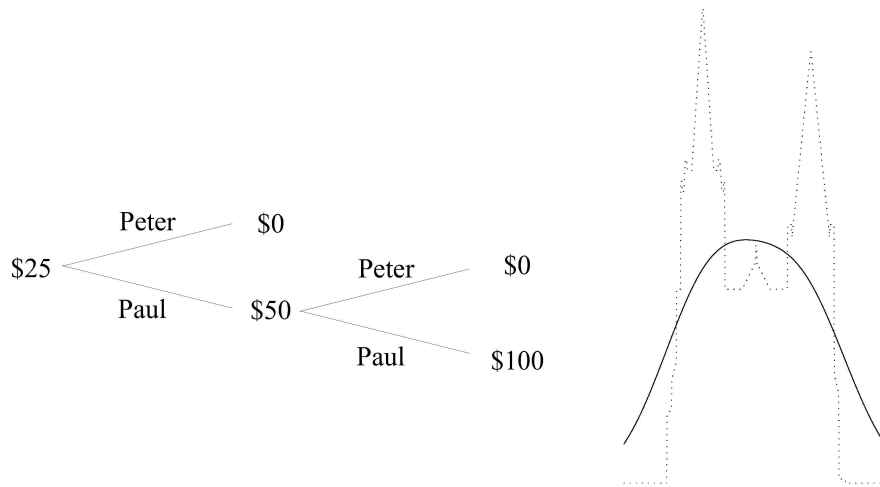


# Confidence intervals for causal effects in sequential decision making

Vladimir Vovk and Ruodu Wang



**The Game-Theoretic Probability and Finance Project**

Working Paper #68

First posted May 25, 2026. Last revised June 25, 2026.

Project web site:

<http://www.probabilityandfinance.com>

# Abstract

We derive confidence intervals and confidence sequences for causal effects in situations where the back-door criterion is applicable. Our tightest confidence intervals hold in the standard setting where the training data consists of IID observations over a system described by a given causal diagram. When interventions are allowed to depend on the past data, our confidence intervals become wider and involve a term coming from the law of the iterated logarithm, even where the number of observations is known in advance. In the sequential setting where the number of observations is not given, our confidence intervals, arranged into a confidence sequence for causal effects, involve more iterated logarithm terms and become even wider.

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Causal effect</b>	<b>1</b>
<b>3</b>	<b>The IID setting</b>	<b>3</b>
<b>4</b>	<b>The adaptive setting with a fixed horizon</b>	<b>5</b>
<b>5</b>	<b>The anytime-valid adaptive setting</b>	<b>6</b>
<b>6</b>	<b>The necessity of the iterated logarithm term in the adaptive setting</b>	<b>7</b>
<b>7</b>	<b>Applications to prediction sets</b>	<b>8</b>
<b>8</b>	<b>Conclusion</b>	<b>9</b>
	<b>References</b>	<b>10</b>
<b>A</b>	<b>Proofs for Sects. 3–5</b>	<b>11</b>
<b>B</b>	<b>Proofs for Sect. 6</b>	<b>16</b>

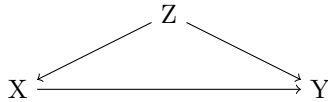


Figure 1: The basic causal graph of this paper

## 1 Introduction

A major limitation of many results in causal inference is that they assume, implicitly or explicitly, IID (independent and identically distributed) observations over a causal system. This limitation is shared by [16], where we derive prediction sets in situations covered by the back-door and front-door criteria [8, Theorems 3.3.2 and 3.3.4] (one section of [16] goes beyond the IID picture but only slightly). It was noted in [15] that the limitation can be easily overcome by applying limit theorems of probability theory, but the results there (such as [15, Theorem 1]) are asymptotic. In this paper we derive finite-sample confidence intervals in natural non-IID settings.

We start in Sect. 2 by giving the definition of the causal effect adapted to the back-door and front-door criteria. In Sect. 3 we consider the simplest IID setting giving the narrowest confidence intervals. Limitations of this setting are pointed out in, e.g., [15, Sect. 1] and [8, end of Sect. 3.6.1].

In Sect. 4 we drop the assumption of IID data, which leads to the appearance of an iterated logarithm term in our confidence intervals. This is developed further in Sect. 5; there we make our confidence intervals anytime-valid, which leads to further iterated logarithm terms. In Sects. 4–5 we only consider the back-door criterion.

This paper was motivated by the difficulty of applying the methods developed in [16] to the most natural setting of sequential causal inference (which we called the “strong interpretation” of causal diagrams). The methods of this paper are completely different, and we are targeting confidence intervals rather than prediction sets, which were targeted in [16]. In Sect. 7 we briefly discuss derivation of prediction sets from our confidence intervals, but this is likely to lead to much more conservative prediction sets in the IID setting as compared with [16].

Finally, Sect. 8 concludes and lists some directions of further research.

## 2 Causal effect

Our running example will be the causal diagram in Figure 1, which we now use to explain our notation (in which we follow mainly Pearl [8]). The variables, such as  $X$ ,  $Y$ , and  $Z$ , in our causal diagrams will always range over finite sets denoted by the corresponding boldface letters, such as  $\mathbf{X}$ ,  $\mathbf{Y}$ , and  $\mathbf{Z}$ , and called their *domains* (equipped with the discrete  $\sigma$ -algebras). However, when

talking specifically about the example in Figure 1, we will usually assume that  $\mathbf{X} = \mathbf{Y} = \mathbf{Z} = \{0, 1\}$ , so that  $X, Y, Z$  are the indicator functions of events.

Let  $P$  be a positive probability measure on  $\mathbf{X} \times \mathbf{Y} \times \mathbf{Z}$  (generating the random variables  $X, Y$ , and  $Z$ ). Suppose it factorizes according to Figure 1:

$$\begin{aligned} P(X = x, Y = y, Z = z) \\ = P(Z = z)P(X = x | Z = z)P(Y = y | X = x, Z = z) \end{aligned} \quad (1)$$

for all  $x \in \mathbf{X}$ ,  $y \in \mathbf{Y}$ , and  $z \in \mathbf{Z}$ . It will be very convenient to use Pearl's [8, Sects. 1.1.4 and 1.1.5] convention and abbreviate  $P(X = x)$  to  $P(x)$ ,  $P(Y = y)$  to  $P(y)$ , etc.; such abbreviated notation will also be used when we have, say,  $\tilde{x}$  or  $x'$  in place of  $x$ . We will also often omit mentioning that  $x \in \mathbf{X}$ ,  $y \in \mathbf{Y}$ , etc. With this convention, we can rewrite (1) as

$$P(x, y, z) = P(z)P(x | z)P(y | x, z).$$

Let  $\tilde{x} \in \mathbf{X}$ . We use Pearl's [8] notation  $\text{do}(X = \tilde{x})$ , usually abbreviated to  $\text{do}(\tilde{x})$ , to signify setting  $X$  to  $\tilde{x}$  (we will define what this means formally only in specific contexts). Let us define, in the context of Figure 1, the *causal effect* of  $X$  on  $Y$  as

$$P(y | \text{do}(\tilde{x})) := \sum_z P(y | \tilde{x}, z)P(z). \quad (2)$$

(The definition in [8, Sect. 3.2] is more general, but it is not our focus in this paper and we will use simpler *ad hoc* definitions.) The interpretation of (2) (and of causal effects in general) is that it is the probability of  $Y = y$  in the mutilated causal model in which the arrow from  $Z$  to  $X$  in Figure 1 has been removed and  $X$  has been set to  $\tilde{x}$ .

The decomposition (1) and these notational conventions generalize to any directed acyclic graph (dag), and Figure 1 can be generalized to the following *back-door criterion*, which is stated in terms of "blocking", as defined in [8, Definition 1.2.3]. If  $X, Y$ , and  $Z$  are disjoint non-empty sets of variables in a dag,  $Z$  is said to satisfy the back-door criterion relative to  $(X, Y)$  if, for any  $X' \in X$  and  $Y' \in Y$ ,

- no vertex in  $Z$  is a descendant of  $X'$ , and
- $Z$  blocks every *back-door path* from  $X'$  to  $Y'$ , i.e., every path between  $X'$  and  $Y'$  that contains an arrow into  $X'$

[8, Definition 3.3.1]. If the back-door criterion is satisfied, the *causal effect* can still be defined as (2) [8, Theorem 3.3.2]. However, now the summing  $\sum_z$  over  $z$  in (2) means summing over all possible values of the variables in  $Z$ :

$$\sum_z := \sum_{z_1 \in \mathbf{Z}^1} \cdots \sum_{z_k \in \mathbf{Z}^k},$$

where  $\mathbf{Z}^i$  is the domain of  $Z^i$ ,  $i = 1, \dots, k$ , and  $Z^1, \dots, Z^k$  are the elements of  $Z$ ,  $Z = \{Z^1, \dots, Z^k\}$ . Moreover,  $\tilde{x}$  should specify the values for all variables in

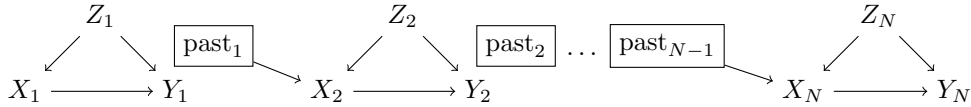


Figure 2: The repeated causal graph

$X$ . In the theorems below we will use the notation

$$|\mathbf{Z}| := |\mathbf{Z}^1| \dots |\mathbf{Z}^k|. \quad (3)$$

Another set of conditions that allows a natural (albeit not obvious) definition of causal effects is known as the *front-door criterion*. Now let  $X$  and  $Y$  be variables and  $Z$  be a non-empty set of variables not containing  $X$  and  $Y$ . Then we say that  $Z$  satisfies the front-door criterion relative to  $(X, Y)$  if

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

[8, Definition 3.3.3]. In this case the *causal effect* is defined by

$$P(y | \text{do}(\tilde{x})) := \sum_z P(z | \tilde{x}) \sum_x P(y | x, z) P(x) \quad (4)$$

[8, Theorem 3.3.4].

It is very important (but does not concern us in this paper) that some of the variables in the causal dag may be unobservable; it's fine as long as these variables do not enter expressions for causal effects, such as (2) or (4).

### 3 The IID setting

We start from the most standard setting where the observations before intervention are IID; see, e.g., [8, the beginning of Sect. 2.2]. Informally, Nature possesses stable causal mechanisms that are organized in the form of a graphical structure. In causal calculus, the structure is a known dag, and the stable causal mechanisms are unknown probability distributions of the variables in the vertices of the dag given their parents. The available observations are generated in the IID fashion. (The IID nature of the observations usually stays implicit, and it appears that in Pearl's book [8] it is made explicitly only in [8, the end of Sect. 3.6.1].)

Figure 2 shows  $N$  repetitions of the causal system represented in Figure 1. Let us ignore the cells labelled  $\text{past}_i$ ,  $i \in \{1, \dots, N-1\}$ , for now (formally, we are assuming that these variables take a fixed known value). This composite causal diagram then represents  $N$  IID observations over the base diagram of

Figure 1. In this section we are interested in Figure 2 with an arbitrary dag as the base diagram, not necessarily the one in Figure 1. This IID picture will give the narrowest confidence intervals out of those derived in this paper.

The following theorem treats the case of the back-door criterion, and it is proved (as all other results in Sects. 3–5) in Appendix A. A confidence interval  $[m - h, m + h]$  will be represented in terms of its *midpoint*  $m$  and *half-width*  $h$ . The number  $N$  of repetitions is fixed, and we set

$$\#xy := |\{n \in [N] : (X_n, Y_n) = (x, y)\}|, \quad (5)$$

where  $[N] := \{1, \dots, N\}$ ; i.e.,  $\#xy$  is the number of times  $(X_n, Y_n) = (x, y)$ . The analogous notation will be used for sequences other than  $xy$ , such as  $xyz$  and  $z$ . See (3) for the definition of  $|\mathbf{Z}|$ .

**Theorem 1.** *Suppose the back-door criterion is satisfied. Let  $\delta > 0$ ,  $\tilde{x} \in \mathbf{X}$ , and  $y \in \mathbf{Y}$ . The following is a  $(1 - \delta)$ -confidence interval for the parameter  $P(y \mid \text{do}(\tilde{x}))$  defined by (2): the midpoint is*

$$\sum_z \hat{p}(y \mid \tilde{x}, z) \hat{p}(z) \quad (6)$$

and the half-width is

$$|\mathbf{Z}| \sqrt{\frac{\ln \frac{4|\mathbf{Z}|}{\delta}}{2N}} + \sum_z \sqrt{\frac{\ln \frac{4|\mathbf{Z}|}{\delta}}{2\#\tilde{x}z}}, \quad (7)$$

where  $\hat{p}(y \mid \tilde{x}, z) := \#\tilde{x}yz / \#\tilde{x}z$  is the standard estimate for  $P(y \mid \tilde{x}, z)$  and  $\hat{p}(z) := \#z/N$  is the standard estimate for  $P(z)$ . (The half-width (7) is understood to be  $\infty$  when  $\#\tilde{x}z = 0$  for some  $z$ .)

In the case of Figure 1 with binary variables, we can replace (7) by

$$2\sqrt{\frac{\ln \frac{6}{\delta}}{2N}} + \sum_{z \in \{0,1\}} \sqrt{\frac{\ln \frac{6}{\delta}}{2\#\tilde{x}z}} \quad (8)$$

(although this does not quite follow from (7)).

The following is the analogue of Theorem 1 for the front-door criterion.

**Theorem 2.** *Suppose the front-door criterion is satisfied. Fix  $\delta > 0$ ,  $\tilde{x}$ , and  $y$ . The following is a  $(1 - \delta)$ -confidence interval for the parameter  $P(y \mid \text{do}(\tilde{x}))$  defined by (4): the midpoint is*

$$\sum_z \hat{p}(z \mid \tilde{x}) \sum_x \hat{p}(y \mid x, z) \hat{p}(x) \quad (9)$$

and the half-width is

$$|\mathbf{X}| |\mathbf{Z}| \sqrt{\frac{\ln \frac{2K}{\delta}}{2N}} + |\mathbf{Z}| \sqrt{\frac{\ln \frac{2K}{\delta}}{2\#\tilde{x}}} + \sum_{x,z} \sqrt{\frac{\ln \frac{2K}{\delta}}{2\#xz}}, \quad (10)$$

where

$$K := |\mathbf{X}| |\mathbf{Z}| + |\mathbf{X}| + |\mathbf{Z}| = (|\mathbf{X}| + 1)(|\mathbf{Z}| + 1) - 1,$$

$\hat{p}(z | \tilde{x}) := \#\tilde{x}z/\#\tilde{x}$  is the standard estimate for  $P(z | \tilde{x})$ ,  $\hat{p}(y | x, z) := \#xyz/\#xz$  is the standard estimate for  $P(y | x, z)$ , and  $\hat{p}(x) := \#x/N$  is the standard estimate for  $P(x)$ .

For the same size  $|\mathbf{Z}|$ , the accuracy (10) that we have for the front-door criterion appears much worse than the accuracy (7) for the back-door one.

## 4 The adaptive setting with a fixed horizon

Under the *strong interpretation* of Figure 2, considered in this section, each box  $\text{past}_n$  stands for the whole past, including the variables  $X_i, Y_i$ , and  $Z_i, i \in [n]$ . Now each  $X_{n+1}, n \in [N - 1]$ , has incoming arrows (not shown explicitly in the figure) from all variables at the previous steps, including  $X_i, Y_i$ , and  $Z_i, i \in [n]$ . As before, we allow repetition of any dag in Figure 2, not just the one in Figure 1.

Our interpretation of Figure 2 is that  $X$  is a decision that has  $Y$  as its result. The decision at step  $n + 1$  may depend on the past decisions and past values of  $Y, Z$ , and other variables. In other words, the decision maker has access to all past observations. In the case of Figure 1, all  $Z_n$  are independent of the past and identically distributed; all  $Y_n$  have the same distribution given  $X_n$  and  $Z_n$ , and they are conditionally independent of the past.

For an integer  $n \geq 2$ , let  $\lfloor\!\!\lfloor n \rfloor\!\!\rfloor$  be the largest integer of the form  $2^k, k \in \{1, 2, \dots\}$ , satisfying  $2^k \leq n$  (and for  $n < 2$ ,  $\lfloor\!\!\lfloor n \rfloor\!\!\rfloor$  is defined as, say, 1). We let  $\text{lb}$  stand for binary logarithm  $\log_2$ , and we will often use it in the context of  $\text{lb}\lfloor\!\!\lfloor n \rfloor\!\!\rfloor = \lfloor \text{lb } n \rfloor$  for a positive integer  $n$ .

It will be useful to extend the notation (5) and set, e.g.,

$$\#_m xy := |\{n \in [m] : (X_n, Y_n) = (x, y)\}|,$$

where  $m$  may be different from  $N$ . Earlier we defined standard estimates such as  $\hat{p}(y | \tilde{x}, z) := \#\tilde{x}yz/\#\tilde{x}z$  (and we will refrain from defining  $\hat{p}$  in other similar contexts in the following theorems). We will also need the modification of  $\hat{p}(y | \tilde{x}, z)$  defined by

$$\hat{p}(y | \tilde{x}, z) := \frac{|\{n \in [N] : \#\tilde{x}nz \leq \lfloor\!\!\lfloor \#\tilde{x}z \rfloor\!\!\rfloor, X_n = \tilde{x}, Y_n = y, Z_n = z\}|}{\lfloor\!\!\lfloor \#\tilde{x}z \rfloor\!\!\rfloor}. \quad (11)$$

In words,  $\hat{p}(y | \tilde{x}, z)$  is the fraction of the first  $\lfloor\!\!\lfloor \#\tilde{x}z \rfloor\!\!\rfloor$  observations with  $X_n = \tilde{x}$  and  $Z_n = z$  for which  $Y_n = y$ . It is also an estimate of  $P(y | \tilde{x}, z)$ , but it might not use all the available data (however, it uses at least one half of the relevant observations). We will also use the notation  $\hat{p}$  in other contexts, such as  $\hat{p}(z | \tilde{x})$ .

Now we have to replace the Hoeffding inequality used in our proof of Theorem 1 in Appendix A by the law of the iterated logarithm, and so we will get our first iterated logarithm term.

**Theorem 3.** *Suppose the back-door criterion is satisfied. Fix  $\delta > 0$ ,  $\tilde{x}$ , and  $y$ ; the time horizon  $N$  is also fixed. Under the strong interpretation, the following is a  $(1 - \delta)$ -confidence interval for the parameter  $P(y \mid \text{do}(\tilde{x}))$  defined by (2): the midpoint is*

$$\sum_z \hat{p}(z) \hat{p}(y \mid \tilde{x}, z) \quad (12)$$

and the half-width is

$$|\mathbf{Z}| \sqrt{\frac{\ln \frac{4|\mathbf{Z}|}{\delta}}{2N}} + \sum_z \sqrt{\frac{2 \ln \text{lb} \llbracket \# \tilde{x} z \rrbracket + \ln \frac{6.6|\mathbf{Z}|}{\delta}}{2 \llbracket \# \tilde{x} z \rrbracket}}. \quad (13)$$

The case  $\# \tilde{x} z \in \{0, 1\}$  in (13) requires special treatment; namely, we set  $\ln \text{lb } 1 := \infty$ , and so (13) is interpreted as  $\infty$  unless  $\# \tilde{x} z \geq 2$  for all  $z$ .

For the binary case of Figure 1, we can replace (13) by

$$2 \sqrt{\frac{\ln \frac{6}{\delta}}{2N}} + \sum_{z \in \{0, 1\}} \sqrt{\frac{2 \ln \text{lb} \llbracket \# \tilde{x} z \rrbracket + \ln \frac{10}{\delta}}{2 \llbracket \# \tilde{x} z \rrbracket}}, \quad (14)$$

similarly to (8).

It is not clear how to extend Theorem 3 to the front-door criterion. Moreover, it is not even obvious that the quantity (4) is well-defined under the strong interpretation, since, e.g., the distribution of  $X$  changes from step to step. (For a demonstration that it is indeed well-defined see, e.g., [8, (3.27)].)

## 5 The anytime-valid adaptive setting

Theorem 3 can be easily extended to the setting in which the time horizon  $N$  is not fixed in advance. Now we would like our results to be anytime valid, with  $N$  ranging over the positive integers  $\{1, 2, \dots\}$ . Since  $N$  is variable, now we will write  $\hat{p}_N(y \mid \tilde{x}, z)$  in place of  $\hat{p}(y \mid \tilde{x}, z)$  defined by (11) and add the lower index  $N$  in other similar places; in particular,  $\hat{p}_N(z) := \# \llbracket N \rrbracket z / \llbracket N \rrbracket$ . Remember that a  $(1 - \delta)$ -confidence sequence is a sequence of (confidence) intervals whose intersection covers the true parameter value with probability at least  $1 - \delta$ .

**Theorem 4.** *Suppose the back-door criterion is satisfied. Let  $\delta > 0$ ,  $\tilde{x} \in \mathbf{X}$ , and  $y \in \mathbf{Y}$ . The following is a  $(1 - \delta)$ -confidence sequence for the parameter  $P(y \mid \text{do}(\tilde{x}))$  defined by (2): the midpoint is*

$$\sum_z \hat{p}_N(z) \hat{p}_N(y \mid \tilde{x}, z)$$

and the half-width is

$$|\mathbf{Z}| \sqrt{\frac{2 \ln \text{lb} \llbracket N \rrbracket + \ln \frac{6.6|\mathbf{Z}|}{\delta}}{2 \llbracket N \rrbracket}} + \sum_z \sqrt{\frac{2 \ln \text{lb} \llbracket \#_N \tilde{x} z \rrbracket + \ln \frac{6.6|\mathbf{Z}|}{\delta}}{2 \llbracket \#_N \tilde{x} z \rrbracket}}. \quad (15)$$

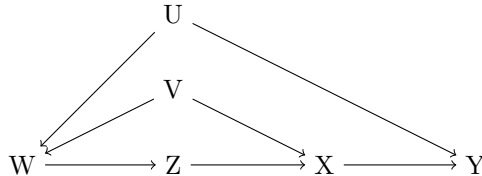


Figure 3: The napkin graph

In the binary case of Figure 1, (15) can be slightly strengthened to

$$2\sqrt{\frac{2 \ln \text{lb}\llbracket N \rrbracket + \ln \frac{10}{\delta}}{2\llbracket N \rrbracket}} + \sum_{z \in \{0,1\}} \sqrt{\frac{2 \ln \text{lb}\llbracket \#_N \tilde{x} z \rrbracket + \ln \frac{10}{\delta}}{2\llbracket \#_N \tilde{x} z \rrbracket}}.$$

Theorem 3–4 may be considered to be finite-sample analogues of Theorem 1 in [15]. Their characteristic feature is the presence of iterated logarithm terms, which are unavoidable (see Sect. 6 for details) and are especially prominent in the anytime-valid adaptive setting.

**Remark 5.** In this paper we only discuss, outside of this remark, causal effects that are representable as arithmetic expressions involving only two arithmetic operations, plus and multiplication (minus could be added for free but is not useful). There are, however, situations in which the causal effect is given in a form involving division, such as the *napkin graph*, shown in Figure 3. It has the following expression for the causal effect of  $X$  on  $Y$ :

$$P(y \mid \text{do}(\tilde{x})) := \frac{\sum_w P(y, \tilde{x} \mid z, w)P(w)}{\sum_w P(\tilde{x} \mid z, w)P(w)} \quad (16)$$

(see [4, Figure 2(c) and (1)]). The expression (16) is a ratio, and our methods encounter even more serious difficulties than in the case of the front-door criterion. An interesting feature of the expression (16) is that its right-hand side involves  $z$  but does not really depend on it, as is clear from its left-hand side; this is an instance of so-called Verma constraints [1]. The presence of  $z$  on the right-hand side of (16) is in a certain sense inevitable; formally,  $Z$  is a “trapdoor variable” as defined in [4, Definition 3] and explained in [4, Sect. 2.3].

## 6 The necessity of the iterated logarithm term in the adaptive setting

The next theorem (proved in Appendix B) shows that an iterated logarithm term is unavoidable already in the fixed-horizon setting of Theorem 3. We assume, without loss of generality,  $|\mathbf{Z}| = 1$ , and so we ignore  $\mathbf{Z}$ . Let us fix a

confidence level  $1 - \delta$  (it can be arbitrarily close to 0),  $\tilde{x} \in \mathbf{X}$ , and  $y \in \mathbf{Y}$ . We are interested in estimating the causal effect  $P(y \mid \text{do}(\tilde{x})) = P(y \mid \tilde{x})$ .

For a positive integer  $N$ , we consider a confidence estimator  $C_N$  for  $P(y \mid \tilde{x})$ ; formally,  $C_N$  maps a sequence  $(x_1, y_1, \dots, x_N, y_N) \in (\mathbf{X} \times \mathbf{Y})^N$  to a (closed) subinterval  $C_N(x_1, y_1, \dots, x_N, y_N)$  of  $[0, 1]$  and satisfies the following natural property of validity: regardless of the underlying probability measure  $P$ ,

$$P(y \mid \tilde{x}) \in C_N(X_1, Y_1, \dots, X_N, Y_N) \quad (17)$$

with probability at least  $1 - \delta$ . If  $C$  is an interval, we let  $|C|$  stand for its length.

**Theorem 6.** *Let  $(C_N)_{N=1}^\infty$  be a family of confidence estimators and let  $f : \{0, 1, \dots\} \rightarrow (0, \infty)$  satisfy*

$$f(n) = o\left(\sqrt{\frac{\ln \ln n}{n}}\right), \quad n \rightarrow \infty.$$

*If  $N$  is sufficiently large, there exists  $(x_1, y_1, \dots, x_N, y_N) \in (\mathbf{X} \times \mathbf{Y})^N$  such that*

$$|C_N(x_1, y_1, \dots, x_N, y_N)| > f(k) \quad (18)$$

*where  $k := \sum_{i=1}^N 1_{\{x_i = \tilde{x}\}}$  is the number of occurrences of  $\tilde{x}$ .*

## 7 Applications to prediction sets

Theorems 1–4 provide confidence intervals for causal effects, whereas in [16] we were interested in prediction sets for  $Y$ . In order to discuss connections between our results here and the strong interpretation in [16], in this section we will state a corollary of the toy version of Theorem 3 for the binary case of Figure 1, with (13) replaced by (14), giving prediction sets; similar corollaries can be easily deduced from Theorems 1–4 as well. We consider the strong interpretation of Figure 1, as in Sect. 4, with  $(X_n, Y_n, Z_n)$ ,  $n \in [N]$ , complemented by another observation  $Y$  with the probabilities of  $Y = y$ ,  $y \in \mathbf{Y}$ , given by the right-hand side of (2) for a fixed  $\tilde{x}$ . Remember that  $\mathbf{X} = \mathbf{Y} = \mathbf{Z} = \{0, 1\}$ .

**Corollary 7.** *Fix  $\delta > 0$ ,  $N$ , and  $\tilde{x} \in \mathbf{X}$ . Then*

$$\Gamma := \left\{ y \in \mathbf{Y} : \sum_{z \in \{0,1\}} \hat{p}(z) \hat{p}(y \mid \tilde{x}, z) + 2\sqrt{\frac{\ln \frac{12}{\delta}}{2N}} + \sum_{z \in \{0,1\}} \sqrt{\frac{2 \ln \text{lb} \llbracket \# \tilde{x} z \rrbracket + \ln \frac{20}{\delta}}{2 \llbracket \# \tilde{x} z \rrbracket}} > \frac{\delta}{2} \right\}$$

*is a  $(1 - \delta)$ -prediction set for  $Y$ .*

*Proof.* We are required to prove that  $Y \notin \Gamma$  with probability at most  $\delta$ . Let us fix  $y \in \mathbf{Y}$  and prove that the probability of the conjunction of  $Y = y$  and  $y \notin \Gamma$  is at most  $\delta/2$ . We will use the confidence interval (12)  $\pm$  (14) with  $\delta/2$  in place of  $\delta$ . Consider two cases:

- Suppose the right-hand-side of (2) exceeds  $\delta/2$ . If  $y \notin \Gamma$ , the right endpoint (12) + (14) (with  $\delta$  replaced by  $\delta/2$ ) of the confidence interval is at most  $\delta/2$ , and the probability of this is at most  $\delta/2$  by the definition of a confidence interval.
- Otherwise,  $Y = y$  with probability at most  $\delta/2$  by the definition of  $Y$ .  $\square$

This procedure is sub-optimal for several reasons. One of them is that Hoeffding’s inequality applied to the Bernoulli model can be greatly improved when the probability of error is small; see, e.g., Vapnik’s [12, Sects. 4.2 and 4.4] use of multiplicative Chernoff inequalities (in what he calls optimistic and pessimistic settings of learning problems).

The analogous corollary of Theorem 1 becomes more comparable with the results that we obtain in [16]. However, the prediction sets derived in [16] are based on e-values (are “e-prediction sets”) whereas the prediction sets here are traditional ones. We expect that methods of this paper lead to much looser results, but their advantage is that they also work for the strong interpretation.

## 8 Conclusion

In this paper we derive confidence intervals and confidence sequences for causal effects in IID and sequential non-IID settings. These are some directions of further research:

- How do we establish analogues of Theorems 3 and 4 for the front-door criterion?
- In this paper we concentrate on upper bounds for achievable widths of confidence intervals, and our results need to be complemented by lower bounds. Theorem 6 is a first step in this direction.
- This paper is based on the standard measure-theoretic probability [5, 10, 11]. To make our results as strong as possible, we could try and present them in the language of game-theoretic probability [9] (this was listed as direction of further research already in [15, Sect. 5]).

## Acknowledgments

In literature search and brainstorming we used OpenAI Codex; in particular, it noticed that Theorem 6 follows from the results of [3]. We take full responsibility for all statements made in this paper.

## References

- [1] Rohit Bhattacharya and Razieh Nabi. On testability of the front-door model via Verma constraints. *Proceedings of Machine Learning Research*, 180:202–212, 2022. UAI 2022.
- [2] Martine Ceberio and Vladik Kreinovich. Greedy algorithms for optimizing multivariate Horner schemes. *ACM SIGSAM Bulletin*, 38:8–15, 2004.
- [3] John Duchi and Saminul Haque. An information-theoretic lower bound in time-uniform estimation. *Proceedings of Machine Learning Research*, 247:1486–1500, 2024. COLT 2024.
- [4] Jouni Helske, Santtu Tikka, and Juha Karvanen. Estimation of causal effects with small data in the presence of trapdoor variables. *Journal of the Royal Statistical Society A*, 184:1030–1051, 2021.
- [5] Andrei N. Kolmogorov. *Grundbegriffe der Wahrscheinlichkeitsrechnung*. Springer, Berlin, 1933. English translation: *Foundations of the Theory of Probability*. Chelsea, New York, 1950.
- [6] J. Kuipers, A. Plaat, J. A. M. Vermaseren, and H. J. van den Herik. Improving multivariate Horner schemes with Monte Carlo tree search. *Computer Physics Communications*, 184:2391–2395, 2013.
- [7] Masashi Okamoto. Some inequalities relating to the partial sum of binomial probabilities. *Annals of the Institute of Statistical Mathematics*, 10:29–35, 1959.
- [8] Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge, Cambridge University Press, second edition, 2009.
- [9] Glenn Shafer and Vladimir Vovk. *Game-Theoretic Foundations for Probability and Finance*. Wiley, Hoboken, NJ, 2019.
- [10] Albert N. Shiryaev. *Probability-1*. Springer, New York, third edition, 2016.
- [11] Albert N. Shiryaev. *Probability-2*. Springer, New York, third edition, 2019.
- [12] Vladimir N. Vapnik. *Statistical Learning Theory*. Wiley, New York, 1998.
- [13] Jean Ville. *Etude critique de la notion de collectif*. Gauthier-Villars, Paris, 1939.
- [14] Václav Voráček. Treatment of statistical estimation problems in randomized smoothing for adversarial robustness. In *NeurIPS 2024*, 2024.
- [15] Vladimir Vovk. Another semantics for Pearl’s action calculus. In Alex Gammerman, editor, *Computational Learning and Probabilistic Reasoning*, chapter 7, pages 127–146. Wiley, Chichester, 1996.

[16] Vladimir Vovk and Ruodu Wang. Conformal e-prediction in the presence of confounding. Technical Report arXiv:2603.11134 [math.ST], arXiv.org e-Print archive, March 2026. For the latest version, see alrw.net, Working Paper 46.

## A Proofs for Sects. 3–5

In the following two propositions we consider an IID Bernoulli sequence  $\xi_1, \xi_2, \dots$  with probability of success  $p$  and use the notation  $\hat{p}_n := \frac{1}{n} \sum_{i=1}^n \xi_i$  for the standard estimate of  $p$ . We start from a standard confidence interval for the probability of success given by Hoeffding’s inequality.

**Proposition 8.** *For a fixed  $n$  and  $\delta > 0$ ,*

$$I := \left\{ p \in [0, 1] : |p - \hat{p}_n| < \sqrt{\frac{\ln \frac{2}{\delta}}{2n}} \right\} \quad (19)$$

is a  $(1 - \delta)$ -confidence interval for  $p$ , in the sense of

$$\mathbb{P}(p \in I) \geq 1 - \delta.$$

*Proof.* By Hoeffding’s inequality (or Okamoto’s earlier result [7, Theorem 1]), for all  $c > 0$ ,

$$\mathbb{P}(|p - \hat{p}_n| \geq c) \leq 2 \exp(-2c^2 n),$$

which gives the confidence interval (19).  $\square$

We will also need the following confidence sequence.

**Proposition 9.** *For each  $\delta > 0$ ,*

$$I_n := \left\{ p : |p - \hat{p}_{\lfloor n \rfloor}| < \sqrt{\frac{2 \ln \lfloor n \rfloor + \ln \frac{3.3}{\delta}}{2 \lfloor n \rfloor}} \right\} \quad (20)$$

is a  $(1 - \delta)$ -confidence sequence, i.e.,

$$\mathbb{P}(\forall n : p \in I_n) \geq 1 - \delta. \quad (21)$$

*Proof.* Similar confidence sequences can be obtained using Ville’s [13] method of continuous mixtures of test martingales or its discrete analogue [9, Sect. 5.1], but we will model our proof on [14, Sect. E]. Fix  $p \in [0, 1]$ .

Since  $\zeta(2) = \pi^2/6$ , we can split the significance level  $\delta$  into the series  $\delta = \sum_{k=1}^{\infty} \delta_k$ , where

$$\delta_k = \frac{6}{\pi^2 k^2} \delta.$$

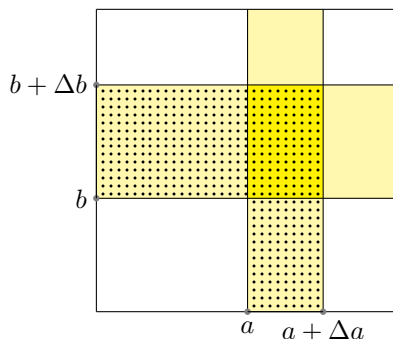


Figure 4: Illustration of an inequality.

Applying (19) with  $n_k = 2^k$  in place of  $n$  and  $\delta_k$  in place of  $\delta$  gives

$$\mathbb{P} \left( |p - \hat{p}_{n_k}| \geq \sqrt{\frac{\ln \frac{\pi^2 k^2}{3\delta}}{2n_k}} \right) \leq \delta_k. \quad (22)$$

Finally, (22) implies (21) since  $\pi^2/3 \approx 3.29 < 3.3$ . (This argument works for  $n \geq 2$ ; otherwise, the inequality in (21) is trivial since our convention, introduced in Sect. 4, is that  $\ln \text{lb } 1 := \infty$ .)  $\square$

Next we need a simple result from interval arithmetic. We are only interested in subintervals of  $[0, 1]$ . Let  $c \pm \Delta c$ , where  $c \in \mathbb{R}$  and  $\Delta c \geq 0$ , stand for the interval

$$c \pm \Delta c := [c - \Delta c, c + \Delta c] \cap [0, 1].$$

For a binary operation  $*$  on the reals (we are mostly interested in addition and multiplication), we define its result on intervals pointwise:

$$I_1 * I_2 := \{p_1 * p_2 : p_1 \in I_1, p_2 \in I_2\} \cap [0, 1].$$

**Lemma 10.** *For any two intervals  $a \pm \Delta a$  and  $b \pm \Delta b$ ,*

$$(a \pm \Delta a) + (b \pm \Delta b) \subseteq (a + b) \pm (\Delta a + \Delta b), \quad (23)$$

$$(a \pm \Delta a) \times (b \pm \Delta b) \subseteq (a \times b) \pm (\Delta a + \Delta b). \quad (24)$$

*Proof.* The inclusion (23) is obvious, so we will only prove (24). The latter inclusion reduces to the conjunction of two inequalities:

$$((a - \Delta a) \vee 0)((b - \Delta b) \vee 0) \geq ab - (\Delta a + \Delta b), \quad (25)$$

$$((a + \Delta a) \wedge 1)((b + \Delta b) \wedge 1) \leq ab + (\Delta a + \Delta b). \quad (26)$$

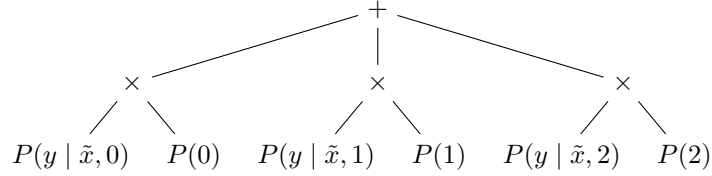


Figure 5: The binary tree representing the polynomial expression (2) over the formal variables  $P(y | \tilde{x}, z)$  and  $P(z)$  in the case of  $\mathbf{Z} = \{0, 1, 2\}$

In the inequality (25) we can assume, without loss of generality,  $a - \Delta a \geq 0$  and  $b - \Delta b \geq 0$ , which makes it obvious. In (26) we can assume, without loss of generality,  $a + \Delta a \leq 1$  and  $b + \Delta b \leq 1$ , which reduces it to

$$(a + \Delta a)(b + \Delta b) - ab \leq \Delta a + \Delta b. \quad (27)$$

The last inequality is illustrated in Figure 4: the dotted area of the plot represents the left-hand side of (27), and the yellow area represents the right-hand side of (27) (with the darker yellow area counted twice).  $\square$

We will distinguish between (multivariate) polynomials and polynomial expressions. A *polynomial expression* is formed from a finite number of formal variables by repeatedly applying the operations of multiplication and addition; we do not allow constants (equivalently, our polynomials and polynomial expressions are over the field  $\{0, 1\}$  and have a zero constant term). A polynomial expression can be represented as a tree (let us call it a *polynomial tree*), not necessarily binary, such as Figure 5 in the case of the polynomial expression given by the right-hand side of (2) for  $\mathbf{Z} = \{0, 1, 2\}$ . The same formal variable may be used several times. A *polynomial* is an equivalence class of polynomial expressions where we do not distinguish polynomial expressions that reduce to each other by applying the usual laws of commutativity, associativity, and distributivity (associativity was already used implicitly when we allowed non-binary multiplications and additions in our trees). Without loss of generality we may assume that the operations at different levels of polynomial trees alternate (so at each level we have the same operation, “+” or “ $\times$ ”, and perhaps some formal variables as leaf nodes; the operations in adjacent levels are different).

**Corollary 11.** *Let  $E$  be a polynomial expression involving  $m$  distinct formal variables. Define  $E^*$  to be the polynomial obtained from  $E$  by replacing all multiplications by additions. Then, for any intervals  $a_i \pm \Delta a_i$ ,  $i = 1, \dots, m$ ,*

$$E(a_1 \pm \Delta a_1, \dots, a_m \pm \Delta a_m) \subseteq E(a_1, \dots, a_m) \pm E^*(\Delta a_i). \quad (28)$$

The expression  $E^*(\Delta a_i)$  in (28) is, of course, the result of substituting  $\Delta a_i$  for the formal variables in  $E^*$  and evaluating the resulting expression.

*Proof of Corollary 11.* We proceed by induction on the height of  $E$  considered as polynomial tree. Repeatedly applying (23), we can extend it to any finite sums. Similarly, repeatedly applying (24), we can extend it to any finite products. This gives the statement (28) when the height of  $E$  is 1. The inductive step is also provided by extensions of (23) and (24) to finite sums and products.  $\square$

*Proof of Theorem 1.* In the proofs of Theorems 1–4 we will use the slightly informal notation exemplified by (6)  $\pm$  (7) being the confidence interval with midpoint (6) and half-width (7).

We obtain the confidence interval (6)  $\pm$  (7) for (2) (involving  $2|\mathbf{Z}|$  probabilities) by combining the confidence intervals (19) for each of the  $2|\mathbf{Z}|$  constituent probabilities. To ensure the overall confidence level  $1 - \delta$ , we replace the  $\delta$  in (19) by  $\delta/(2|\mathbf{Z}|)$ . By Corollary 11 applied to (2), we then indeed obtain the overall  $(1 - \delta)$ -confidence interval with midpoint (6) and semi-width

$$\sum_z \left( \sqrt{\frac{\ln \frac{4|\mathbf{Z}|}{\delta}}{2N}} + \sqrt{\frac{\ln \frac{4|\mathbf{Z}|}{\delta}}{2\#\tilde{x}z}} \right), \quad (29)$$

i.e., (7). The two square roots in (29) are the half-widths of the confidence intervals for  $P(z)$  and  $P(y | \tilde{x}, z)$ , respectively, given by Proposition 8.  $\square$

To derive (8) for Figure 1 with binary variables, notice that in the binary case we only need confidence intervals for three constituent probabilities, since an interval estimate for  $P(Z = 0)$  gives one for  $P(Z = 1)$  and vice versa. This allows us to replace  $\delta$  by  $\delta/3$  rather than  $\delta/4$  in (19).

*Proof of Theorem 2.* The proof is similar to that of Theorem 1. We regard (4) as a multivariate polynomial with formal variables  $P(z | \tilde{x})$  (indexed by  $z \in \mathbf{Z}$ ),  $P(y | x, z)$  (indexed by  $(x, z) \in \mathbf{X} \times \mathbf{Z}$ ), and  $P(x)$  (indexed by  $x \in \mathbf{X}$ ). It is clear that (9) is the midpoint for the confidence interval given by Corollary 11, so we only needed to show that (10) is the resulting half-width.

The total number of distinct formal variables in (4) is  $K := |\mathbf{X}| |\mathbf{Z}| + |\mathbf{X}| + |\mathbf{Z}|$ :

- there are  $|\mathbf{Z}|$  of  $P(z | \tilde{x})$ ,
- there are  $|\mathbf{X}| |\mathbf{Z}|$  of  $P(y | x, z)$ ,
- and there are  $|\mathbf{X}|$  of  $P(x)$ .

To ensure that (19) are simultaneous confidence intervals for all of them at confidence level  $1 - \delta$ , we replace the  $\delta$  in (19) by  $\delta/K$ .

Corollary 11 applied to the polynomial expression (4) gives the half-width

$$\sum_z \left( \sqrt{\frac{\ln \frac{2K}{\delta}}{2\#\tilde{x}}} + \sum_x \left( \sqrt{\frac{\ln \frac{2K}{\delta}}{2\#xz}} + \sqrt{\frac{\ln \frac{2K}{\delta}}{2N}} \right) \right)$$

$$= |\mathbf{X}| |\mathbf{Z}| \sqrt{\frac{\ln \frac{2K}{\delta}}{2N}} + |\mathbf{Z}| \sqrt{\frac{\ln \frac{2K}{\delta}}{2\#\tilde{x}}} + \sum_{x,z} \sqrt{\frac{\ln \frac{2K}{\delta}}{2\#xz}}, \quad (30)$$

i.e., it gives us (10). The square roots in (30) are again coming from Proposition 8.  $\square$

**Remark 12.** As far as the number of the operations “+” and “ $\times$ ” is concerned, the expanded form of a polynomial is typically less efficient than what we get by applying multivariate Horner schemes (see, e.g., [2]); there are several other methods for optimizing the number of operations (see, e.g., [6]). Namely, applying a multivariate Horner scheme leaves the same number of additions and reduces the number of multiplications. Since the right-hand side of (4) is already in the Horner form obtained from the expanded form

$$P(y | \text{do}(\tilde{x})) = \sum_{x,z} P(z | \tilde{x}) P(y | x, z) P(x) \quad (31)$$

of (4) by starting from the formal variables  $P(z | \tilde{x})$ , using it leads to a tighter confidence interval than using the expanded form (which we will spell out at the end of this remark). The multivariate Horner scheme depends on the order in which we apply it to different variables, and

$$P(y | \text{do}(\tilde{x})) = \sum_x P(x) \sum_z P(y | x, z) P(z | \tilde{x}) \quad (32)$$

is what we obtain in place of (4) when we start from the formal variables  $P(x)$ . Using (32) will give a different half-width from (30), namely

$$\begin{aligned} & \sum_x \left( \sqrt{\frac{\ln \frac{2K}{\delta}}{2N}} + \sum_z \left( \sqrt{\frac{\ln \frac{2K}{\delta}}{2\#xz}} + \sqrt{\frac{\ln \frac{2K}{\delta}}{2\#\tilde{x}}} \right) \right) \\ &= |\mathbf{X}| \sqrt{\frac{\ln \frac{2K}{\delta}}{2N}} + |\mathbf{X}| |\mathbf{Z}| \sqrt{\frac{\ln \frac{2K}{\delta}}{2\#\tilde{x}}} + \sum_{x,z} \sqrt{\frac{\ln \frac{2K}{\delta}}{2\#xz}}. \quad (33) \end{aligned}$$

We cannot say *a priori* which is larger, (30) or (33). It is easy to check that (30) is less than (33) if and only if

$$\frac{\#\tilde{x}}{N} < \left( \frac{1 - 1/|\mathbf{X}|}{1 - 1/|\mathbf{Z}|} \right)^2.$$

Therefore, the half-width (30) looks likely to be better than (33) overall; in particular, (30) is less than (33) when  $\#\tilde{x}/N \leq 1/4$  or  $|\mathbf{X}| \geq |\mathbf{Z}|$ .

The expanded form (31) gives

$$\sum_{x,z} \left( \sqrt{\frac{\ln \frac{2K}{\delta}}{2N}} + \sqrt{\frac{\ln \frac{2K}{\delta}}{2\#\tilde{x}}} + \sqrt{\frac{\ln \frac{2K}{\delta}}{2\#xz}} \right)$$

$$= |\mathbf{X}| |\mathbf{Z}| \sqrt{\frac{\ln \frac{2K}{\delta}}{2N}} + |\mathbf{X}| |\mathbf{Z}| \sqrt{\frac{\ln \frac{2K}{\delta}}{2\#\tilde{x}}} + \sum_{x,z} \sqrt{\frac{\ln \frac{2K}{\delta}}{2\#xz}},$$

which is worse than both (30) and (33).

*Proof of Theorem 3.* We apply Proposition 8 to estimating  $P(z)$  and Proposition 9 to estimating  $P(y | \tilde{x}, z)$  in (2). For Proposition 8 to be applicable, the values of  $Z_n$  at different steps in Figure 2 should be IID, and this follows from no vertex in  $Z_n$  being a descendant of  $X_n$  (one of the two conditions in the definition of the back-door criterion). For Proposition 9 to be applicable, we need the values of  $Y_n$  at the steps where  $X_n = \tilde{x}$  and  $Z_n = z$  to be IID, and this follows from  $Y_n$  and the (possibly non-IID) variables at the previous steps being  $d$ -separated by  $X_n \cup Z_n$ . To check the  $d$ -separation [8, Sect. 1.2.3], notice that each path from a variable at the previous step to  $Y_n$  is blocked by  $X_n$  if it has an arrow emanating from  $X_n$  and is blocked by  $Z_n$  if it has an arrow entering  $X_n$  from a variable at step  $n$ .

Since the total number of distinct formal variables in (2) is  $2|\mathbf{Z}|$ , we replace the  $\delta$  in (19) and (20) by  $\delta/(2|\mathbf{Z}|)$ . Since (12) is obviously the midpoint of the confidence interval given by Corollary 11, we only check that (13) is its half-width.

Plugging the half-width of the confidence interval for  $P(z)$  given by Proposition 8 and the half-width of the confidence interval for  $P(y | \tilde{x}, z)$  given by Proposition 9 into the polynomial expression (2) with the multiplications replaced by additions, we obtain the overall half-width

$$\sum_z \left( \sqrt{\frac{\ln \frac{4|\mathbf{Z}|}{\delta}}{2N}} + \sqrt{\frac{2 \ln \text{lb}[\#\tilde{x}z] + \ln \frac{6.6|\mathbf{Z}|}{\delta}}{2[\#\tilde{x}z]}} \right)$$

i.e., (13). □

In the case of Figure 1 with binary variables, we obtain (14) if we again replace  $\delta$  by  $\delta/3$  rather than  $\delta/4$  (and round up 9.9 to 10).

*Proof of Theorem 4.* The proof of Theorem 4 is analogous, and the only difference is that we use the confidence sequence (20) for estimating  $P(z)$  as well. □

## B Proofs for Sect. 6

In this appendix we prove Theorem 6. We assume, without loss of generality, that  $\mathbf{X} = \mathbf{Y} = \{0, 1\}$ ,  $\tilde{x} = 1$ , and  $y = 1$ . The underlying probability measure  $P$  defines the probability distribution of  $X_n$  given the past  $X_i, Y_i, i = 1, \dots, n-1$ , and also the probabilities  $p_1 := P(Y_n = 1 | X_n = 1)$  and  $p_0 := P(Y_n = 1 | X_n = 0)$ .

Theorem 6 will be deduced from the following special case of Proposition 11 in [3].

**Proposition 13.** *Let  $(I_n)$  be a  $(1 - \delta)$ -confidence sequence for the Bernoulli model. Then, for infinitely many  $n$ , there exists  $(y_1, \dots, y_n) \in \{0, 1\}^n$  such that  $|I_n(y_1, \dots, y_n)| > f(n)$ .*

*Proof.* It suffices to prove the statement of the proposition for the Bernoulli submodel in which the probability of success  $p$  is restricted to  $p \in (1/4, 3/4)$ . We just need to check the conditions of [3, Proposition 11]. The key condition is that

$$\text{KL}(B_p, B_{p'}) = O((p - p')^2),$$

where KL stands for Kullback–Leibler distance, and it follows from Taylor’s formula:

$$\text{KL}(B_p, B_{p+x}) = -p \ln \left( 1 + \frac{x}{p} \right) - (1-p) \ln \left( 1 - \frac{x}{1-p} \right) = \frac{x^2}{2p(1-p)} + O(x^3),$$

where  $O$  is uniform in  $p$ . □

To prove Theorem 6, we argue indirectly. Suppose such a family  $(C_N)$  of confidence estimators exists. First we turn it into a family of finite confidence sequences for the Bernoulli model. Define the “padded intervals”

$$I_n^N(y_1, \dots, y_n) := C_N(1, y_1, \dots, 1, y_n, 0, 0, \dots, 0, 0), \quad n = 0, \dots, N$$

(we apply  $C_N$  to the  $n$  pairs  $(1, y_i)$ ,  $i = 1, \dots, n$ , followed by  $N - n$  pairs  $(0, 0)$ ).

**Lemma 14.** *For each  $N$ , the padded intervals form a confidence sequence of length  $N$  for the Bernoulli model: for each  $p \in (0, 1)$ ,  $\mathbb{P}(p \in \cap_n I_n^N(\xi_1, \dots, \xi_n)) \geq 1 - \delta$ , where  $\xi_1, \dots, \xi_N$  are IID and Bernoulli with probability of success  $p$ .*

*Proof.* Given  $p$ , define the underlying probability measure as follows: let  $p_1 := p$  and  $p_0 := 0$ ; set  $X_n := 1$ ,  $n = 1, 2, \dots$ , and generate  $Y_n$  (setting  $Y_n := 1$  with probability  $p$ ) until

$$p \notin I_n^N(Y_1, \dots, Y_n). \tag{34}$$

As soon as (34) happens, start setting  $X_n := 0$  (and, therefore,  $Y_n := 0$  a.s.). It remains to notice that (34) means that (17) is violated. □

Suppose (18) is violated for all  $(x_1, y_1, \dots, x_N, y_N)$  for arbitrarily large  $N$ ; fix a sequence  $N_j$ ,  $j = 1, 2, \dots$ , of such successful horizons  $N$  (they are successful for the confidence estimator: it achieves accuracy  $f$ ). Therefore, we always have

$$|C_{N_j}(x_1, y_1, \dots, x_{N_j}, y_{N_j})| \leq f(x_1 + \dots + x_{N_j}).$$

In particular, we have

$$|I_n^{N_j}(y_1, \dots, y_n)| \leq f(n) \tag{35}$$

for all padded intervals.

Now let us pass to a limit along subsequences. The set of all binary strings  $y_1, \dots, y_n$  is countable, and the closed subintervals of  $[0, 1]$  form a compact set in the natural topology (where convergence is defined as the convergence of

both end-points). By a diagonal compactness argument, there is a subsequence, which we still denote  $N_j$ , such that for every  $n$  and every  $y_1, \dots, y_n$ ,

$$I_n^{N_j}(y_1, \dots, y_n) \rightarrow I_n(y_1, \dots, y_n)$$

for some limit interval  $I_n(y_1, \dots, y_n)$ . Namely, we find a subsequence of  $N_j$ , denoted  $N'_j$ , such that  $I_n^{N'_j}(y_1, \dots, y_n)$  converges for the first  $y_1, \dots, y_n$ , then we find a subsequence  $N''_j$  of  $N'_j$  such that  $I_n^{N''_j}(y_1, \dots, y_n)$  converges for the second  $y_1, \dots, y_n$ , etc.; the resulting subsequence  $\bar{N}_j$  is formed by taking the first element of  $N'_j$ , the second element of  $N''_j$ , etc.

Since  $I_n^{N_j}$  form finite confidence sequences,  $I_n$  will form a confidence sequence for the Bernoulli model. By (35), we always have  $|I_n| \leq f(n)$ . This contradicts Proposition 13.