.005, 0.409]$
	Male-Male $\cup$ Male-Female $\cup$ Female-Male	$[-0.539, 0.003]$	$[-0.003, 0.415]$
4	Male-Male	$[-0.525, -0.010]$	$[0.010, 0.397]$
	Female-Female	$[-0.531, -0.001]$	$[0.002, 0.404]$
5	Male-Male	$[-0.523, -0.013]$	$[0.013, 0.395]$
	Female-Female	$[-0.530, -0.003]$	$[0.003, 0.403]$
	Male-Female		
	Female-Male		

(a) The column headed No. states the experiment number. The column headed  $\mathcal{R}_Z$  states the events that form the points of support of  $Z$ . For example, in experiment 1  $Z$  is a random variable that takes the value  $z$  when the event Male-Male  $\cup$  Female-Female occurs and the value  $z'$  when this does not occur. The event *Male-Female* is the event that the oldest two children in a household are male and female, respectively. The columns headed Bound are the admissible set of values of  $ACE_n(D \rightarrow Y)$ . A superscript  $-$  is written when the bounds are conditional on  $ACE_n(D \rightarrow Y) \leq 0$ , and a superscript  $+$  is written when the bounds are conditional on  $ACE_n(D \rightarrow Y) \geq 0$ . 0.950 confidence regions for  $ACE_n(D \rightarrow Y)$  are in blue and are constructed from one-sided 0.975 confidence regions for the lower bound and the upper bound.

Table 4:  $\mathcal{R}_Z$  and admissible sets of values of  $ACE_n(D \rightarrow Y)$ .

No.	$\mathcal{R}_Z$	Bound
		$ACE_n^-(D \rightarrow Y)$
6	Multiple birth	$[-0.256, -0.052]$
	Single birth	$[-0.277, -0.025]$
7	Multiple birth	$[-0.242, -0.057]$
	Single birth $\cap$ (Male-Male $\cup$ Female-Female)	$[-0.263, -0.029]$
	Single birth $\cap$ (Male-Female $\cup$ Female-Male)	
8	Multiple birth	$[-0.240, -0.061]$
	Single birth $\cap$ Male-Male	
	Single birth $\cap$ Female-Female	$[-0.262, -0.031]$
	Single birth $\cap$ Male-Female	
	Single birth $\cap$ Female-Male	

(a) The column headed No. states the experiment number. The column headed  $\mathcal{R}_Z$  states the events that form the points of support of  $Z$ . For example, in experiment 1  $Z$  is a random variable that takes the value  $z$  when the event Male-Male  $\cup$  Female-Female occurs and the value  $z'$  when this does not occur. The event *Male-Female* is the event that the oldest two children in a household are male and female, respectively. The columns headed Bound are the admissible set of values of  $ACE_n(D \rightarrow Y)$ . A superscript  $-$  is written when the bounds are conditional on  $ACE_n(D \rightarrow Y) \leq 0$ , and a superscript  $+$  is written when the bounds are conditional on  $ACE_n(D \rightarrow Y) \geq 0$ . 0.950 confidence regions for  $ACE_n(D \rightarrow Y)$  are in blue and are constructed from one-sided 0.975 confidence regions for the lower bound and the upper bound.

## References

- Abrevaya, J., Y.-C. Hsu, and R. P. Lieli (2013). Estimating conditional average treatment effects. Technical report, Working paper.
- Angrist, J. D. (2014). Angrist data archive. [Online; accessed 25-July-2014].
- Angrist, J. D. and W. N. Evans (1998). Children and their parents' labor supply: Evidence from exogenous variation in family size. *American Economic Review* 88(3), 450 – 477.
- Artstein, Z. (1983). Distributions of random sets and random selections. *Israel Journal of Mathematics* 46(4), 313 – 324.
- Balke, A. and J. Pearl (1997). Bounds on treatment effects from studies with imperfect compliance. *Journal of the American Statistical Association* 92(439), 1171 – 1176.
- Bugni, F. A. (2010). Bootstrap inference in partially identified models defined by moment inequalities: Coverage of the identified set. *Econometrica* 78(2), 735 – 753.
- Chernozhukov, V., S. Lee, and A. M. Rosen (2013). Intersection bounds: Estimation and inference. *Econometrica* 81(2), 667 – 737.
- Chesher, A. (2005). Nonparametric identification under discrete variation. *Econometrica* 73(5), 1525 – 1550.
- Chesher, A. (2010). Instrumental variable models for discrete outcomes. *Econometrica* 78(2), 575 – 601.
- Chesher, A. and A. M. Rosen (2013). What do instrumental variable models deliver with discrete dependent variables?. *American Economic Review* 103(3), 557 – 562.
- Chesher, A., A. M. Rosen, and K. Smolinski (2013). An instrumental variable model of multiple discrete choice. *Quantitative Economics* 4(2), 157 – 196.
- Heckman, J. J. and E. Vytlacil (2005). Structural equations, treatment effects, and econometric policy evaluation1. *Econometrica* 73(3), 669 – 738.
- Hurwicz, L. (1950). Generalization of the concept of identification. *Statistical Inference in Dynamic Economic Models* (T. Koopmans, ed.). Cowles Commission, Monograph 10, 245 – 257.
- Imbens, G. W. and J. D. Angrist (1994). Identification and estimation of local average treatment effects. *Econometrica* 62(2), 467 – 475.
- Khan, S. and E. Tamer (2010). Irregular identification, support conditions, and inverse weight estimation. *Econometrica* 78(6), 2021 – 2042.
- Kitagawa, T. (2009). Identification region of the potential outcome distributions under instrument independence. *Econometrica* (Revise and resubmit).
- Koopmans, T. C. and O. Reiersøl (1950). The identification of structural characteristics. *Annals of Mathematical Statistics* 21(2), 165 – 181.
- Manski, C. F. (1988). *Analog Estimation Methods in Econometrics*. New York and London: Chapman and Hall.
- Manski, C. F. (2013). *Public Policy in an Uncertain World: Analysis and Decisions*. Cambridge and London: Harvard University Press.
- Molchanov, I. (2005). *Theory of Random Sets*. Springer.
- Pearl, J. (1995). On the testability of causal models with latent and instrumental variables. In *Proceedings of the Eleventh conference on Uncertainty in artificial intelligence*, pp. 435 – 443. Morgan Kaufmann Publishers Inc.

- Shaikh, A. M. and E. J. Vytlačil (2011). Partial identification in triangular systems of equations with binary dependent variables. *Econometrica* 79(3), 949 – 955.
- Strotz, R. H. and H. O. Wold (1960). Recursive vs. nonrecursive systems: An attempt at synthesis (part i of a triptych on causal chain systems). *Econometrica: Journal of the Econometric Society*, 417 – 427.