

# $\epsilon$ -Identifiability of Causal Quantities

Ang Li, Scott Mueller and Judea Pearl

Cognitive Systems Laboratory, Department of Computer Science,  
University of California, Los Angeles,  
Los Angeles, California, USA.  
{angli, scott, judea}@cs.ucla.edu

## Abstract

Identifying the effects of causes and causes of effects is vital in virtually every scientific field. Often, however, the needed probabilities may not be fully identifiable from the data sources available. This paper shows how partial identifiability is still possible for several probabilities of causation. We term this  $\epsilon$ -identifiability and demonstrate its usefulness in cases where the behavior of certain subpopulations can be restricted to within some narrow bounds. In particular, we show how unidentifiable causal effects and counterfactual probabilities can be narrowly bounded when such allowances are made. Often those allowances are easily measured and reasonably assumed. Finally,  $\epsilon$ -identifiability is applied to the unit selection problem.

## 1 Introduction

Both Effects of Causes (EoC) and Causes of Effects (CoE) play an important role in several fields, such as health science, social science, and business. For example, the causal effects identified by the adjustment [Pearl, 1993] formula helps decision-maker avoid randomized controlled trial using purely observational data. For another example, probabilities of causation have been proven critical in personalized decision-making [Mueller and Pearl, 2022]. Besides, a linear combination of probabilities of causation has been used to solve the unit selection problem defined by Li and Pearl [Li and Pearl, 2022b; Li and Pearl, 2019; Li and Pearl, 2022d]. Causal quantities can also increase the accuracy of machine learning models by combining causal quantities with the model’s label [Li *et al.*, 2020].

The causal quantities have been studied for decades. Pearl first defined the causal quantities such as causal effects [Pearl, 1993], probability of necessity and sufficiency (PNS), probability of sufficiency (PS), and probability of necessity (PN) [Pearl, 1999] and their identifiability [Pearl, 2009] using the structural causal model (SCM) [Galles and Pearl, 1998; Halpern, 2000]. Pearl also proposed the identification conditions of the causal effects (i.e., back-door and front-door criteria) [Pearl, 1993]. Pearl, Bareinboim, etc. have studied more conditions for identifying the causal effects [Bareinboim and Pearl, 2012;

Shpitser and Pearl, 2009]. If the causal effects are not identifiable, the informative bounds are given by Li and Pearl using non-linear programming [Li and Pearl, 2022c]. Then, Tian and Pearl proposed the identification conditions of the binary probabilities of causation (i.e., monotonicity) [Tian and Pearl, 2000]. If the probabilities of causation are not identifiable, Tian and Pearl [Tian and Pearl, 2000] also have informative tight bounds for them using Balke’s Linear programming [Balke and Pearl, 1997]. Mueller, Li, and Pearl [Mueller *et al.*, 2021], as well as Dawid [Dawid *et al.*, 2017], increased those bounds using additional covariate information and the corresponding causal structure. Recently, Li and Pearl also proposed the theoretical work for non-binary probabilities of causation [Li and Pearl, 2022a].

In real-world applications, decision-makers are more likely to have identifiable cases (i.e., the causal quantities have point estimations) because the bounds under unidentifiable cases may be less informative (e.g.,  $0.1 \leq \text{PNS} \leq 0.9$ ). Besides, estimating the bounds often requires a combination of experimental and observational data. So we wonder if something is sitting between the identifiable and the bounds. Inspired by the idea of the confidence interval, in this paper, we proposed the definition of  $\epsilon$ -identifiability, in which more conditions of  $\epsilon$ -identifiability can be found while the estimations of the causal quantities are still near point estimations.

## 2 Preliminaries

Here, we review the definition of PNS, PS, and PN defined by Pearl [Pearl, 1999], as well as the definition of identifiable and the conditions for identifying PNS, PS, and PN [Tian and Pearl, 2000]. Besides, we review the tight bounds of PNS, PS, and PN when they are unidentifiable [Tian and Pearl, 2000]. Readers who are familiar with the above knowledge may skip this section.

Similarly to any works mentioned above, we used the causal language of the SCM [Galles and Pearl, 1998; Halpern, 2000]. The introductory counterfactual sentence “Variable  $Y$  would have the value  $y$ , had  $X$  been  $x$ ” in this language is denoted by  $Y_x = y$ , and shorted as  $y_x$ . We have two types of data: experimental data, which is in the form of causal effects (denoted as  $P(y_x)$ ), and observational data, which is in the form of a joint probability function (denoted as  $P(x, y)$ ).

First, the definition of identifiable for any causal quantities defined using SCM is as follows:

**Definition 1 (Identifiability).** Let  $Q(M)$  be any computable quantity of a class of SCM  $M$  that is compatible with graph  $G$ . We say that  $Q$  is identifiable in  $M$  if, for any pairs of models  $M_1$  and  $M_2$  from  $M$ ,  $Q(M_1) = Q(M_2)$  whenever  $P_{M_1}(v) = P_{M_2}(v)$ , where  $P(v)$  is the statistical data over the set  $V$  of observed variables. If our observations are limited and permit only a partial set  $F_M$  of features (of  $P_M(v)$ ) to be estimated, we define  $Q$  to be identifiable from  $F_M$  if  $Q(M_1) = Q(M_2)$  whenever  $F_{M_1} = F_{M_2}$ . [Pearl, 2009]

Second, the definitions of three binary probabilities of causation defined using SCM are as follow [Pearl, 1999]:

**Definition 2 (Probability of necessity (PN)).** Let  $X$  and  $Y$  be two binary variables in a causal model  $M$ , let  $x$  and  $y$  stand for the propositions  $X = \text{true}$  and  $Y = \text{true}$ , respectively, and  $x'$  and  $y'$  for their complements. The probability of necessity is defined as the expression

$$\begin{aligned} PN &\triangleq P(Y_{x'} = \text{false} | X = \text{true}, Y = \text{true}) \\ &\triangleq P(y'_{x'} | x, y) \end{aligned}$$

**Definition 3 (Probability of sufficiency (PS)).** Let  $X$  and  $Y$  be two binary variables in a causal model  $M$ , let  $x$  and  $y$  stand for the propositions  $X = \text{true}$  and  $Y = \text{true}$ , respectively, and  $x'$  and  $y'$  for their complements. The probability of sufficiency is defined as the expression

$$PS \triangleq P(y_x | y', x')$$

**Definition 4 (Probability of necessity and sufficiency (PNS)).** Let  $X$  and  $Y$  be two binary variables in a causal model  $M$ , let  $x$  and  $y$  stand for the propositions  $X = \text{true}$  and  $Y = \text{true}$ , respectively, and  $x'$  and  $y'$  for their complements. The probability of necessity and sufficiency is defined as the expression

$$PNS \triangleq P(y_x, y'_{x'})$$

Third, we review the identification conditions for causal effects [Pearl, 1993; Pearl, 1995].

**Definition 5 (Back-door criterion).** Given an ordered pair of variables  $(X, Y)$  in a directed acyclic graph  $G$ , a set of variables  $Z$  satisfies the back-door criterion relative to  $(X, Y)$ , if no node in  $Z$  is a descendant of  $X$ , and  $Z$  blocks every path between  $X$  and  $Y$  that contains an arrow into  $X$ .

If a set of variables  $Z$  satisfies the back-door criterion for  $X$  and  $Y$ , the causal effects of  $X$  on  $Y$  are identifiable and given by the adjustment formula:

$$P(y_x) = \sum_z P(y|x, z)P(z). \quad (1)$$

**Definition 6 (Front-door criterion).** A set of variables  $Z$  is said to satisfy the front-door criterion relative to an ordered pair of variables  $(X, Y)$  if:

- $Z$  intercepts all directed paths from  $X$  to  $Y$ ;
- there is no back-door path from  $X$  to  $Z$ ; and
- all back-door paths from  $Z$  to  $Y$  are blocked by  $X$ .

If a set of variables  $Z$  satisfies the front-door criterion for  $X$  and  $Y$ , and  $P(x, Z) > 0$ , then the causal effects of  $X$  on  $Y$  are identifiable and given by the adjustment formula:

$$P(y_x) = \sum_z P(z|x) \sum_{x'} P(y|x', z)P(x').$$

If causal effects are not identifiable, Tian and Pearl [Tian and Pearl, 2000] provided the following bounds for the causal effects.

$$P(x, y) \leq P(y_x) \leq 1 - P(x, y'). \quad (2)$$

Finally, we review the identification conditions for PNS, PS, and PN [Tian and Pearl, 2000].

**Definition 7. (Monotonicity)** A Variable  $Y$  is said to be monotonic relative to variable  $X$  in a causal model  $M$  iff

$$y'_x \wedge y_{x'} = \text{false}.$$

**Theorem 8.** If  $Y$  is monotonic relative to  $X$ , then PNS, PN, and PS are all identifiable, and

$$PNS = P(y_x) - P(y_{x'}),$$

$$PN = \frac{P(y) - P(y_{x'})}{P(x, y)},$$

$$PS = \frac{P(y_x) - P(y)}{P(x', y')}.$$

If PNS, PN, and PS are not identifiable, informative bounds are given by Tian and Pearl [Tian and Pearl, 2000].

$$\max \left\{ \begin{array}{l} 0, \\ P(y_x) - P(y_{x'}), \\ P(y) - P(y_{x'}), \\ P(y_x) - P(y) \end{array} \right\} \leq \text{PNS} \quad (3)$$

$$\min \left\{ \begin{array}{l} P(y_x), \\ P(y'_{x'}), \\ P(x, y) + P(x', y'), \\ P(y_x) - P(y_{x'}) + \\ P(x, y') + P(x', y) \end{array} \right\} \geq \text{PNS} \quad (4)$$

$$\max \left\{ \begin{array}{l} 0, \\ \frac{P(y) - P(y_{x'})}{P(x, y)} \end{array} \right\} \leq \text{PN} \quad (5)$$

$$\min \left\{ \begin{array}{l} 1, \\ \frac{P(y'_{x'}) - P(x', y')}{P(x, y)} \end{array} \right\} \geq \text{PN} \quad (6)$$

$$\max \left\{ \begin{array}{l} 0, \\ \frac{P(y') - P(y'_x)}{P(x', y')} \end{array} \right\} \leq \text{PS} \quad (7)$$

$$\min \left\{ \begin{array}{l} 1, \\ \frac{P(y_x) - P(x, y)}{P(x', y')} \end{array} \right\} \geq \text{PS} \quad (8)$$

The identification conditions mentioned above (i.e., back-door and front-door criteria and monotonicity) are robust. However, it may still be hard to achieve in real-world applications. In this work, we extend the definition of identifiability, in which a sufficiently small interval is allowed. By the new definition, the estimates of causal quantities are still near point estimations, and more conditions for identifiability could be discovered. If nothing is specified, the discussion in this paper will be restricted to binary treatment and effect (i.e.,  $X$  and  $Y$  are binary).

### 3 Main Results

First, we extend the definition of identifiability, which we call  $\epsilon$ -identifiability.

**Definition 9** ( $\epsilon$ -Identifiability). *Let  $Q(M)$  be any computable quantity of a class of SCM  $M$  that is compatible with graph  $G$ . We say that  $Q$  is  $\epsilon$ -identifiable in  $M$  (and  $\epsilon$ -identified to  $q$ ) if, there exists  $q$  s.t. for any model  $m$  from  $M$ ,  $Q(m) \in [q - \epsilon, q + \epsilon]$  with statistical data  $P_M(v)$ , where  $P(v)$  is the statistical data over the set  $V$  of observed variables. If our observations are limited and permit only a partial set  $F_M$  of features (of  $P_M(v)$ ) to be estimated, we define  $Q$  to be  $\epsilon$ -identifiable from  $F_M$  if  $Q(m) \in [q - \epsilon, q + \epsilon]$  with statistical data  $F_M$ .*

With the above definition, the causal quantity is at a maximum distance of  $\epsilon$  from its true value. We will use the infix operator symbol  $\approx_\epsilon$  to represent its left-hand side being within  $\epsilon$  of its right-hand side:

$$r \approx_\epsilon q \iff r \in [q - \epsilon, q + \epsilon]. \quad (9)$$

The following sections explicate conditions for  $\epsilon$ -identifiability of causal effects, PNS, PS, and PN.

#### 3.1 $\epsilon$ -Identifiability of Causal Effects

The causal effect  $P(Y_X)$  can be  $\epsilon$ -identified with information about the observational joint distribution  $P(X, Y)$ . This can be seen by rewriting Equation (2) as:

$$P(x, y) \leq P(y_x) \leq P(x, y) + P(x'). \quad (10)$$

Here,  $P(y_x)$  is  $\epsilon$ -identified to  $P(x, y) + \epsilon$  when  $P(x') \leq 2\epsilon$ . This  $\epsilon$ -identification indicates a lower bound of  $P(x, y)$  and an upper bound of  $P(x, y) + 2\epsilon$ . Since  $P(x') \leq 2\epsilon$ , these bounds are equivalent to (10). Notably, only  $P(x, y)$  and an upper bound on  $P(x')$  are necessary to  $\epsilon$ -identify  $P(y_x)$ . This is generalized in Theorem 10, without any assumptions of the causal structure.

**Theorem 10.** *The causal effect  $P(Y_X)$  is  $\epsilon$ -identified as follows:*

$$P(y_x) \approx_\epsilon P(x, y) + \epsilon \quad \text{if } P(x') \leq 2\epsilon, \quad (11)$$

$$P(y'_x) \approx_\epsilon P(x, y') + \epsilon \quad \text{if } P(x') \leq 2\epsilon, \quad (12)$$

$$P(y_{x'}) \approx_\epsilon P(x', y) + \epsilon \quad \text{if } P(x) \leq 2\epsilon, \quad (13)$$

$$P(y'_{x'}) \approx_\epsilon P(x', y') + \epsilon \quad \text{if } P(x) \leq 2\epsilon. \quad (14)$$

*Proof.* See Appendix 8.1.  $\square$

When the complete distribution  $P(X, Y)$  is known, Theorem 10 provides no extra precision over Equation (10). Its power comes from when only part of the distribution is known and only an upper bound on  $P(X)$  is available or able to be assumed.

Knowledge of a causal structure can aid  $\epsilon$ -identification. In Figure 1, there is a binary confounder  $U$ . If the full joint distribution  $P(X, Y, U)$  was available, the causal effect  $P(Y_X)$  could be computed simply through the backdoor adjustment formula. In the absence of the full joint distribution, Theorem 11 allows  $\epsilon$ -identification of  $P(y_x)$  with only knowledge of  $P(x)$  and the conditional probability  $P(y|x)$  as well as an upper bound on  $P(u)$ .

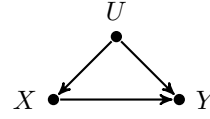


Figure 1: The causal graph, where  $X$  is a binary treatment,  $Y$  is a binary effect, and  $U$  is a binary confounder.

**Theorem 11.** *Given the causal graph in Figure 1 and  $P(u) \leq P(x) - c$  for some constant  $c$ , where  $0 < c \leq P(x)$ ,*

$$P(y_x) \approx_\epsilon P(y|x) + \frac{P(x) - c}{2cP(x) + P(x) + c} \cdot \epsilon$$

$$\text{if } P(u) \leq \frac{2cP(x)}{2cP(x) + P(x) + c} \cdot \epsilon. \quad (15)$$

*Specifically, if  $P(x) \geq 0.5$ , then the causal effect  $P(y_x)$  is  $\epsilon$ -identified to  $P(y|x) + \frac{\epsilon}{13}$  if  $P(u) < \frac{4}{13}\epsilon$ .*

*Proof.* See Appendix 8.2.  $\square$

Note that  $x \in \{x, x'\}$ ,  $y \in \{y, y'\}$ , and  $u \in \{u, u'\}$  in Theorem 11. The constant  $c$  should be maximized satisfying both  $c \leq P(x) - P(u)$  and the condition in Equation (15) for a given  $\epsilon$ . The larger  $c$  is, the closer  $P(y_x)$  is  $\epsilon$ -identified to  $P(y|x)$ . This needs to be balanced with minimizing  $\epsilon$ .

As an example, if  $P(x) \geq 0.5$  and  $P(u) \leq 0.1$ , then the causal effect  $P(y_x)$  is  $\epsilon$ -identified to  $P(y|x) + \frac{\epsilon}{13}$  if  $P(u) \leq \frac{4}{13}\epsilon$ .

Essentially,  $P(y_x)$  is  $\epsilon$ -identified to  $P(y|x)$  plus some fraction of  $\epsilon$  when  $P(u)$  is sufficiently small. Therefore, the causal effect  $P(y_x)$  is near  $P(y|x)$  if  $P(U)$  is specific (i.e.,  $P(u)$  or  $P(u')$  is minimal). In this case, Theorem 11 can be advantageous over the backdoor adjustment formula to compute  $P(y_x)$ , even when data on  $X$ ,  $Y$ , and  $U$  are available, because  $P(Y|X, U)$ , required for the adjustment formula, is impractical to estimate with  $P(U)$  close to 0.

#### 3.2 $\epsilon$ -Identifiability of PNS

Even though Tian and Pearl derived tight bounds on PNS [Tian and Pearl, 2000], the PNS can be potentially further narrowed when taking into account particular upper bound assumptions on causal effects or observational probabilities. This can be seen by analyzing the bounds of PNS in Equations (3) and (4). Picking any of the arguments to the max function of the lower bound and any of the arguments to the min function of the upper bound, we can make a condition that the range of those two values is less than  $2\epsilon$ . For example, let us pick the second argument of the max function,  $P(y_x) - P(y_{x'})$ , and the first argument of the min function,  $P(y_x)$ :

$$P(y_x) - [P(y_x) - P(y_{x'})] \leq 2\epsilon,$$

$$P(y_{x'}) \leq 2\epsilon. \quad (16)$$

Equation (16) is the assumption and the PNS is the  $\epsilon$ -identified to  $\epsilon$  above the lower bound or  $\epsilon$  below the upper bound:

$$\text{PNS} \approx_\epsilon P(y_x) - P(y_{x'}) + \epsilon, \text{ or} \quad (17)$$

$$\text{PNS} \approx_\epsilon P(y_x) - \epsilon. \quad (18)$$

Since it is assumed that  $P(y_{x'}) \leq 2\epsilon$ , Equation (17) is equivalent to Equation (18). The complete set of  $\epsilon$ -identifications and associated conditions are stated in Theorem 12.

**Theorem 12.** *The PNS is  $\epsilon$ -identified as follows:*

$$PNS \approx_{\epsilon} \epsilon \quad \text{if } P(y_x) \leq 2\epsilon, \quad (19)$$

$$PNS \approx_{\epsilon} \epsilon \quad \text{if } P(y'_{x'}) \leq 2\epsilon, \quad (20)$$

$$PNS \approx_{\epsilon} \epsilon \quad \text{if } P(x, y) + P(x', y') \leq 2\epsilon, \quad (21)$$

$$PNS \approx_{\epsilon} \epsilon \quad \text{if } P(y_x) - P(y_{x'}) + P(x, y') + P(x', y) \leq 2\epsilon, \quad (22)$$

$$PNS \approx_{\epsilon} P(y_x) - \epsilon \quad \text{if } P(y_{x'}) \leq 2\epsilon, \quad (23)$$

$$PNS \approx_{\epsilon} P(y'_{x'}) - \epsilon \quad \text{if } P(y'_x) \leq 2\epsilon, \quad (24)$$

$$PNS \approx_{\epsilon} P(y_x) - P(y_{x'}) + \epsilon \quad \text{if } P(x, y') + P(x', y) \leq 2\epsilon, \quad (25)$$

$$PNS \approx_{\epsilon} P(y_x) - P(y_{x'}) + \epsilon \quad \text{if } P(y_{x'}) - P(y_x) + P(x, y) + P(x', y') \leq 2\epsilon, \quad (26)$$

$$PNS \approx_{\epsilon} P(x, y) - P(x', y') - \epsilon \quad \text{if } P(y_{x'}) - P(y_x) + P(x, y) + P(x', y') \leq 2\epsilon, \quad (27)$$

$$PNS \approx_{\epsilon} P(y'_{x'}) - \epsilon \quad \text{if } P(y') \leq 2\epsilon, \quad (28)$$

$$PNS \approx_{\epsilon} P(y_x) - \epsilon \quad \text{if } P(y_x) + P(y_{x'}) - P(y) \leq 2\epsilon, \quad (29)$$

$$PNS \approx_{\epsilon} P(y) - P(y_{x'}) + \epsilon \quad \text{if } P(y_x) + P(y_{x'}) - P(y) \leq 2\epsilon, \quad (30)$$

$$PNS \approx_{\epsilon} P(x, y) + P(x', y') - \epsilon \quad \text{if } P(x', y') + P(y_{x'}) - P(x', y) \leq 2\epsilon, \quad (31)$$

$$PNS \approx_{\epsilon} P(y) - P(y_{x'}) + \epsilon \quad \text{if } P(x', y') + P(y_{x'}) - P(x', y) \leq 2\epsilon, \quad (32)$$

$$PNS \approx_{\epsilon} P(y) - P(y_{x'}) + \epsilon \quad \text{if } P(x', y) + P(y'_{x'}) - P(x', y') \leq 2\epsilon, \quad (33)$$

$$PNS \approx_{\epsilon} P(y_x) - \epsilon \quad \text{if } P(y) \leq 2\epsilon, \quad (34)$$

$$PNS \approx_{\epsilon} P(y'_{x'}) - \epsilon \quad \text{if } P(y'_{x'}) - P(y_x) + P(y) \leq 2\epsilon, \quad (35)$$

$$PNS \approx_{\epsilon} P(y) - P(y_{x'}) + \epsilon \quad \text{if } P(y'_{x'}) - P(y_x) + P(y) \leq 2\epsilon, \quad (36)$$

$$PNS \approx_{\epsilon} P(x, y) + P(x', y') - \epsilon \quad \text{if } P(x, y) + P(y'_{x'}) - P(x, y') \leq 2\epsilon, \quad (37)$$

$$PNS \approx_{\epsilon} P(y_x) - P(y) + \epsilon \quad \text{if } P(x, y) + P(y'_{x'}) - P(x, y') \leq 2\epsilon, \quad (38)$$

$$PNS \approx_{\epsilon} P(y_x) - P(y) + \epsilon \quad \text{if } P(x', y) + P(y'_{x'}) - P(x', y') \leq 2\epsilon. \quad (39)$$

*Proof.* See Appendix 8.3.  $\square$

Note that in the above theorem, eight conditions consist solely of experimental probabilities or solely of observational probabilities. This potentially eliminates the need for some types of studies, at least partially, even when estimating a counterfactual quantity such as PNS. For example, if a decision-maker knows that  $P(y)$  is large ( $P(y) \geq 0.95$ ), they can immediately conclude  $PNS \approx_{0.05} P(y'_{x'}) - 0.05$  without knowing the specific value of  $P(y)$ . Thus, only a control group study would be sufficient.

### 3.3 $\epsilon$ -Identifiability of PN and PS

Tian and Pearl derived tight bounds on PN and PS in addition to PNS. Similar to the derivation of Theorem 12, we can potentially narrow those bounds by taking into account upper bound assumptions on causal effects or observational probabilities. The set of  $\epsilon$ -identifications and associated conditions are stated in Theorems 13 and 14.

**Theorem 13.** *The PN is  $\epsilon$ -identified as follows:*

$$PN \approx_{\epsilon} \epsilon \quad \text{if } P(y'_{x'}) - P(x', y') \leq 2\epsilon P(x, y), \quad (40)$$

$$PN \approx_{\epsilon} 1 - \epsilon \quad \text{if } P(y_{x'}) - P(x', y) \leq 2\epsilon P(x, y), \quad (41)$$

$$PN \approx_{\epsilon} \frac{P(y) - P(y_{x'})}{P(x, y)} + \epsilon \quad \text{if } P(y_{x'}) - P(x', y) \leq 2\epsilon P(x, y), \quad (42)$$

$$PN \approx_{\epsilon} \frac{P(y'_{x'}) - P(x', y')}{P(x, y)} - \epsilon \quad \text{if } P(x, y') \leq 2\epsilon P(x, y), \quad (43)$$

$$PN \approx_{\epsilon} \frac{P(y) - P(y_{x'})}{P(x, y)} + \epsilon \quad \text{if } P(x, y') \leq 2\epsilon P(x, y). \quad (44)$$

*Proof.* See Appendix 8.4.  $\square$

Table 1: Results of an observational study with 1500 individuals who have access to the medicine, where 1260 individuals chose to receive the medicine and 240 individuals chose not to.

	Take the medicine	Take no medicine
Recovered	780	210
Not recovered	480	30

**Theorem 14.** *The PS is  $\epsilon$ -identified as follows:*

$$PS \approx_{\epsilon} \epsilon \quad \text{if } P(y_x) - P(x, y) \leq 2\epsilon P(x', y'), \quad (45)$$

$$PS \approx_{\epsilon} 1 - \epsilon \quad \text{if } P(y'_x) - P(x, y') \leq 2\epsilon P(x', y'), \quad (46)$$

$$PS \approx_{\epsilon} \frac{P(y') - P(y'_x)}{P(x', y')} + \epsilon \quad \text{if } P(y'_x) - P(x, y') \leq 2\epsilon P(x', y'), \quad (47)$$

$$PS \approx_{\epsilon} \frac{P(y_x) - P(x, y)}{P(x', y')} - \epsilon \quad \text{if } P(x', y) \leq 2\epsilon P(x', y'), \quad (48)$$

$$PS \approx_{\epsilon} \frac{P(y') - P(y'_x)}{P(x', y')} + \epsilon \quad \text{if } P(x', y) \leq 2\epsilon P(x', y'). \quad (49)$$

*Proof.* See Appendix 8.5.  $\square$

## 4 Examples

Here, we illustrate how to apply  $\epsilon$ -Identifiability in real applications by two simulated examples.

### 4.1 Causal Effects of Medicine

Consider a medicine manufacturer who wants to know the causal effect of a new medicine on a disease. They conducted an observational study where 1500 patients were given access to the medicine; the results of the study are summarized in Table 1. In addition, the expert from the medicine manufacturer acknowledged that family history is the only confounder of taking medicine and recovery, and the family history of the disease is extremely rare; only 1% of the people have the family history.

Let  $X = x$  denote that a patient chose to take the medicine, and  $X = x'$  denote that a patient chose not to take the medicine. Let  $Y = y$  denote that a patient recovered, and  $Y = y'$  denote that a patient did not recover. Let  $U = u$  denote that a patient has the family history, and  $U = u'$  denote that a patient has no family history.

To obtain the causal effect of the medicine (i.e., using adjustment formula (1)), we have to know the observational data associated with family history, which is difficult to obtain.

Fortunately, from Table 1, we obtain that  $P(x) = 0.84$  and  $P(y|x) = 0.62$ . We also have the prior that  $P(u) = 0.01$ . Since  $0.01 = P(u) \leq P(x) - 0.8$  (let  $c = 0.8$ ) and  $0.01 = P(u) < \frac{2c * 0.025P(x)}{2cP(x) + P(x) + c} = 0.0113$ , we can apply Theorem 11 to obtain that  $P(y_x)$  is 0.025-identified to  $P(y|x) + \frac{P(x) - c}{2cP(x) + P(x) + c} * 0.025 = 0.62$ . This means the causal

effect of the medicine is very close to 0.62 (i.e., 0.025 close), which can not be 0.025 far from 0.62. Then the medicine manufacturer can conclude that the causal effect of the medicine is roughly 0.62 without knowing the observational data associated with the family history.

Or even simpler, note that  $P(x) = 0.84 > 0.5$  and  $P(u) = 0.01 < 0.1$ ,  $P(u) = 0.01 < \frac{4}{13} * 0.035 = 0.0108$ . We obtain that  $P(y_x)$  is 0.035-identified to  $P(y|x) + \frac{0.035}{13} = 0.62$ . The decision-maker can make the same conclusion as above.

### 4.2 PNS of Flu Shot

Consider a newly invented flu shot. After a vaccination company introduced a new flu shot, the number of people infected by flu reached the lowest point in 20 years (i.e., less than 5% of people infected by flu). The government concluded that the new flu shot is the key to success. However, some anti-vaccination associations believe it is because people's physical quality increases yearly. Therefore, they all want to know how many percentages of people are uninfected because of the flu shot. The PNS of the flu shot (i.e., the percentage of individuals who would not infect by the flu if they had taken the flu shot and would infect otherwise) is indeed what they want.

Let  $X = x$  denote that an individual has taken the flu shot and  $X = x'$  denote that an individual has not taken the flu shot. Let  $Y = y$  denote an individual infected by the flu and  $Y = y'$  denote an individual not infected by the flu.

If they want to apply the bounds of PNS in Equations (3) and (4), they must conduct both experimental and observational studies. However, note that  $P(y) < 0.05$ , one could apply Equation (34) in Theorem 12, which PNS is 0.025-identified to  $P(y_x) - 0.025$  (i.e., PNS is very close to  $P(y_x)$ ). Thus, according to [Li *et al.*, 2022], only an experimental study for the treated group with a sample size of 385 is adequate for estimating PNS.

## 5 $\epsilon$ -Identifiability in Unit Selection Problem

One utility of the causal quantities is the unit selection problem [Li and Pearl, 2022b; Li and Pearl, 2019], in which Li and Pearl defined an objective causal function to select a set of individuals that have the desired mode of behavior.

Let  $X$  denote the binary treatment and  $Y$  denote the binary effect. According to Li and Pearl, individuals were divided into four response types: Complier (i.e.,  $P(y_x, y'_{x'})$ ), always-taker (i.e.,  $P(y_x, y_{x'})$ ), never-taker (i.e.,  $P(y'_{x'}, y'_{x'})$ ), and defier (i.e.,  $P(y'_{x'}, y_{x'})$ ). Suppose the payoff of selecting a complier, always-taker, never-taker, and defier is  $\beta, \gamma, \theta, \delta$ , respectively (i.e., benefit vector). The objective function (i.e., benefit function) that optimizes the composition of the four types over the selected set of individuals  $c$  is as follows:

$$f(c) = \beta P(y_x, y'_{x'}|c) + \gamma P(y_x, y_{x'}|c) + \theta P(y'_{x'}, y'_{x'}|c) + \delta P(y'_{x'}, y_{x'}|c).$$

Li and Pearl provided two types of identifiability conditions for the benefit function. One is about the response type such that there is no defier in the population (i.e., monotonicity). Another is about the benefits vector's relations, such that  $\beta + \delta = \gamma + \theta$  (i.e., gain equality). These two conditions are helpful

Table 2: Results of an experimental study with 1500 randomly selected customers were forced to apply the discount, and 1500 randomly selected customers were forced not to.

	Discount	No discount
Bought the purchase	900	750
No purchase	600	750

but still too specific and challenging to satisfy in real-world applications. If the benefit function is not identifiable, it can be bounded using experimental and observational data. Here in this paper, we extend the gain equality to the  $\epsilon$ -identifiability as stated in the following theorem.

**Theorem 15.** *Given a causal diagram  $G$  and distribution compatible with  $G$ , let  $C$  be a set of variables that does not contain any descendant of  $X$  in  $G$ , then the benefit function  $f(c) = \beta P(y_x, y_{x'}|c) + \gamma P(y_x, y_{x'}|c) + \theta P(y_{x'}, y_{x'}|c) + \delta P(y_{x'}, y_{x'}|c)$  is  $\frac{|\beta - \gamma - \theta + \delta|}{2}$ -identified to  $(\gamma - \delta)P(y_x|c) + \delta P(y_{x'}|c) + \theta P(y_{x'}|c) + \frac{\beta - \gamma - \theta + \delta}{2}$ .*

One critical use case of the above theorem is that decision-makers usually only care about the sign (gain or lose) of the benefit function. Decision-makers can apply the above theorem before conducting any observational study to see if the sign of the benefit function can be determined, as we will illustrate in the next section.

### 5.1 Example: Non-immediate Profit

Consider the most common example in [Li and Pearl, 2019]. A sale company proposed a discount on a purchase in order to increase the total non-immediate profit. The company assessed that the profit of offering the discount to complier, always-taker, never-taker, and defier is \$100, -\$60, \$0, -\$140, respectively. Let  $X = x$  denote that a customer applied the discount, and  $X = x$  denote that a customer did not apply the discount. Let  $Y = y$  denote that a customer bought the purchase and  $Y = y'$  denote that a customer did not. The benefit function is then (here  $c$  denote all customers)

$$f(c) = 100P(y_x, y_{x'}|c) - 60P(y_x, y_{x'}|c) + 0P(y_{x'}, y_{x'}|c) - 140P(y_{x'}, y_{x'}|c).$$

The company conducted an experimental study where 1500 randomly selected customers were forced to apply the discount, and 1500 randomly selected customers were forced not to. The results are summarized in Table 2. The experimental data reads  $P(y_x|c) = 0.6$  and  $P(y_{x'}|c) = 0.5$ .

Before conducting any observational study, one can conclude that the benefit function is 10-identified to -12 using Theorem 15. This result indicates that the benefit function is at most 10 away from -12; thus, the benefit function is negative regardless of the observational data. The decision-maker then can easily conclude that the discount should not offer to the customers.

## 6 Discussion

We have defined the  $\epsilon$ -identifiability of causal quantities and provided a list of  $\epsilon$ -identifiable conditions for causal effects,

PNS, PN, and PS. We still have some further discussions about the topic.

First, all conditions except Theorem 11 are conditions from observational or experimental data. In other words, if some of the observational or experimental distributions satisfied a particular condition, then the causal quantities are  $\epsilon$ -identifiable. These conditions are advantageous in real-world applications as no specific causal graph is needed. However, we still love to discover more graphical conditions of  $\epsilon$ -identifiability, such as back-door or front-door criterion.

Second, the bounds of PNS, PS, PN, and the benefit function can be narrowed by covariates information with their causal structure [Dawid *et al.*, 2017; Li and Pearl, 2022d; Mueller *et al.*, 2021]. The  $\epsilon$ -identifiability can also be extended if covariates information and their causal structure are available, which should be an exciting direction in the future.

Third, monotonicity is defined using a causal quantity, and in the meantime, monotonicity is also an identifiable condition for other causal quantities (e.g., PNS). Thus, another charming direction is how the  $\epsilon$ -identifiability of monotonicity affects the  $\epsilon$ -identifiability of other causal quantities.

## 7 Conclusion

In this paper, we defined the  $\epsilon$ -identifiability of causal quantities, which is easier to satisfy in real-world applications. We provided the  $\epsilon$ -identifiability conditions for causal effects, PNS, PS, and PN. We further illustrated the use cases of the proposed conditions by simulated examples.

## Acknowledgements

This research was supported in parts by grants from the National Science Foundation [#IIS-2106908 and #IIS-2231798], Office of Naval Research [#N00014-21-1-2351], and Toyota Research Institute of North America [#PO-000897].

## References

- [Balke and Pearl, 1997] Alexander A Balke and Judea Pearl. Probabilistic counterfactuals: Semantics, computation, and applications. Technical report, UCLA Dept. of Computer Science, 1997.
- [Bareinboim and Pearl, 2012] E. Bareinboim and J. Pearl. Causal inference by surrogate experiments:  $z$ -identifiability. In Nando de Freitas and Kevin Murphy, editors, *Proceedings of the Twenty-Eighth Conference on Uncertainty in Artificial Intelligence*, pages 113–120, Corvallis, OR, 2012. AUAI Press.
- [Dawid *et al.*, 2017] Philip Dawid, Monica Musio, and Rossella Murtas. The probability of causation. *Law, Probability and Risk*, (16):163–179, 2017.
- [Galles and Pearl, 1998] David Galles and Judea Pearl. An axiomatic characterization of causal counterfactuals. *Foundations of Science*, 3(1):151–182, 1998.
- [Halpern, 2000] Joseph Y Halpern. Axiomatizing causal reasoning. *Journal of Artificial Intelligence Research*, 12:317–337, 2000.

- [Li and Pearl, 2019] Ang Li and Judea Pearl. Unit selection based on counterfactual logic. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pages 1793–1799. International Joint Conferences on Artificial Intelligence Organization, 7 2019.
- [Li and Pearl, 2022a] A. Li and J. Pearl. Probabilities of causation with non-binary treatment and effect. Technical Report R-516, <[http://ftp.cs.ucla.edu/pub/stat\\_ser/r516.pdf](http://ftp.cs.ucla.edu/pub/stat_ser/r516.pdf)>, Department of Computer Science, University of California, Los Angeles, CA, 2022.
- [Li and Pearl, 2022b] A. Li and J. Pearl. Unit selection with nonbinary treatment and effect. Technical Report R-517, <[http://ftp.cs.ucla.edu/pub/stat\\_ser/r517.pdf](http://ftp.cs.ucla.edu/pub/stat_ser/r517.pdf)>, Department of Computer Science, University of California, Los Angeles, CA, 2022.
- [Li and Pearl, 2022c] Ang Li and Judea Pearl. Bounds on causal effects and application to high dimensional data. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 5773–5780, 2022.
- [Li and Pearl, 2022d] Ang Li and Judea Pearl. Unit selection with causal diagram. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 5765–5772, 2022.
- [Li *et al.*, 2020] Ang Li, Suming J. Chen, Jingzheng Qin, and Zhen Qin. Training machine learning models with causal logic. In *Companion Proceedings of the Web Conference 2020*, pages 557–561, 2020.
- [Li *et al.*, 2022] A. Li, R. Mao, and J. Pearl. Probabilities of causation: Adequate size of experimental and observational samples. Technical Report R-518, <[http://ftp.cs.ucla.edu/pub/stat\\_ser/r518.pdf](http://ftp.cs.ucla.edu/pub/stat_ser/r518.pdf)>, Department of Computer Science, University of California, Los Angeles, CA, 2022.
- [Mueller and Pearl, 2022] Mueller and Pearl. Personalized decision making – a conceptual introduction. Technical Report R-513, Department of Computer Science, University of California, Los Angeles, CA, 2022.
- [Mueller *et al.*, 2021] S. Mueller, A. Li, and J. Pearl. Causes of effects: Learning individual responses from population data. Technical Report R-505, <[http://ftp.cs.ucla.edu/pub/stat\\_ser/r505.pdf](http://ftp.cs.ucla.edu/pub/stat_ser/r505.pdf)>, Department of Computer Science, University of California, Los Angeles, CA, 2021. Forthcoming, Proceedings of IJCAI-2022.
- [Pearl, 1993] J Pearl. Aspects of graphical models connected with causality. *Proceedings of the 49th Session of the international Statistical Institute, Italy*, pages 399–401, 1993.
- [Pearl, 1995] Judea Pearl. Causal diagrams for empirical research. *Biometrika*, 82(4):669–688, 1995.
- [Pearl, 1999] Judea Pearl. Probabilities of causation: Three counterfactual interpretations and their identification. *Synthese*, pages 93–149, 1999.
- [Pearl, 2009] Judea Pearl. *Causality*. Cambridge university press, 2nd edition, 2009.
- [Shpitser and Pearl, 2009] I. Shpitser and J Pearl. Effects of treatment on the treated: Identification and generalization. In *Proceedings of the Twenty-Fifth Conference on Uncertainty in Artificial Intelligence*, pages 514–521. AUAI Press, Montreal, Quebec, 2009.
- [Tian and Pearl, 2000] Jin Tian and Judea Pearl. Probabilities of causation: Bounds and identification. *Annals of Mathematics and Artificial Intelligence*, 28(1-4):287–313, 2000.

## 8 Appendix

### 8.1 Proof of Theorem 10

*Proof.* From Equation (2) we have,

$$P(x, y) \leq P(y_x) \leq 1 - P(x, y').$$

Let  $1 - P(x, y') - P(x, y) \leq 2\epsilon$ , we obtain  $P(x') \leq 2\epsilon$ . Therefore,  $P(y_x)$  is  $\epsilon$ -identified to  $P(x, y) + \epsilon$  if  $P(x') \leq 2\epsilon$ , Equation (11) holds. Similarly, we can substitute  $x, y$  with  $x', y'$ , respectively. Equations (12) to (14) hold.  $\square$

### 8.2 Proof of Theorem 11

*Proof.* First, by adjustment formula in Equation (1), we have,

$$P(y_x) = P(y|x, u)P(u) + P(y|x, u')P(u').$$

Thus,

$$\begin{aligned} & P(y_x) \\ \geq & P(y|x, u')P(u') \\ = & P(y|x, u')(1 - P(u)) \\ = & \frac{P(x, y, u')}{P(x, u')}(1 - P(u)) \\ \geq & \frac{P(x, y) - P(u)}{P(x)}(1 - P(u)) \\ = & P(y|x) - P(y|x)P(u) - \frac{P(u)}{P(x)} + \frac{P^2(u)}{P(x)} \\ \geq & P(y|x) - P(u) - \frac{P(u)}{P(x)} \\ = & P(y|x) - (1 + \frac{1}{P(x)})P(u). \end{aligned}$$

Also if  $P(x) \geq P(u) + c$  for some constant  $c > 0$ , we have,

$$\begin{aligned} & P(y_x) \\ \leq & P(u) + P(y|x, u')(1 - P(u)) \\ \leq & P(u) + \frac{P(x, y, u')}{P(x, u')}(1 - P(u)) \\ \leq & P(u) + \frac{P(x, y)}{P(x) - P(u)}(1 - P(u)) \\ \leq & P(u) + \frac{P(x, y)}{P(x) - P(u)} \\ = & P(u) + \frac{P(x, y)}{P(x)(1 - \frac{P(u)}{P(x)})} \\ = & P(u) + \frac{P(x, y)(1 - \frac{P(u)}{P(x)}) + P(y|x)P(u)}{P(x)(1 - \frac{P(u)}{P(x)})} \\ = & P(u) + P(y|x) + \frac{P(y|x)P(u)}{P(x) - P(u)} \\ \leq & P(y|x) + P(u) + \frac{P(u)}{P(x) - P(u)} \\ \leq & P(y|x) + P(u) + \frac{P(u)}{c} \\ = & P(y|x) + P(u)(1 + \frac{1}{c}) \end{aligned}$$

Therefore, we have,

$$P(y|x) - (1 + \frac{1}{P(x)})P(u) \leq P(y_x) \leq P(y|x) + (1 + \frac{1}{c})P(u).$$

Let

$$(1 + \frac{1}{c})P(u) + (1 + \frac{1}{P(x)})P(u) \leq 2\epsilon.$$

We have,

$$\begin{aligned} & P(u) \\ \leq & \frac{2}{2 + \frac{1}{c} + \frac{1}{P(x)}}\epsilon \\ = & \frac{2cP(x)}{2cP(x) + P(x) + c}\epsilon. \end{aligned}$$

Then we know that if  $P(u) \leq \frac{2cP(x)}{2cP(x) + P(x) + c}\epsilon$ ,

$$\begin{aligned} P(y|x) - (1 + \frac{1}{P(x)})\frac{2cP(x)}{2cP(x) + P(x) + c}\epsilon & \leq P(y_x), \\ P(y|x) + (1 + \frac{1}{c})\frac{2cP(x)}{2cP(x) + P(x) + c}\epsilon & \geq P(y_x), \\ P(y|x) - \frac{2cP(x) + 2c}{2cP(x) + P(x) + c}\epsilon & \leq P(y_x), \\ P(y|x) + \frac{2cP(x) + 2P(x)}{2cP(x) + P(x) + c}\epsilon & \geq P(y_x). \end{aligned}$$

Therefore,  $P(y_x)$  is  $\epsilon$ -identified to  $P(y|x) - \frac{2cP(x) + 2c}{2cP(x) + P(x) + c}\epsilon + \epsilon = P(y|x) + \frac{P(x) - c}{2cP(x) + P(x) + c}\epsilon$ .

Besides, if  $P(x) \geq 0.5$  and  $P(u) \leq 0.1$ , let  $c = 0.4$ , we have

$$\begin{aligned} P(y|x) - (1 + \frac{1}{P(x)})P(u) & \leq P(y_x), \\ P(y|x) + (1 + \frac{1}{c})P(u) & \geq P(y_x), \\ P(y|x) - (1 + \frac{1}{0.5})P(u) & \leq P(y_x), \\ P(y|x) + (1 + \frac{1}{0.4})P(u) & \geq P(y_x). \end{aligned}$$

$$P(y|x) - 3P(u) \leq P(y_x) \leq P(y|x) + 3.5P(u).$$

Let  $3.5P(u) + 3P(u) \leq 2\epsilon$ , we have  $P(u) \leq \frac{4}{13}\epsilon$ , and

$$P(y|x) - \frac{12}{13}\epsilon \leq P(y_x) \leq P(y|x) + \frac{14}{13}\epsilon.$$

Therefore,  $P(y_x)$  is  $\epsilon$ -identified to  $P(y|x) - \frac{12}{13}\epsilon + \epsilon = P(y|x) + \frac{\epsilon}{13}$ .  $\square$

### 8.3 Proof of Theorem 12

*Proof.* From the bounds of PNS in Equations (3) and (4) is as follows:

$$\begin{aligned} & \max \left\{ \begin{array}{l} 0, \\ P(y_x) - P(y_{x'}), \\ P(y) - P(y_{x'}), \\ P(y_x) - P(y) \end{array} \right\} \leq \text{PNS} \\ & \min \left\{ \begin{array}{l} P(y_x), \\ P(y_{x'}), \\ P(x, y) + P(x', y'), \\ P(y_x) - P(y_{x'}) + \\ + P(x, y') + P(x', y) \end{array} \right\} \geq \text{PNS}. \end{aligned}$$



Let  $P(y_x) - 0 \leq 2\epsilon$ , we obtain that PNS is  $\epsilon$ -identified to  $\epsilon$  if  $P(y_x) \leq 2\epsilon$ , Equation (19) holds. □

Similarly, the rest of 20 equations can be obtained by letting

$$\begin{aligned}
P(y'_{x'}) - 0 &\leq 2\epsilon, \\
P(x, y) + P(x', y') - 0 &\leq 2\epsilon, \\
P(y_x) - P(y_{x'}) + P(x, y') + P(x', y) - 0 &\leq 2\epsilon, \\
P(y_x) - (P(y_x) - P(y_{x'})) &\leq 2\epsilon, \\
P(y'_{x'}) - (P(y_x) - P(y_{x'})) &\leq 2\epsilon, \\
P(x, y) + P(x', y') - (P(y_x) - P(y_{x'})) &\leq 2\epsilon, \\
P(y_x) - P(y_{x'}) + P(x, y') + P(x', y) - \\
(P(y_x) - P(y_{x'})) &\leq 2\epsilon, \\
P(y_x) - (P(y) - P(y_{x'})) &\leq 2\epsilon, \\
P(y'_{x'}) - (P(y) - P(y_{x'})) &\leq 2\epsilon, \\
P(x, y) + P(x', y') - (P(y) - P(y_{x'})) &\leq 2\epsilon, \\
P(y_x) - P(y_{x'}) + P(x, y') + P(x', y) - \\
(P(y) - P(y_{x'})) &\leq 2\epsilon, \\
P(y_x) - (P(y_x) - P(y)) &\leq 2\epsilon, \\
P(y'_{x'}) - (P(y_x) - P(y)) &\leq 2\epsilon, \\
P(x, y) + P(x', y') - (P(y_x) - P(y)) &\leq 2\epsilon, \\
P(y_x) - P(y_{x'}) + P(x, y') + P(x', y) - \\
(P(y_x) - P(y)) &\leq 2\epsilon.
\end{aligned}$$

□

#### 8.4 Proof of Theorem 13

*Proof.* From the bounds of PN in Equations (5) and (6) is as follows:

$$\max \left\{ \frac{0, P(y) - P(y_{x'})}{P(x, y)} \right\} \leq \text{PN} \leq \min \left\{ \frac{1, P(y'_{x'}) - P(x', y')}{P(x, y)} \right\}$$

Let  $\frac{P(y'_{x'}) - P(x', y')}{P(x, y)} - 0 \leq 2\epsilon$ , we obtain that PN is  $\epsilon$ -identified to  $\epsilon$  if  $P(y'_{x'}) - P(x', y') \leq 2P(x, y)\epsilon$ , Equation (40) holds. Similarly, the rest of 4 equations can be obtained by letting

$$\begin{aligned}
1 - \frac{P(y) - P(y_{x'})}{P(x, y)} &\leq 2\epsilon, \\
\frac{P(y'_{x'}) - P(x', y')}{P(x, y)} - \frac{P(y) - P(y_{x'})}{P(x, y)} &\leq 2\epsilon.
\end{aligned}$$

□

#### 8.5 Proof of Theorem 14

*Proof.* From the bounds of PS in Equations (7) and (8) is as follows:

$$\max \left\{ \frac{0, P(y') - P(y'_x)}{P(x', y')} \right\} \leq \text{PS} \leq \min \left\{ \frac{1, P(y_x) - P(x, y)}{P(x', y')} \right\}$$

Let  $\frac{P(y_x) - P(x, y)}{P(x', y')} - 0 \leq 2\epsilon$ , we obtain that PS is  $\epsilon$ -identified to  $\epsilon$  if  $P(y_x) - P(x, y) \leq 2P(x', y')\epsilon$ , Equation (45). Similarly, the rest of 4 conditions can be obtained by letting

$$\begin{aligned}
1 - \frac{P(y') - P(y'_x)}{P(x', y')} &\leq 2\epsilon, \\
\frac{P(y_x) - P(x, y)}{P(x', y')} - \frac{P(y') - P(y'_x)}{P(x', y')} &\leq 2\epsilon.
\end{aligned}$$

#### 8.6 Proof of Theorem 15

*Proof.*

$$\begin{aligned}
f(c) &= \beta P(y_x, y'_{x'}|c) + \gamma P(y_x, y_{x'}|c) + \\
&\theta P(y'_{x'}, y'_{x'}|c) + \delta P(y'_{x'}, y_{x'}|c) \\
&= \beta P(y_x, y'_{x'}|c) + \gamma [P(y_x|c) - P(y_x, y'_{x'}|c)] + \\
&\theta [P(y'_{x'}) - P(y_x, y'_{x'}|c)] + \delta P(y'_{x'}, y_{x'}|c) \\
&= \gamma P(y_x|c) + \theta P(y'_{x'}|c) + (\beta - \gamma - \theta) P(y_x, y'_{x'}|c) + \\
&\delta P(y'_{x'}, y_{x'}|c). \tag{50}
\end{aligned}$$

Note that, we have,

$$P(y'_{x'}, y_{x'}|c) = P(y_x, y'_{x'}|c) - P(y_x|c) + P(y_{x'}|c). \tag{51}$$

Substituting Equation (51) into Equation (50), we have,

$$\begin{aligned}
f(c) &= \gamma P(y_x|c) + \theta P(y'_{x'}|c) + (\beta - \gamma - \theta) P(y_x, y'_{x'}|c) + \\
&\delta P(y'_{x'}, y_{x'}|c) \\
&= \gamma P(y_x|c) + \theta P(y'_{x'}|c) + (\beta - \gamma - \theta) P(y_x, y'_{x'}|c) + \\
&\delta [P(y_x, y'_{x'}|c) - P(y_x|c) + P(y_{x'}|c)] \\
&= (\gamma - \delta) P(y_x|c) + \delta P(y_{x'}|c) + \theta P(y'_{x'}|c) + \\
&(\beta - \gamma - \theta + \delta) P(y_x, y'_{x'}|c).
\end{aligned}$$

Case 1: If  $\beta - \gamma - \theta + \delta \geq 0$ ,

$$\begin{aligned}
f(c) &\leq (\gamma - \delta) P(y_x|c) + \delta P(y_{x'}|c) + \theta P(y'_{x'}|c) + \\
&\frac{\beta - \gamma - \theta + \delta}{2} + \frac{|\beta - \gamma - \theta + \delta|}{2} \\
&= (\gamma - \delta) P(y_x|c) + \delta P(y_{x'}|c) + \theta P(y'_{x'}|c) + \\
&\beta - \gamma - \theta + \delta.
\end{aligned}$$

and,

$$\begin{aligned}
f(c) &\geq (\gamma - \delta) P(y_x|c) + \delta P(y_{x'}|c) + \theta P(y'_{x'}|c) + \\
&\frac{\beta - \gamma - \theta + \delta}{2} - \frac{|\beta - \gamma - \theta + \delta|}{2} \\
&= (\gamma - \delta) P(y_x|c) + \delta P(y_{x'}|c) + \theta P(y'_{x'}|c).
\end{aligned}$$

Therefore,  $f(c)$  is  $\frac{|\beta - \gamma - \theta + \delta|}{2}$ -identified to  $(\gamma - \delta) P(y_x|c) + \delta P(y_{x'}|c) + \theta P(y'_{x'}|c) + \frac{\beta - \gamma - \theta + \delta}{2}$ .

Case 2: If  $\beta - \gamma - \theta + \delta < 0$ ,

$$\begin{aligned}
f(c) &\leq (\gamma - \delta) P(y_x|c) + \delta P(y_{x'}|c) + \theta P(y'_{x'}|c) + \\
&\frac{\beta - \gamma - \theta + \delta}{2} + \frac{|\beta - \gamma - \theta + \delta|}{2} \\
&= (\gamma - \delta) P(y_x|c) + \delta P(y_{x'}|c) + \theta P(y'_{x'}|c).
\end{aligned}$$

and,

$$\begin{aligned}
f(c) &\geq (\gamma - \delta) P(y_x|c) + \delta P(y_{x'}|c) + \theta P(y'_{x'}|c) + \\
&\frac{\beta - \gamma - \theta + \delta}{2} - \frac{|\beta - \gamma - \theta + \delta|}{2} \\
&= (\gamma - \delta) P(y_x|c) + \delta P(y_{x'}|c) + \theta P(y'_{x'}|c) + \\
&\beta - \gamma - \theta + \delta.
\end{aligned}$$

Therefore,  $f(c)$  is  $\frac{|\beta-\gamma-\theta+\delta|}{2}$ -identified to  $(\gamma - \delta)P(y_x|c) + \delta P(y_{x'}|c) + \theta P(y_{x'}|c) + \frac{\beta-\gamma-\theta+\delta}{2}$ .  $\square$