

# Exact Anytime-valid Confidence Intervals for Contingency Tables and Beyond

Rosanne J. Turner<sup>a,b,\*</sup>, Peter D. Grünwald<sup>a,c</sup>

<sup>a</sup>*CWI, Amsterdam, part of NWO-I, Netherlands*

<sup>b</sup>*University Medical Center Utrecht, Brain Center, Netherlands*

<sup>c</sup>*Leiden University, Department of Mathematics, Netherlands*

---

## Abstract

E-variables are tools for retaining type-I error guarantee with optional stopping. We extend E-variables for sequential two-sample tests to general null hypotheses and anytime-valid confidence sequences. We provide implementations for estimating risk difference, relative risk and odds-ratios in contingency tables.

*Keywords:* confidence sequences, contingency tables, anytime-valid, effect size

---

## 1. Introduction

We consider a setting where we collect samples from two distinct groups, denoted  $a$  and  $b$ . In both groups, data come in sequentially and are i.i.d. We thus have two data streams,  $Y_{1,a}, Y_{2,a}, \dots$  i.i.d.  $\sim P_{\theta_a}$  and  $Y_{1,b}, Y_{2,b}, \dots$  i.i.d.  $\sim P_{\theta_b}$  where we assume that  $\theta_a, \theta_b \in \Theta$ ,  $\{P_\theta : \theta \in \Theta\}$  representing some parameterized underlying family of distributions, all assumed to have a probability density or mass function denoted by  $p_\theta$  on some outcome space  $\mathcal{Y}$ .

E-variables (Grünwald et al., 2022; Vovk and Wang, 2021) are a tool for constructing tests that keep their Type-I error control under optional stopping and continuation. Previously, Turner et al. (2021) developed E-variables for testing equality of both data streams, i.e. with null hypothesis  $\vec{\Theta}_0 := \{(\theta_a, \theta_b) \in \Theta^2 : \theta_a = \theta_b\}$ . Here we first generalize these E-variables to more general null hypotheses in which we may have  $\theta_a \neq \theta_b$ . We then use these generalized E-variables to construct *anytime-valid* confidence sequences; these provide confidence sets that remain valid under optional stopping (Darling and Robbins, 1967; Howard et al., 2021).

As in (Turner et al., 2021), we first design E-variables for a *single block* of data  $(Y_a^{n_a}, Y_b^{n_b})$ , where a block is a set of data consisting of  $n_a$  outcomes  $Y_a^{n_a} = (Y_{a,1}, \dots, Y_{a,n_a})$  in group  $a$  and  $n_b$  outcomes  $Y_b^{n_b} = (Y_{b,1}, \dots, Y_{b,n_b})$  in group  $b$ , for some pre-specified  $n_a$  and  $n_b$ . An *E-variable* is then, by definition,

---

\*Corresponding author, Rosanne.Turner@cwi.nl, National Research Institute for mathematics and computer science in the Netherlands (CWI), Science Park 123, 1098 XG Amsterdam, The Netherlands. Declarations of interest: none.

any nonnegative random variable  $S = s'(Y_a^{n_a}, Y_b^{n_b})$  such that

$$\sup_{(\theta_a, \theta_b) \in \vec{\Theta}_0} \mathbf{E}_{Y_a^{n_a} \sim P_{\theta_a}, Y_b^{n_b} \sim P_{\theta_b}} [s'(Y_a^{n_a}, Y_b^{n_b})] \leq 1. \quad (1)$$

Turner et al. (2021) first defined such an E-variable for  $\vec{\Theta}_0 = \{(\theta_a, \theta_b) \in \Theta^2 : \theta_a = \theta_b\}$  so that it would tend to have high power against a given simple alternative  $\vec{\Theta}_1 = \{(\theta_a^*, \theta_b^*)\}$ . Their E-variable is of the following simple form (with  $n = n_a + n_b$ ):

$$s'(Y_a^{n_a}, Y_b^{n_b}) = \frac{p_{\theta_a^*}(Y_a^{n_a})}{\prod_{i=1}^{n_a} (\frac{n_a}{n} p_{\theta_a^*}(Y_{a,i}) + \frac{n_b}{n} p_{\theta_b^*}(Y_{a,i}))} \cdot \frac{p_{\theta_b^*}(Y_b^{n_b})}{\prod_{i=1}^{n_b} (\frac{n_a}{n} p_{\theta_a^*}(Y_{b,i}) + \frac{n_b}{n} p_{\theta_b^*}(Y_{b,i}))}. \quad (2)$$

These E-variables can be extended to sequences of blocks  $Y_{(1)}, Y_{(2)}, \dots$  by multiplication, and can be extended to composite alternatives by sequentially learning  $(\theta_a^*, \theta_b^*)$  from the data, for example via a Bayesian prior on  $\vec{\Theta}_1$ . The  $n_a$  and  $n_b$  used for the  $j$ -th block  $Y_{(j)}$  are allowed to depend on past data, but they must be fixed before the first observation in block  $j$  occurs. For simplicity, in this note we only consider the case with  $n_a$  and  $n_b$  that remain fixed throughout; extension to the general case is straightforward.

By a general property of E-variables, at each point in time, the running product of block E-variables observed so far is itself an E-variable, and the random process of the products is known as a *test martingale* (Grünwald et al., 2022; Shafer, 2021). An E-variable-based test at level  $\alpha$  is a test which, in combination with any stopping rule  $\tau$ , reports ‘reject’ if and only if the product of E-values corresponding to all blocks that were observed at the stopping time and have already been completed, is larger than  $1/\alpha$ . Such a test has a type-I error probability bounded by  $\alpha$  irrespective of the stopping time  $\tau$  that was used; see the aforementioned references for much more detailed introductions and, for example (Henzi and Ziegel, 2021), for a practical application.

In case  $\{P_\theta : \theta \in \Theta\}$  is convex, the E-variable (2) has the so-called GRO- (*growth-rate-optimality*) property: it maximizes, over all E-variables (i.e. over all nonnegative random variables  $S = s'(Y_a^{n_a}, Y_b^{n_b})$  satisfying (1)) the logarithmic growth rate

$$\mathbf{E}_{Y_a^{n_a} \sim P_{\theta_a^*}, Y_b^{n_b} \sim P_{\theta_b^*}} [\log S], \quad (3)$$

which implies that, under  $(\theta_a^*, \theta_b^*)$ , the expected number of data points before the null can be rejected is minimized (Grünwald et al., 2022).

Below, in Theorem 1 in section 2, which generalizes Theorem 1 in Turner et al. (2021), we extend (2) to the case of general null hypotheses,  $\vec{\Theta}_0 \subset \Theta^2$ , allowing for the case that the elements of  $\vec{\Theta}_0$  have two

different components, and provide a condition under which it has the GRO property. From then onwards we focus on what we call ‘the  $2 \times 2$  contingency table setting’ in which both streams are Bernoulli,  $\theta_j$  denoting the probability of 1 in group  $j$ . For this case, Theorem 2 gives a simplified expression for the E-variable and shows that the GRO property holds if  $\vec{\Theta}_0 \subset [0, 1]^2$  is convex. Then we will extend this E-variable to deal with composite  $\vec{\Theta}_1$  and use this to define anytime-valid confidence sequences. We illustrate these through simulations. All proofs are in Appendix A.

## 2. General Null Hypotheses

In this section, we first construct an E-variable for general null hypotheses that generalizes (2). We then instantiate the new result to the  $2 \times 2$  case. The following development and results require  $\{P_\theta : \theta \in \Theta\}$  to be ‘nondegenerate’ in the sense that there exists  $\theta \in \Theta$  such that for all  $\theta' \in \Theta$ ,  $D(P_\theta \| P_{\theta'}) < \infty$ . This mild condition holds, for example, for exponential families; we tacitly assume nondegeneracy from now on.

Our goal is thus to define an E-variable for a block of  $n = n_a + n_b$  data points with  $n_g$  points in group  $g$ ,  $g \in \{a, b\}$ . For notational convenience we define, for  $\theta_a, \theta_b \in \Theta$ ,  $P_{\theta_a, \theta_b}$  as the joint distribution of  $Y_a^{n_a} \sim P_{\theta_a}$  and  $Y_b^{n_b} \sim P_{\theta_b}$ , so that  $p_{\theta_a, \theta_b}(y_a^{n_a}, y_b^{n_b}) = \prod_{i=1}^{n_a} p_{\theta_a}(y_{a,i}) \prod_{i=1}^{n_b} p_{\theta_b}(y_{b,i})$  so that we can write the null hypothesis as  $\mathcal{H}_0 := \{P_{\theta_a, \theta_b} : (\theta_a, \theta_b) \in \vec{\Theta}_0\}$ . Our strategy will be to first develop an E-variable for a *modified* setting in which there is only a single outcome, falling with probability  $n_a/n$  in group  $a$  and  $n_b/n$  in group  $b$ . To this end, for  $\vec{\theta} = (\theta_a, \theta_b)$ , we define  $p'_{\vec{\theta}}(Y|a) := p_{\theta_a}(y)$ ,  $p'_{\vec{\theta}}(Y|b) := p_{\theta_b}(y)$ , all distributions with a ‘ referring to the modified setting with just one outcome. We let  $\mathcal{W}^\circ(\vec{\Theta}_0)$  be the set of all distributions on  $\vec{\Theta}_0$  with finite support. For  $W \in \mathcal{W}^\circ(\vec{\Theta}_0)$ , we define  $p'_W(Y|g) = \int p'_{\vec{\theta}}(Y|g) dW(\vec{\theta})$ . We set  $p'_W(y^k|g) := \prod_{i=1}^k p'_W(y_i|g)$ . We further define, for given alternative  $\vec{\Theta}_1 = \{(\theta_a^*, \theta_b^*)\}$ ,  $p^\circ(\cdot|g)$ ,  $g \in \{a, b\}$  to be, if it exists, the conditional probability density satisfying

$$\mathbf{E}_{G \sim Q'} \mathbf{E}_{Y \sim P_{\theta_g^*}} [-\log p^\circ(Y|G)] = \inf_{W \in \mathcal{W}^\circ(\vec{\Theta}_0)} \mathbf{E}_{G \sim Q'} \mathbf{E}_{Y \sim P_{\theta_g^*}} [-\log p'_W(Y|G)] \quad (4)$$

with  $Q'(G)$  the distribution for  $G \in \{a, b\}$  with  $Q'(G = a) = n_a/n$ . Clearly we can rephrase (4) equivalently as:

$$D(Q'(G, Y) \| P^\circ(G, Y)) = \inf_{W \in \mathcal{W}^\circ(\vec{\Theta}_0)} D(Q'(G, Y) \| P'_W(G, Y)), \quad (5)$$

where  $D$  is the KL divergence. Here we extended the conditional distributions  $P'_W(Y|G)$  and  $P^\circ(Y|G)$  (corresponding to densities  $p'_W(Y|G)$  and  $p^\circ(Y|G)$ ) to a joint distribution by setting  $P'_W(G, Y) := Q'(G)P'_W(Y|G)$  (and similarly for  $P^\circ$ ) and we extended  $Q'(G, Y) := Q'(G)P_{\theta_g^*}(Y)$ . We have now constructed a modified

null hypothesis  $\mathcal{H}'_0 = \{P'_{\vec{\theta}}(G, Y) : \vec{\theta} \in \vec{\Theta}_0\}$  of joint distributions for a single ‘group’ outcome  $G \in \{a, b\}$  and ‘data’ outcome  $Y \in \mathcal{Y}$ . We let  $\bar{\mathcal{H}}'_0 = \{P_W(G, Y) : W \in \mathcal{W}^\circ(\vec{\Theta}_0)\}$  be the convex hull of  $\mathcal{H}'_0$ .

The  $p^\circ$  satisfying (5) is commonly called the *reverse information projection* of  $Q'$  onto  $\bar{\mathcal{H}}'_0$ . Li (1999) shows that  $p^\circ$  always exists under our nondegeneracy condition, though in some cases it may represent a sub-distribution (integrating to strictly less than one); see (Grünwald et al., 2022, Theorem 1) (re-stated for convenience in the supplementary material) who, building on Li’s work, established a general relation between reverse information projection and E-variables. Part 1 of that theorem establishes that if the minimum in (4) (or (5)) is achieved by some  $W^\circ \in \mathcal{W}^\circ$  then  $p^\circ(\cdot|\cdot) = p'_{W^\circ}(\cdot|\cdot)$  and, with  $\vec{\theta}^* = (\theta_a, \theta_b)$ , for all  $\vec{\theta} \in \vec{\Theta}_0$ ,

$$\mathbf{E}_{G \sim Q'} \mathbf{E}_{Y \sim P'_{\vec{\theta}}|G} \left[ \frac{p'_{\vec{\theta}^*}(Y|G)}{p^\circ(Y|G)} \right] = \mathbf{E}_{G \sim Q'} \mathbf{E}_{Y \sim P'_{\vec{\theta}^*}|G} \left[ \frac{p'_{\vec{\theta}^*}(G, Y)}{p^\circ(G, Y)} \right] \leq 1. \quad (6)$$

This expresses that  $p'_{\vec{\theta}^*}(Y|G)/p^\circ(Y|G)$  is an E-variable for our modified problem, in which within a single block we observe a single outcome in group  $g$ , with  $g$  chosen with probability  $n_g/n$ . If we were to interpret the E-variable of the modified problem as in (6) as a likelihood ratio for a single outcome, its corresponding likelihood ratio for a single block of data in our original problem with  $n_g$  outcomes in group  $g$  would be:

$$s(y_a^{n_a}, y_b^{n_b}; n_a, n_b, (\theta_a^*, \theta_b^*); \vec{\Theta}_0) := \frac{p'_{(\theta_a^*, \theta_b^*)}(y_a^{n_a}|a)p'_{(\theta_a^*, \theta_b^*)}(y_b^{n_b}|b)}{p^\circ(y_a^{n_a}|a)p^\circ(y_b^{n_b}|b)} = \frac{p_{\theta_a^*}(y_a^{n_a})p_{\theta_b^*}(y_b^{n_b})}{p^\circ(y_a^{n_a}|a)p^\circ(y_b^{n_b}|b)}. \quad (7)$$

The following theorem expresses that this ‘extension’ of the E-variable in the modified problem gives us an E-variable in our original problem:

**Theorem 1.**  $S_{[n_a, n_b, \theta_a^*, \theta_b^*; \vec{\Theta}_0]} := s(Y_a^{n_a}, Y_b^{n_b}; n_a, n_b, (\theta_a^*, \theta_b^*); \vec{\Theta}_0)$  as in (7) is an E-variable, i.e. with  $s'(\cdot) = s(\cdot; n_a, n_b, (\theta_a^*, \theta_b^*); \vec{\Theta}_0)$ , we have (1). Moreover, if  $\mathcal{H}'_0 = \{P'_{\vec{\theta}} : \vec{\theta} \in \vec{\Theta}_0\}$  (the null hypothesis for the modified problem) is a convex set of distributions and  $\mathcal{Y}$  is finite (so that  $\mathcal{H}'_0 = \bar{\mathcal{H}}'_0$ ) and furthermore  $\mathcal{H}'_0$  is compact in the weak topology, then (a)  $p^\circ(\cdot|\cdot) = p'_{\vec{\theta}}(\cdot|\cdot)$  for some  $\vec{\theta} \in \vec{\Theta}_0$  and (b)  $S_{[n_a, n_b, \theta_a^*, \theta_b^*; \vec{\Theta}_0]}$  is the  $(\theta_a^*, \theta_b^*)$ -GRO E-variable for the original problem, maximizing (3) among all E-variables.

In the case that  $\mathcal{H}'_0$  is not convex and compact, we do not have a simple expression for  $p^\circ$  in general, and we may have to find it numerically by minimizing (4). In the  $2 \times 2$  table (Bernoulli  $\Theta$ ) case though, there are interesting  $\mathcal{H}_0$  for which the corresponding  $\mathcal{H}'_0$  is convex, and we shall now see that this leads to major simplifications.

2.1. General Convex  $\vec{\Theta}_0$  for the  $2 \times 2$  contingency table

In this subsection and the next,  $\{P_{\theta_a, \theta_b}\}$  refers to the  $2 \times 2$  model again, with  $\mathcal{Y} = \{0, 1\}$  and  $\theta$  denoting the probability of 1. We now let  $\vec{\Theta}_0$  be any closed convex subset of  $[0, 1]^2$  that contains a point in the interior of  $[0, 1]^2$ . Again, note that the corresponding  $\mathcal{H}_0 = \{P_{\vec{\theta}} : \vec{\theta} \in \vec{\Theta}_0\}$  need not be convex; still,  $\mathcal{H}'_0$ , the null hypothesis for the modified problem as defined above, must be convex if  $\vec{\Theta}_0$  is convex, and this will allow us to design E-variables for such  $\vec{\Theta}_0$ . Let  $\mathcal{H}_1 = \{P_{\theta_a^*, \theta_b^*}\}$  with  $(\theta_a^*, \theta_b^*)$  in the interior of  $[0, 1]^2$ , and let

$$\text{KL}(\theta_a, \theta_b) := D(P_{\theta_a^*, \theta_b^*}(Y_a^{n_a}, Y_b^{n_b}) \| P_{\theta_a, \theta_b}(Y_a^{n_a}, Y_b^{n_b})) = \sum_{y_a^{n_a} \in \{0, 1\}^{n_a}, y_b^{n_b} \in \{0, 1\}^{n_b}} p_{\theta_a^*}(y_a^{n_a}) p_{\theta_b^*}(y_b^{n_b}) \log \frac{p_{\theta_a^*}(y_a^{n_a}) p_{\theta_b^*}(y_b^{n_b})}{p_{\theta_a}(y_a^{n_a}) p_{\theta_b}(y_b^{n_b})} \quad (8)$$

stand for the KL divergence between  $P_{\theta_a^*, \theta_b^*}$  and  $P_{\theta_a, \theta_b}$  restricted to a single block (note that in the previous subsection, KL divergence was defined for a single outcome  $Y$ ). The following result builds on Theorem 1:

**Theorem 2.**  $\min_{(\theta_a, \theta_b) \in \vec{\Theta}_0} \text{KL}(\theta_a, \theta_b)$  is uniquely achieved by some  $(\theta_a^\circ, \theta_b^\circ)$ . If  $(\theta_a^*, \theta_b^*) \in \vec{\Theta}_0$ , then  $(\theta_a^\circ, \theta_b^\circ) = (\theta_a^*, \theta_b^*)$ . Otherwise,  $(\theta_a^\circ, \theta_b^\circ)$  lies on the boundary of  $\vec{\Theta}_0$ , but not on the boundary of  $[0, 1]^2$ . The E-variable (7) is given by the distribution  $W$  that puts all its mass on  $(\theta_a^\circ, \theta_b^\circ)$ , i.e.

$$s(y_a^{n_a}, y_b^{n_b}; n_a, n_b, (\theta_a^*, \theta_b^*); \vec{\Theta}_0) = \frac{p_{\theta_a^*}(y_a^{n_a}) p_{\theta_b^*}(y_b^{n_b})}{p_{\theta_a^\circ}(y_a^{n_a}) p_{\theta_b^\circ}(y_b^{n_b})} \quad (9)$$

is an E-variable. Moreover, this is the  $(\theta_a^*, \theta_b^*)$ -GRO E-variable relative to  $\vec{\Theta}_0$ .

We can extend this E-variable to the case of a composite  $\mathcal{H}_1 = \{P_{\theta_a, \theta_b} : (\theta_a, \theta_b) \in \vec{\theta}_1\}$  by *learning* the true  $(\theta_a^*, \theta_b^*) \in \vec{\theta}_1$  from the data (Turner et al., 2021). We thus replace, for each  $j = 1, 2, \dots$ , for the block  $Y^{(j)}$  consisting of  $n_a$  points  $Y_{(j), a, 1}, \dots, Y_{(j), a, n_a}$  in group  $a$  and  $n_b$  points  $Y_{(j), b, 1}, \dots, Y_{(j), b, n_b}$  in group  $b$ , the ‘true’  $\theta_g^*$  for  $g \in \{a, b\}$  by an estimate  $\check{\theta}_g | Y^{(j-1)}$  based on the previous  $j - 1$  data blocks. The E-variable corresponding to  $m$  blocks of data then becomes

$$S_{[n_a, n_b, W_1; \vec{\Theta}_0]}^{(m)} = \prod_{j=1}^m \prod_{i=1}^{n_a} \frac{p_{\check{\theta}_a | Y^{(j-1)}}(Y_{(j), a, i})}{p_{\check{\theta}_a^\circ | Y^{(j-1)}}(Y_{(j), a, i})} \prod_{i=1}^{n_b} \frac{p_{\check{\theta}_b | Y^{(j-1)}}(Y_{(j), b, i})}{p_{\check{\theta}_b^\circ | Y^{(j-1)}}(Y_{(j), b, i})} \quad (10)$$

where, for  $g \in \{a, b\}$ ,  $\check{\theta}_g | Y^{(j-1)}$  can be an arbitrary estimator (function from  $Y^{(j-1)}$  to  $\theta_g$ ) and  $(\check{\theta}_a^\circ | Y^{(j-1)}, \check{\theta}_b^\circ | Y^{(j-1)})$  is defined to achieve  $\min_{(\theta_a, \theta_b) \in \vec{\Theta}_0} D(P_{\check{\theta}_a | Y^{(j-1)}, \check{\theta}_b | Y^{(j-1)}}(Y_a^{n_a}, Y_b^{n_b}) \| P_{\theta_a, \theta_b}(Y_a^{n_a}, Y_b^{n_b}))$ . No matter what estimator we choose, (10) gives us an E-variable. In Section 3, as in (Turner et al.,

2021), we implement this estimator by fixing a prior  $W$  and using the Bayes posterior mean,  $\check{\theta}_g|Y^{(j-1)} := \mathbf{E}_{\theta_g \sim W|Y^{(j-1)}}[\theta_g]$ . Let us now illustrate Theorem 2 for two choices of  $\vec{\Theta}_0$ .

$\vec{\Theta}_0$  with linear boundary. First, we let  $\vec{\Theta}_0(s, c)$ , for  $s \in \mathbf{R}, c \in \mathbf{R}$ , stand for any straight line through  $[0, 1]^2$ :  $\vec{\Theta}_0(s, c) := \{(\theta_a, \theta_b) \in [0, 1]^2 : \theta_b = s + c\theta_a\}$ . This can be extended to  $\vec{\Theta}_0(\leq s, c) := \bigcup_{s' \leq s} \vec{\Theta}_0(s', c)$  and similarly to  $\vec{\Theta}_0(\geq s, c) := \bigcup_{s' \geq s} \vec{\Theta}_0(s', c)$ . For example, we could take  $\vec{\Theta}_0 = \vec{\Theta}_0(s, c)$  to be the solid line in Figure 1(a) (which would correspond to  $s = 0.1, c = 1$ ), or the whole area underneath the line ( $\vec{\Theta}_0(\leq s, c)$ ) including the line itself, or the whole area above it including the line itself ( $\vec{\Theta}_0(\geq s, c)$ ).

Now consider a  $\vec{\Theta}_0(s, c)$  that has nonempty intersection with the interior of  $[0, 1]^2$  and that is separated from the point alternative  $(\theta_a^*, \theta_b^*)$ , i.e.  $\min_{(\theta_a, \theta_b) \in \vec{\Theta}_0} \text{KL}(\theta_a, \theta_b) > 0$ . Simple differentiation gives that the minimum is achieved by the unique  $(\theta_a^\circ, \theta_b^\circ) \in \vec{\Theta}_0$  satisfying:

$$n_a \left( -\frac{\theta_a^*}{\theta_a^\circ} + \frac{1 - \theta_a^*}{1 - \theta_a^\circ} \right) + n_b \cdot c \cdot \left( -\frac{\theta_b^*}{\theta_b^\circ} + \frac{1 - \theta_b^*}{1 - \theta_b^\circ} \right) = 0, \quad (11)$$

which can now be plugged into the E-variable (9) if the alternative is the simple alternative, or otherwise into its sequential form (10). In the basic case in which  $\vec{\Theta}_0 = \{(\theta_a, \theta_b) \in [0, 1]^2 : \theta_a = \theta_b\}$ , the solution to (11) reduces to the familiar  $\theta_a^\circ = \theta_b^\circ = (n_a\theta_a^* + n_b\theta_b^*)/n$  from Turner et al. (2021).

If  $(\theta_a^*, \theta_b^*)$  lies above the line  $\vec{\Theta}_0(s, c)$ , then by Theorem 2,  $\min_{(\theta_a, \theta_b) \in \vec{\Theta}_0(\leq s, c)} \text{KL}(\theta_a, \theta_b)$  must lie on  $\vec{\Theta}_0(s, c)$ . Theorem 2 gives that it must be achieved by the  $(\theta_a^\circ, \theta_b^\circ)$  satisfying (11). Similarly, if  $(\theta_a^*, \theta_b^*)$  lies below the line  $\vec{\Theta}_0(s, c)$ , then  $\min_{(\theta_a, \theta_b) \in \vec{\Theta}_0(\geq s, c)} \text{KL}(\theta_a, \theta_b)$  is again achieved by the  $(\theta_a^\circ, \theta_b^\circ)$  satisfying (11).

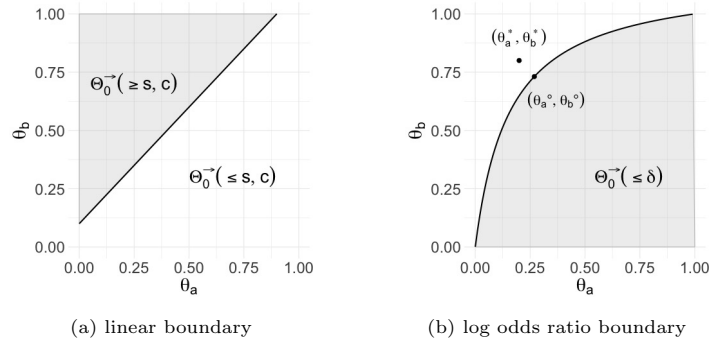


Figure 1: Examples of null hypothesis parameter spaces for two types of boundaries.

$\vec{\Theta}_0$  with log odds ratio boundary. Similarly, we can consider  $\vec{\Theta}_0(\delta)$ ,  $\vec{\Theta}_0(\leq \delta)$ ,  $\vec{\Theta}_0(\geq \delta)$  that correspond to a given log odds effect size  $\delta$ . That is, we now take

$$\begin{aligned}\vec{\Theta}_0(\delta) &:= \left\{ (\theta_a, \theta_b) \in [0, 1]^2 : \log \frac{\theta_b(1-\theta_a)}{(1-\theta_b)\theta_a} = \delta \right\} \\ \vec{\Theta}_0(\leq \delta) &:= \left\{ (\theta_a, \theta_b) \in [0, 1]^2 : \log \frac{\theta_b(1-\theta_a)}{(1-\theta_b)\theta_a} \leq \delta \right\} \\ \vec{\Theta}_0(\geq \delta) &:= \left\{ (\theta_a, \theta_b) \in [0, 1]^2 : \log \frac{\theta_b(1-\theta_a)}{(1-\theta_b)\theta_a} \geq \delta \right\}.\end{aligned}$$

For example, we could now take  $\vec{\Theta}_0 = \vec{\Theta}_0(\leq \delta)$  to be the area under the curve (including the curve boundary itself) in Figure 1(b), which would correspond to  $\delta = 2$ . Now let  $\delta$  and point alternative  $(\theta_a^*, \theta_b^*)$  be such that  $\delta > 0$  and  $\vec{\Theta}_0(\leq \delta)$  is separated from  $(\theta_a^*, \theta_b^*)$ , i.e.  $\min_{(\theta_a, \theta_b) \in \vec{\Theta}_0(\leq \delta)} \text{KL}(\theta_a, \theta_b) > 0$ . Let  $(\theta_a^\circ, \theta_b^\circ) := \arg \min_{(\theta_a, \theta_b) \in \vec{\Theta}_0(\leq \delta)} \text{KL}(\theta_a, \theta_b)$ . As Figure 1(b) suggests,  $\vec{\Theta}_0(\leq \delta)$  is convex. Theorem 2 now tells us that  $\min_{(\theta_a, \theta_b) \in \vec{\Theta}_0(\leq \delta)} \text{KL}(\theta_a, \theta_b)$  is achieved by  $(\theta_a^\circ, \theta_b^\circ)$ . Plugging these into (9) thus gives us an E-variable.  $(\theta_a^\circ, \theta_b^\circ)$  can easily be determined numerically. Similarly, if  $\delta < 0$ ,  $\vec{\Theta}_0(\geq \delta)$  is convex and closed and if  $(\theta_a^*, \theta_b^*)$  is separated from  $\vec{\Theta}_0(\geq \delta)$ , the  $(\theta_a^\circ, \theta_b^\circ)$  minimizing KL on  $\vec{\Theta}_0(\delta)$  gives an E-variable relative to  $\vec{\Theta}_0(\geq \delta)$ .

### 3. Anytime-Valid Confidence for the $2 \times 2$ case

We will now use the E-variables defined above to construct anytime-valid confidence sequences. Let  $\delta = \delta(\theta_a, \theta_b)$  be a notion of effect size such as the log odds ratio (see above) or absolute risk  $\theta_b - \theta_a$  or relative risk  $\theta_b/\theta_a$ . A  $(1 - \alpha)$ -anytime-valid (AV) confidence sequence (Darling and Robbins, 1967; Howard et al., 2021) is a sequence of random (i.e. determined by data) subsets  $\text{CS}_{\alpha,(1)}, \text{CS}_{\alpha,(2)}, \dots$  of  $\Gamma$ , with  $\text{CS}_{\alpha,(m)}$  being a function of the first  $m$  data blocks  $Y^{(m)}$ , such that for all  $(\theta_a, \theta_b) \in [0, 1]^2$ ,

$$P_{\theta_a, \theta_b} (\exists m \in \mathbf{N} : \delta(\theta_a, \theta_b) \notin \text{CS}_{\alpha,(m)}) \leq \alpha.$$

We first consider the case in which for all values  $\gamma \in \Gamma$  that  $\delta$  can take,  $\vec{\Theta}_0(\gamma) := \{(\theta_a, \theta_b) \in [0, 1]^2 : \delta(\theta_a, \theta_b) = \gamma\}$  is a convex set, as it will be for absolute and relative risk. Fix a prior  $W_1$  on  $[0, 1]^2$ . Based on (10) we can make an *exact* (nonasymptotic) AV confidence sequence

$$\text{CS}_{\alpha,(m)} = \left\{ \delta : S_{[n_a, n_b, W_1; \vec{\Theta}_0(\delta)]}^{(m)} \leq \frac{1}{\alpha} \right\} \quad (12)$$

where  $S_{[n_a, n_b, W_1; \vec{\Theta}_0(\delta)]}^{(m)}$  is defined as in (10) and is a valid E-variable by Theorem 2. To see that  $(\text{CS}_{\alpha, (m)})_{m \in \mathbf{N}}$  really is an AV confidence sequence, note that, by definition of the  $\text{CS}_{\alpha, (m)}$ , we have

$P_{\theta_a, \theta_b}(\exists m \in \mathbf{N} : \delta(\theta_a, \theta_b) \notin \text{CS}_{\alpha, (m)})$  is given by

$$P_{\theta_a, \theta_b} \left( \exists m \in \mathbf{N} : S_{[n_a, n_b, W_1; \vec{\Theta}_0(\delta)]}^{(m)} \geq \frac{1}{\alpha} \right) \leq \alpha,$$

by Ville’s inequality (Grünwald et al., 2022; Turner et al., 2021). Here the  $\text{CS}_{\alpha, (m)}$  are not necessarily intervals, but, potentially losing some information, we can make a AV confidence sequence consisting of intervals by defining  $\text{CI}_{\alpha, (m)}$  to be the smallest interval containing  $\text{CS}_{\alpha, (m)}$ . We can also turn any confidence sequences  $(\text{CS}_{\alpha, (m)})_{m \in \mathbf{N}}$  into an alternative AV confidence sequence with sets  $\text{CS}'_{\alpha, (m)}$  that are always a subset of  $\text{CS}_{\alpha, (m)}$  by taking the *running intersection*

$$\text{CS}'_{\alpha, (m)} := \bigcap_{j=1..m} \text{CS}_{\alpha, (j)}.$$

In this form, the confidence sequences  $\text{CS}'_{\alpha, (m)}$  can be interpreted as *the set of  $\delta$ ’s that have not yet been rejected* in a setting in which, for each null hypothesis  $\vec{\Theta}_0(\delta)$  we stop and reject as soon as the corresponding E-variable exceeds  $1/\alpha$ . The running intersection can also be applied to the intervals  $(\text{CI}_{\alpha, (m)})_{m \in \mathbf{N}}$ .

To simplify calculations, it is useful to take  $W_1$  a prior under which  $\theta_a$  and  $\theta_b$  have independent beta distributions with parameters  $\alpha_a, \beta_a, \alpha_b, \beta_b$ . We can, if we want, infuse some prior knowledge or hopes by setting these parameters to certain values — our confidence sequences will be valid irrespective of our choice (Howard et al., 2021). In case no such knowledge can be formulated (as in the simulations below), we advocate the prior, which, among all priors of the simple form asymptotically achieves the REGROW criterion (a criterion related to minimax log-loss regret, see (Grünwald et al., 2022)), i.e for the case  $n_a = n_b = 1$  we set  $W_1$  to an independent beta prior on  $\theta_a$  and  $\theta_b$  with  $\gamma = 0.18$  as was empirically found to be the ‘best’ value (Turner et al., 2021).

*Log Odds Ratio Effect Size.* The situation is slightly trickier if we take the log odds ratio as effect size, for  $\vec{\Theta}_0(\delta)$  is then not convex. Without convexity, Theorem 2 cannot be used and hence the validity of AV confidence sequences as constructed above breaks down. We can get nonasymptotic anytime-valid confidence sequences after all as follows. First, we consider a one-sided AV confidence sequence for the submodel of



positive effect sizes  $\{(\theta_a, \theta_b) : \delta(\theta_a, \theta_b) \geq 0\}$ , defining

$$\text{CS}_{\alpha, (m)}^+ = \{\delta \geq 0 : S_{[n_a, n_b, W_1; \vec{\Theta}_0(\leq \delta)]}^{(m)} \leq \alpha^{-1}, \}$$

where we note that  $\vec{\Theta}_0(\leq \delta)$  is convex (since  $\delta \geq 0$ ) and also contains  $(\theta_a, \theta_b)$  with  $\delta(\theta_a, \theta_b) < 0$ . This confidence sequence can give a lower bound on  $\delta$ . Analogously, we consider a one-sided AV confidence sequence for the submodel  $\{(\theta_a, \theta_b) : \delta(\theta_a, \theta_b) \leq 0\}$ , defining

$$\text{CS}_{\alpha, (m)}^- = \{\delta \leq 0 : S_{[n_a, n_b, W_1; \vec{\Theta}_0(\geq \delta)]}^{(m)} \leq \alpha^{-1}\},$$

and derive an upper bound on  $\delta$ . By Theorem 2, both sequences  $(\text{CS}_{\alpha, (m)}^+)_{m=1,2,\dots}$  and  $(\text{CS}_{\alpha, (m)}^-)_{m=1,2,\dots}$  are AV confidence sequences for the submodels with  $\delta \geq 0$  and  $\delta \leq 0$  respectively. Defining  $\text{CS}_{\alpha, (m)} = \text{CS}_{\alpha, (m)}^+ \cup \text{CS}_{\alpha, (m)}^-$ , we find, for  $(\theta_a, \theta_b)$  with  $\delta(\theta_a, \theta_b) > 0$ ,

$$P_{\theta_a, \theta_b}(\exists m \in \mathbf{N} : \delta(\theta_a, \theta_b) \notin \text{CS}_{\alpha, (m)}) = P_{\theta_a, \theta_b}(\exists m \in \mathbf{N} : \delta(\theta_a, \theta_b) \notin \text{CS}_{\alpha, (m)}^+) \leq \alpha,$$

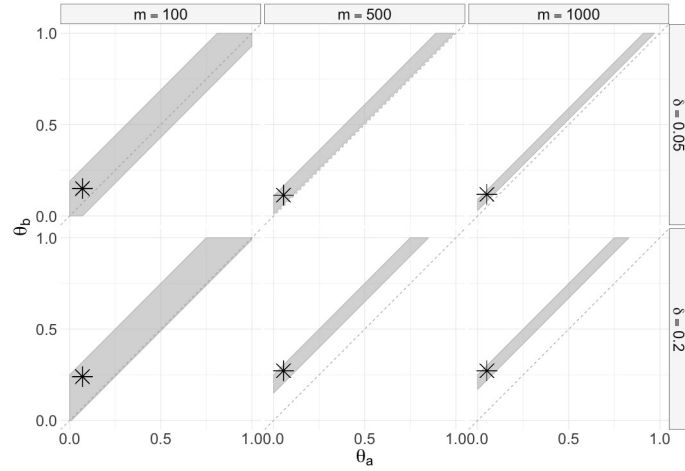
and analogously for  $(\theta_a, \theta_b)$  with  $\delta(\theta_a, \theta_b) < 0$ . We have thus arrived at a confidence sequence that works for all  $\delta$ , positive or negative.

### 3.1. Simulations

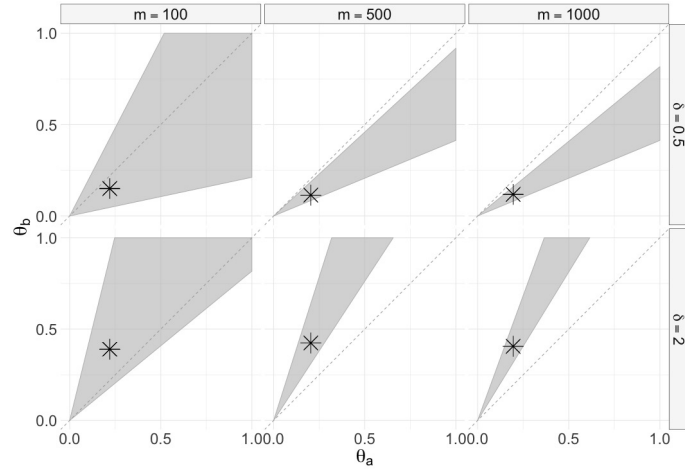
In this section some numerical examples of confidence sequences for the two types of effect sizes are given. All simulations were run with code available in our software package (Ly et al., 2022).

*Risk difference.* Risk difference is defined as the difference between success probabilities in the two streams:  $\delta = \theta_b - \theta_a$ . Figure 2 shows running intersections of confidence sequences with  $\delta$  as the risk difference for simulations for various distributions and stream lengths. These sequences are constructed by testing null hypotheses based on  $\vec{\Theta}_0(s, c)$ , with  $c = 1$  and  $s = \delta$ .  $\text{CI}_{\alpha, (m)}$  for the risk difference on  $\vec{\Theta}_0$  is an interval, corresponding to the ‘beam’ of  $(\theta_a, \theta_b) \in [0, 1]^2$  bounded by the lines  $\theta_b = \theta_a + \delta_L$  and  $\theta_b = \theta_a + \delta_R$  with  $\delta_L > \delta_R$  being values such that  $S_{[n_a, n_b, W_1; \vec{\Theta}_0(\delta_L)]}^{(m)} = S_{[n_a, n_b, W_1; \vec{\Theta}_0(\delta_R)]}^{(m)} = 1/\alpha$ . Figure B.1 in the Appendix illustrates that the running intersection indeed improves the confidence sequence, albeit slightly.

*Relative risk.* Relative risk is defined as the ratio between the success probabilities in group  $b$  and  $a$ :  $\delta = \theta_b/\theta_a$ . Hence, confidence sequences for this effect size measure can again be constructed using the linear



(a) Risk difference

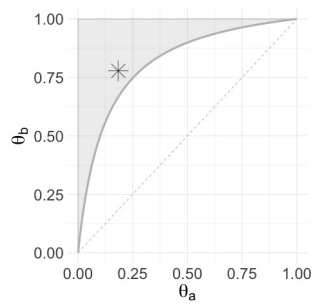


(b) Relative risk

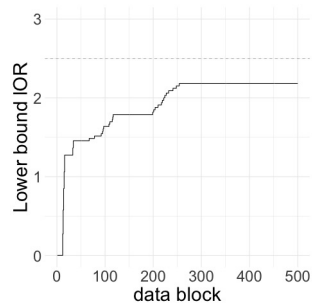
Figure 2: Depiction of parameter space with running intersection of confidence sequence for data generated under various effect sizes, at different time points  $m$  in a data stream. The asterisks indicate the maximum likelihood estimator at that time point. The significance threshold was set to 0.05. The design was balanced, with data block sizes  $n_a = 1$  and  $n_b = 1$ .

boundary form  $\vec{\Theta}_0(s, c)$  again, but now with  $s = 0$  and  $c = \delta$ . Figure 2 shows running intersections of confidence sequences with  $\delta$  as the relative risk.

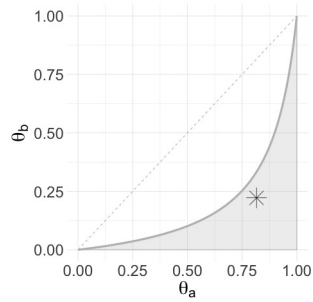
*Log odds ratio boundary.* If the maximum likelihood estimate based on  $Y^{(m)}$  lies in the upper left corner as in Figure 3(a), the confidence sets  $CS_{(m)}$  we get at time  $m$  have a one-sided shape such as the shaded region, or the shaded region in Figure 3(c), if the estimate lies in the lower right corner. Again, we can improve these confidence sequences by taking the running intersection; running intersections over time are illustrated in Figures 3(b) and 3(d).



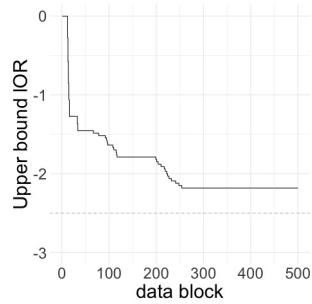
(a)  $CS^+$  at  $n = 500$ , true IOR 2.5



(b) Running lower bound  $CS^+$ , true IOR 2.5



(c)  $CS^-$  at  $n = 500$ , true IOR -2.5



(d) Running upper bound  $CS^-$ , true IOR -2.5

Figure 3: One-sided confidence sequences for odds ratios. 500 data blocks were generated under  $P_{\theta_a, \theta_b}$  with  $\theta_a = 0.2$  and log of the odds ratio (IOR) 2.5 for figures a and b, and  $\theta_a = 0.8$  and IOR -2.5 for figures c and d. The asterisks indicate the maximum likelihood estimator at  $n = 500$ . The significance threshold was set to 0.05. The design was balanced, with data block sizes  $n_a = 1$  and  $n_b = 1$ . Note that  $CS^-$  is empty for (a) and (b) and  $CS^+$  for (c) and (d) in these confidence sequences.

## 4. Conclusion

We have shown how E-variables for data streams can be extended to general null hypotheses and non-asymptotic always-valid confidence sequences. We specifically implemented the confidence sequences for the  $2 \times 2$  contingency tables setting; the resulting confidence sequences are efficiently computed and show quick convergence in simulations. For estimating risk differences or relative risk ratios between proportions in two groups, to our knowledge, such exact confidence sequences did not yet exist. For the log odds ratio we could also have used the sequential probability ratio (SPR) in Wald’s SPR test (Wald, 1945) test, which can be re-interpreted as a (product of) E-variables (Grünwald et al., 2022). However, the SPR does not satisfy the GRO property making it sub-optimal (see also (Adams, 2020)); moreover, as should be clear from the development, our method for constructing confidence sequences can be implemented for any effect size notion with convex rejection sets  $\vec{\Theta}_0(\leq \delta)$  and  $\vec{\Theta}_0(\geq \delta)$ , not just the log odds ratio. A main goal for future work is to use Theorem 2 to provide such sequences for sequential two-sample settings that go beyond the  $2 \times 2$  table.

## Acknowledgements

Funding: this work is part of the Enabling Personalized Interventions (EPI) project, which is supported by the Dutch Research Council (NWO) in the Commit2 - Data –Data2Person program under contract 628.011.028. Declarations of interest: none.

## References

- Reuben Adams. Safe hypothesis tests for  $2 \times 2$  contingency table. Master’s thesis, Delft Technical University, 2020.
- D.A. Darling and H. Robbins. Confidence sequences for mean, variance, and median. *Proc. National Academy of Sciences USA*, 58(1):66, 1967.
- Peter Grünwald, Rianne de Heide, and Wouter Koolen. Safe testing. *accepted, pending minor revision, for publication in Journal of the Royal Statistical Society: Series B*, 2022.
- Alexander Henzi and Johanna F. Ziegel. Valid sequential inference on probability forecast performance. *Biometrika*, 2021.
- Steven R Howard, Aaditya Ramdas, Jon McAuliffe, and Jasjeet Sekhon. Uniform, nonparametric, non-asymptotic confidence sequences. *Annals of Statistics*, 2021.
- J.Q. Li. *Estimation of Mixture Models*. PhD thesis, Yale University, New Haven, CT, 1999.
- Alexander Ly, Rosanne Turner, and Judith Ter Schure. R-package `safestats`, 2022. CRAN.
- Glenn Shafer. The language of betting as a strategy for statistical and scientific communication. *Journal of the Royal Statistical Society, Series A*, 2021.

Rosanne Turner, Alexander Ly, and Peter Grünwald. Generic e-variables for exact sequential k-sample tests that allow for optional stopping. *arXiv preprint arXiv:2106.02693*, 2021.

Vladimir Vovk and Ruodu Wang. E-values: Calibration, combination, and applications. *Annals of Statistics*, 2021.

Abraham Wald. Sequential tests of statistical hypotheses. *The Annals of Mathematical Statistics*, 16(2):117–186, 1945.

# Appendix

## A. Proofs

Both proofs below use Theorem 1 of Grünwald et al. (2022) and a direct corollary (called Corollary 2 by Grünwald et al. (2022)), which we re-state here, for convenience, combined as a single statement. Recall that we use notation  $P_W := \int P_{\vec{\theta}} dW(\vec{\theta})$ .

*Theorem (Theorem 1 of Grünwald et al. (2022)).* Let  $Y$  be a random variable taking values in a set  $\mathcal{Y}$ . Suppose  $Q$  is a probability distribution for  $Y$  with density  $q$  that is strictly positive on all of  $\mathcal{Y}$  and let  $\mathcal{H}_0 = \{P_{\vec{\theta}} : \vec{\theta} \in \vec{\Theta}_0\}$  be a set of distributions for  $Y$  where each  $P_{\vec{\theta}}$  has density  $p_{\vec{\theta}}$ . Let  $\mathcal{W}_0$  be the set of all distributions on  $\vec{\Theta}_0$ . Assume  $\inf_{W_0 \in \mathcal{W}_0(\vec{\Theta}_0)} D(Q \| P_{W_0}) < \infty$ . Then (a) there exists a (potentially sub-) distribution  $P_0^*$  with density  $p_0^*$  such that

$$S^* := \frac{q(Y)}{p_0^*(Y)}$$

is an E-variable ( $p_0^*$  is called the *Reverse Information Projection (RIPr)* of  $q$  onto  $\{p_W : W \in \mathcal{W}_0\}$ ). Moreover, (b),  $S^*$  satisfies

$$\sup_{S \in \mathcal{E}(\vec{\Theta}_0)} \mathbf{E}_{Y \sim Q}[\log S] = \mathbf{E}_{Y \sim Q}[\log S^*] = \inf_{W_0 \in \mathcal{W}_0(\vec{\Theta}_0)} D(Q \| P_{W_0}) = D(Q \| P_0^*). \quad (\text{A.1})$$

(where  $\mathcal{E}(\vec{\Theta}_0)$  is the set of all E-variables relative to null hypothesis  $\mathcal{H}_0$ ) and  $S^*$  is thus the  $Q$ -GRO E-variable for  $Y$ . If the minimum is achieved by some  $W_0^*$ , i.e.  $D(Q \| P_0^*) = D(Q \| P_{W_0^*})$ , then  $P_0^* = P_{W_0^*}$ . Moreover, (c), if there exists an E-variable  $S$  of the form  $q(Y)/p_{W_0}(Y)$  for some  $W_0 \in \mathcal{W}_0$  then  $W_0$  must achieve the infimum in (A.1) and  $S$  must be essentially equal to  $S^*$  in the sense that for all  $P \in \mathcal{H}_0 \cup \{Q\}$ ,  $P(S^* = q(Y)/p_{W_0}(Y)) = 1$ . Similarly (d), if there exists a  $W_0^* \in \mathcal{W}_0$  that achieves the infimum in (A.1) then  $S = q(Y)/p_{W_0^*}(Y)$  is an E-variable and  $S$  is again essentially equal to  $S^*$ .

*Proof of Theorem 1. Part 1* The real idea behind the proof is the formulation of the modified testing problem in which only a single outcome per block is observed. This we already did in the main text. Linking the two is simply the last, very simple step, with analogies to the proof of Part 1 of Theorem 1 in Turner et al. (2021).

Let  $n_a, n_b \in \mathbf{N}$ ,  $n := n_a + n_b$  and let  $u, v \in \mathbf{R}^+$ . Suppose that  $n_a u + n_b v \leq n$ . Then  $u^{n_a} v^{n_b} \leq 1$ , which follows immediately from applying Young's inequality to  $u^{n_a/n}, v^{n_b/n}$  but can also be derived directly by writing  $v$  as function of  $u$  and differentiating  $\log(u^{n_a} v^{n_b})$  to  $u$ .

Further, by independence, for  $(\theta_a, \theta_b) \in \vec{\Theta}_0$ ,

$$\begin{aligned}
& \mathbf{E}_{Y_a^{n_a} \sim P_{\theta_a}, Y_b^{n_b} \sim P_{\theta_b}} [s'(Y_a^{n_a}, Y_b^{n_b})] = \\
& \mathbf{E}_{Y_a^{n_a} \sim P_{\theta_a}} \left[ \frac{p_{\theta_a^*}(Y_a^{n_a})}{p^\circ(Y_a^{n_a}|a)} \right] \cdot \mathbf{E}_{Y_b^{n_b} \sim P_{\theta_b}} \left[ \frac{p_{\theta_b^*}(Y_b^{n_b})}{p^\circ(Y_b^{n_b}|b)} \right] = \\
& \left( \mathbf{E}_{Y \sim P_{\theta_a}} \left[ \frac{p_{\theta_a^*}(Y)}{p^\circ(Y|a)} \right] \right)^{n_a} \cdot \left( \mathbf{E}_{Y \sim P_{\theta_b}} \left[ \frac{p_{\theta_b^*}(Y)}{p^\circ(Y|b)} \right] \right)^{n_b} = \\
& \left( \mathbf{E}_{Y \sim P'_{\theta^*|a}} \left[ \frac{p'_{\theta^*}(Y|a)}{p^\circ(Y|a)} \right] \right)^{n_a} \cdot \left( \mathbf{E}_{Y \sim P'_{\theta^*|b}} \left[ \frac{p'_{\theta^*}(Y|b)}{p^\circ(Y|b)} \right] \right)^{n_b}. \tag{A.2}
\end{aligned}$$

Combining the two facts stated above, (6) implies that the latter quantity is bounded by 1.

*Part 2* By lower-semicontinuity of the KL divergence in its second argument (Posner's theorem, used as in Grünwald et al. (2022)) the infimum in (4) is achieved by some prior distribution  $W^\circ$  so that by Theorem 1 of Grünwald et al. (2022) (part (b) in the formulation above),  $p^\circ(\cdot | \cdot) = p'_{W^\circ}(\cdot | \cdot)$  and hence also  $P^\circ(G, Y) = P'_{W^\circ}(G, Y)$ . By convexity of  $\mathcal{H}'_0$  and finiteness of the support of  $P'_\theta(G, Y)$ , there must be some  $\vec{\theta}$  such that  $P'_{W^\circ}(G, Y) = P_{\vec{\theta}}(G, Y)$  and hence also  $p'_{W^\circ}(\cdot | \cdot) = p'_{\vec{\theta}}(\cdot | \cdot)$ , which shows (a). This means that we have now created an E-variable for the original problem which can be written as  $p_{\theta_a^*, \theta_b^*} / p_{W_0}$  with  $p_{W_0}$  a prior distribution on  $\vec{\theta}_0$  (namely, the one that puts mass 1 on  $\vec{\theta}$ ). (b) is then an immediate consequence of Theorem 1 of Grünwald et al. (2022) (part (c) in the formulation above). (note that we *cannot* draw this conclusion if  $\mathcal{H}'_0$  is not convex; for then the distribution  $p'_{W^\circ}$  may not correspond to the distribution  $p_{W^\circ}$  in the original problem — this correspondence is only guaranteed if  $p'_{W^\circ}$  coincides with some  $p'_{\vec{\theta}}$ ).

*Proof of Theorem 2.* Recall that we assume that  $\vec{\Theta}_0$  is convex and compact. We set  $\text{KL}'(\theta_a, \theta_b) := D(P'_{\theta_a^*, \theta_b^*} \| P'_{\theta_a, \theta_b})$  where  $D$  is the KL divergence as in (5), i.e. for the modified setting in which  $P'_{\theta_a, \theta_b}$  is a distribution on a single outcome, as discussed before Theorem 1. For the  $2 \times 2$  model this KL divergence can be written explicitly as

$$D(P'_{\theta_a^*, \theta_b^*} \| P'_{\theta_a, \theta_b}) = \frac{n_a}{n} \sum_{y_a \in \{0,1\}} p_{\theta_a^*}(y_a) \log \frac{p_{\theta_a^*}(y_a)}{p_{\theta_a}(y_a)} + \frac{n_b}{n} \sum_{y_b \in \{0,1\}} p_{\theta_b^*}(y_b) \log \frac{p_{\theta_b^*}(y_b)}{p_{\theta_b}(y_b)} \tag{A.3}$$

From (8) we now see that  $n\text{KL}'(\theta_a, \theta_b) = \text{KL}(\theta_a, \theta_b)$ . We will prove the theorem with KL replaced by  $\text{KL}'$  and  $\mathcal{H}_0$  by  $\mathcal{H}'_0$ ; since the two KL's agree up to a constant factor of  $n$ , all results transfer to the KL mentioned in the theorem statement.

Since  $\vec{\Theta}_0$  is compact in the Euclidean topology and all distributions in  $\mathcal{H}'_0$  can be represented as 2-dimensional vectors, i.e. they have common and finite support, we must have that  $\mathcal{H}_0$  is compact in the

weak topology so we can use the lower-semicontinuity of KL divergence in its second argument (Posner's theorem) as in (Grünwald et al., 2022) to give us that the minimum KL divergence  $\min \text{KL}'(\theta_a, \theta_b)$  is achieved by some  $(\theta_a^\circ, \theta_b^\circ)$ . Since KL divergence is strictly convex in its second argument and  $\mathcal{H}'_0$  is convex (this is the place where we need to use  $\text{KL}'$  rather than  $\text{KL}$ :  $\mathcal{H}_0$  may *not* be convex!), the minimum must be achieved uniquely. Since KL divergence  $\text{KL}'(\theta_a, \theta_b)$  is nonnegative and 0 only if  $(\theta_a, \theta_b) = (\theta_a^*, \theta_b^*)$ , it follows that  $(\theta_a^\circ, \theta_b^\circ) = (\theta_a^*, \theta_b^*)$  if  $\min \text{KL}(\theta_a, \theta_b) = 0$ . Otherwise, since we assume  $(\theta_a^*, \theta_b^*)$  to be in the interior of  $[0, 1]^2$ ,  $\text{KL}(\theta_a, \theta_b) = \infty$  iff  $(\theta_a, \theta_b)$  lies on the boundary of  $[0, 1]^2$ . Thus,  $(\theta_a^\circ, \theta_b^\circ)$  must lie in the interior of  $[0, 1]^2$  as well.  $(\theta_a^\circ, \theta_b^\circ)$  cannot lie in the interior of  $\vec{\Theta}_0$  though: for any point  $(\theta_a, \theta_b)$  in the interior of  $\vec{\Theta}_0$  we can draw a line segment between this point and  $(\theta_a^*, \theta_b^*)$ . Differentiation along that line gives that  $\text{KL}'(\theta_a, \theta_b)$  monotonically decreases as we move towards  $(\theta_a^*, \theta_b^*)$ , so the minimum within the closed set  $\vec{\Theta}_0$  must lie on its boundary.

It remains to show that (9) is the  $(\theta_a^*, \theta_b^*)$ -GRO E-variable relative to  $\mathcal{H}_0$ . To see this, note that, by convexity of  $\mathcal{H}'_0$ , from Theorem 1, we must have that the GRO E-variable for this original problem is of the form

$$\frac{p_{\theta_a^*}(y_a^{n_a})p_{\theta_b^*}(y_b^{n_b})}{p_{\theta_a^+}(y_a^{n_a})p_{\theta_b^+}(y_b^{n_b})}$$

for some  $(\theta_a^+, \theta_b^+)$ . The result then follows again by Theorem 1 of Grünwald et al. (2022) (part (c) in the formulation above): this shows that the distribution  $W_0$  that puts mass 1 on  $(\theta_a^+, \theta_b^+)$  minimizes, among all distributions  $W$  on  $\vec{\Theta}_0$ ,  $D(P_{\theta_a^*, \theta_b^*} \| P_W)$ . Since the set of such distributions includes all distributions that put mass 1 on *some*  $(\theta_a, \theta_b) \in \vec{\Theta}_0$ , we must have that  $(\theta_a^+, \theta_b^+) = (\theta_a^\circ, \theta_b^\circ)$ .



## B. Extended simulation results

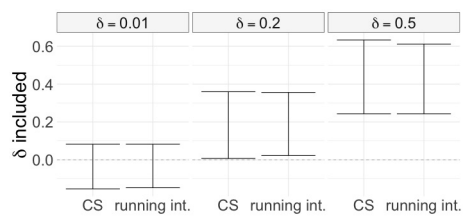


Figure B.1: Confidence sequence with and without running intersection, for data generated under  $P_{\theta_a, \theta_a + \delta}$  with  $\theta_a = 0.05$ , for a data stream of length 100. The significance threshold was set to 0.05. The design was balanced, with data block sizes  $n_a = 1$  and  $n_b = 1$ .