

Two-Sample Tests that are Safe under Optional Stopping, with an Application to Contingency Tables

Rosanne Turner^{a,b}, Alexander Ly^{a,c} and Peter Grünwald^{a,d*}

October 12, 2021

Abstract

We develop E-variables for testing whether two data streams come from the same source or not, and more generally, whether the difference between the sources is larger than some minimal effect size. These E-variables lead to tests that remain safe, i.e. keep their Type-I error guarantees, under flexible sampling scenarios such as optional stopping and continuation. In special cases our E-variables also have an optimal ‘growth’ property under the alternative. We illustrate the generic construction through the special case of 2×2 contingency tables, where we also allow for the incorporation of different restrictions on a composite alternative. Comparison to p-value analysis in simulations and a real-world example show that E-variables, through their flexibility, often allow for early stopping of data collection, thereby retaining similar power as classical methods.

*The authors gratefully acknowledge Reuben Adams, Rianne de Heide, Wouter Koolen, Muriel Perez and Judith ter Schure for useful conversations and in particular Adams and De Heide for performing experiments that inspired the E-variables presented here. 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.

^aCWI, Amsterdam, ^bUniversity Medical Center Utrecht, Brain Center, ^cUniversity of Amsterdam, Department of Psychology and ^dLeiden University, Department of Mathematics

1 Introduction

We develop hypothesis tests that are robust under flexible sampling scenarios, in which one is allowed to engage in optional continuation and optional stopping. We focus on the setting with data coming from two groups, the goal being to test whether the underlying distributions are the same or not. Our methodology is based on the notions of E-variables and test martingales. While to some extent going back as far as Darling and Robbins [1967], interest in these concepts has exploded only very recently, in part in relation to the ongoing replicability crisis in the applied sciences [Howard et al., 2021, Ramdas et al., 2020, Vovk and Wang, 2021, Shafer, 2021, Grünwald et al., 2019, Pace and Salvan, 2019, Manole and Ramdas, 2021, Henzi and Ziegel, 2021].

We collect samples from two distinct groups, denoted a and b . In the general setup we assume that in both groups, data are i.i.d. and come in sequentially — even though, as explained underneath (1.1) below, our approach can also be fruitfully used in the fixed design case. 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. $\sim P_{\theta_b}$ with $\theta_a, \theta_b \in \Theta$, $\{P_\theta : \theta \in \Theta\}$ representing some parameterized underlying family of distributions, all assumed to have a probability density or mass function denoted by p_θ on some outcome space \mathcal{Y} . We will use notation $P_{(\theta_a, \theta_b)}$ (density $p_{(\theta_a, \theta_b)}$) to represent the joint distribution of both streams. We consider the testing scenario, in which the null hypothesis \mathcal{H}_0 expresses that $\theta_a = \theta_b$ and the alternative \mathcal{H}_1 expresses that $d(\theta_a, \theta_b) > \delta$ for some divergence measure d and some effect size $\delta \geq 0$. We design a family of tests for this scenario that preserve type-I error guarantees under optional stopping. Hence, if the level α -test is performed and the null hypothesis holds true, the probability that the null will *ever* be rejected is bounded by α . While our tests can be implemented for arbitrary $\{P_\theta : \theta \in \Theta\}$, we extensively illustrate them on a simple, classical problem: 2×2 contingency tables. We provide simulations showing that if a standard fixed-design method for this scenario, the p-value resulting from Fisher’s exact test, is (ab)used with optional stopping, the type-I error blows up; in contrast, our tests retain type-I error guarantee while, due to the optional stopping, having power competitive with Fisher’s p-value.

Our test depends on the choice of a prior distribution on the alternative $\mathcal{H}_1 = \{P_{(\theta_a, \theta_b)} : (\theta_a, \theta_b) \in \Theta_1\}$ with $\Theta_1 \subset \{(\theta_a, \theta_b) : \theta_a, \theta_b \in \Theta\}$. The choice of prior does not affect the type-I error safety guarantee, hence it is fine, even from a frequentist point of view, if such a prior is chosen based on vague prior knowledge. Still, the prior affects how fast one will tend to reject the null if it is indeed false. For the case that no clear prior knowledge is available, one may use the prior that is optimal in terms of worst-case power or the related GRO criterion (definition in Section 1.1).

E-Variable Perspective; Block-wise Approach; Optional Continuation In its simplest form, an *E-variable* is a nonnegative random variable S such that under all distributions P in the null hypothesis,

$$\mathbf{E}_P[S] \leq 1. \tag{1.1}$$

Our test works by first designing E-variables for a *single block* of data, and then later extending these to sequences of blocks $Y_{(1)}, Y_{(2)}, \dots$ by multiplication. A block is a set of data consisting of n_a outcomes in group a and n_b outcomes in group b , for some pre-specified n_a and n_b . 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 (this rule can be loosened to some extent, see Section 2.1).

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*. An E-variable-based test at level α is then a test with, in combination with any stopping rule τ , reports ‘reject’ if and only if the product of E-values corresponding to all blocks that were observed so far and have already been completed, is larger than $1/\alpha$. The definition of τ may, and often will, be unknown to the user — the user only needs to get the signal to stop and can then report the product E-variable. We note though that if one stops ‘in the middle’ of an as-yet incomplete block, the data of that last block cannot yet be taken into account. A classical paired one-sample test corresponds to the special case with $n_a = n_b = 1$ and data coming in in the order a, b, a, b, \dots .

We can combine E-variables from different trials that share a common null (but may be defined relative to a different alternative) by multiplication, and still retain type-I error control. If we used p -values rather than E-variables we would have to resort to e.g. Fisher’s method for combining p -values, which, in contrast to multiplication of e -values, is invalid if there is a dependency between the (decision to perform) tests. With E-variables, such dependencies pose no problems for error control. Thus, in our setting, even if the design (i.e. n_a and n_b) is fixed in advance and optional stopping plays no role, we might still want to use the E-variable based tests described in this paper rather than a classic p -value based approach, since it allows us to do optional continuation over many experiments/studies (essentially, doing a meta-analysis [ter Schure et al., 2021]) while keeping type-I error control.

E-variables and test martingales are explained in more detail in Section 1.1 below, but we refer to Grünwald et al. [2019], Shafer [2021] for an extensive introduction to E -variables, their use in ‘optional continuation’ over several studies, and their enlightening *betting* interpretation (indeed, Shafer refers to E -variables as *betting scores*). The general story that emerges from these papers as well as, for example, [Vovk and Wang, 2021, Ramdas et al., 2020] is that E -variables and test martingales are the ‘right’ generalization of likelihood ratios to the case that both \mathcal{H}_0 and \mathcal{H}_1 can be composite and combination of data from several trials may be required.

Relevance of the 2×2 application Even in this age of big data and huge models, the lowly 2×2 model is still used as heavily as ever in clinical trials, psychological studies and so on — areas heavily plagued by the *reproducibility crisis* [Pace and Salvan, 2019]. In a by-now notorious questionnaire [John et al., 2012], more than 55% of the interviewed psychologists admitted to the practice of ‘adding data until the results look good’. While classical methods lose their type-I error guarantee if one does this (Figure 2 in Section 5),

our E-value based tests allow for it, while, due to the option of stopping early, remaining competitive in terms of sample sizes needed to obtain a desired power. We illustrate the practical advantage of our test in Section 6 using the recent real-world example of the SWEPIIS trial which was stopped early for harm [Wennerholm et al., 2019]. Their analysis being based on a p -value (by definition designed for fixed sampling plan), the question whether there was indeed sufficient evidence available to stop early is very hard to answer, since the sampling plan was not followed so that the p -value that led them to stop was by definition incorrectly calculated. This also makes it very difficult to combine the test results with results from earlier or future data while keeping anything like error control. We show that with our E-value based methodology we would have obtained sufficient evidence to stop for harm after the same number of events had occurred. Additionally, this E-value, even though based on a stopped trial, can be effortlessly combined with E-values from other trials while retaining error guarantees. Also, our results are of interest beyond mere testing: the E-variables we develop in this paper can be used to obtain *anytime-valid confidence intervals* [Howard et al., 2021] that also remain valid under optional stopping. We will report on this extension elsewhere.

An additional advantage of focusing on the 2×2 setting is that it is arguably the simplest and clearest example in which there is a nuisance parameter (the proportion under the null) that does not admit a group invariance. Nuisance parameters that satisfy such an invariance (such as the variance in the t -test, or the grand mean in the two-sample t -test) are quite straightforward to turn into E-variables and test martingales via the method of maximal invariants, as explained by Grünwald et al. [2019] and already put into practice by e.g. Robbins [1970], Lai [1976]. The present paper shows that the proportion under the null can also be handled in a clean and simple manner.

Finally, as explained below, our work appears to be quite different from existing sequential and Bayesian approaches. Thus, more than 85 years after *the lady tasting tea*, we are able to still say something quite new about the age-old problem of contingency table testing.

Related Work Sequential tests for the 2×2 setting that can be turned into test martingales (and would then be safe to use under optional stopping) have been suggested before [Barnard, 1946, Siegmund, 2013, Section V.2]. Yet, such earlier tests were based either on generalized likelihood ratios (which in general do not satisfy (1.1), hence they do not provide E-variables) or on skipping data points in which both groups have the same outcomes, which — as our (unreported) experiments confirm — is quite wasteful. Our E-variable based tests are entirely different in nature, and, in contrast to earlier approaches, are all *exact and nonasymptotic*. In fact our tests are more closely related to, yet still different from, Bayes factor tests: in the case of simple null hypotheses, E-variable based tests coincide with Bayes factors [Grünwald et al., 2019]. However, in the 2×2 setting the null is not simple, and while the Bayes factor is a ratio of two Bayes marginal likelihoods, our E-variables are ratios of more general, ‘prequential’ [Dawid, 1984] likelihood ratios. In some special cases, the numerator is still a Bayes marginal likelihood, but the denominator, in the 2×2 setting, almost never is. Thus, while similar in ‘look’, our approach is in the

end quite different from the default Bayes factors for tests of two proportions that were proposed by Kass and Vaidyanathan [1992] and by Jamil et al. [2017], the latter based on early work by Gunel and Dickey [1974]. To illustrate, in Appendix C (Supplementary Material) we show that none of the variants of the Gunel-Dickey Bayes factor that are applicable in our set-up yield valid E-variables.

Another, very recent, approach that bears some similarity to ours are the two-sample tests from Manole and Ramdas [2021]. They focus on a nonparametric setting and (in addition to many related results) provide always-valid tests and confidence intervals that avoid the block structure and thus allow stopping at any sample size. Their test martingales satisfy optimality properties as the sample size gets large. Instead, we focus on the parametric case and, for this case, manage to derive E-variables that are equal to or closely approximate the optimal (as measured according to the GRO criterion) E-variables, thus optimizing for the small-sample case (in principle, our tests could be used in a nonparametric setting as well, but since they rely on using a prior on the alternative, the test martingales of Manole and Ramdas [2021] might be easier to use in that case). Another general nonparametric two-sample approach with a sequential flavor (but without optional stopping error guarantees) is Lhéritier and Cazals [2018].

Contents In the remainder of this introductory section, we formally introduce E-variables, optional stopping and the concept of GRO-optimality. In Section 2 we propose almost GRO-optimal E-variables for tests of two streams in general. In Sections 3 and 4 we specifically show how these general E-variables can be applied in the setting of a test of two proportions, with and without restrictions on the alternative hypothesis. In Sections 5 and 6 we show through simulations and a real-world example comparisons of various E-variables and Fisher’s exact test with respect to power, and we end with a conclusion. All proofs are in Appendix A (Supplementary Material).

1.1 E-Variables and Test Martingales, Safety and Optimality

We first need to extend the notion of E-variable to random processes:

Definition 1. Let $\{Y_{(j)}\}_{j \in \mathbf{N}}$, with all $Y_{(j)}$ taking values in some set \mathcal{Y} , represent a discrete-time random process. Let \mathcal{H}_0 be a collection of distributions for the process $\{Y_{(j)}\}_{j \in \mathbf{N}}$. For all $j \in \mathbf{N}$, let $S_{(j)}$ be a non-negative random variable that is adapted to $\sigma(Y^{(j)})$, with $Y^{(j)} = (Y_{(1)}, \dots, Y_{(j)})$, i.e. there exists a function s such that $S_{(j)} = s(Y^{(j)})$.

1. We say that $S_{(j)}$ is an E-variable for $Y_{(j)}$ conditionally on $Y^{(j-1)}$ if for all $P \in \mathcal{H}_0$,

$$\mathbf{E}_P [S_{(j)} \mid Y_{(1)}, \dots, Y_{(j-1)}] \leq 1. \quad (1.2)$$

That is, for each $y^{(j-1)} \in \mathcal{Y}^{j-1}$, all $P_0 \in \mathcal{H}_0$, (1.1) holds with $S = s(y_{(1)}, \dots, y_{(j-1)}, Y_{(j)})$ and P set to $P_0 \mid Y^{(j-1)} = y^{(j-1)}$.

2. If, for each j , $S_{(j)}$ is an E -variable conditional on $Y_{(1)}, \dots, Y_{(j-1)}$, then we call the process $\{S_{(j)}\}_{j \in \mathbf{N}}$ a conditional E -variable process relative to the given \mathcal{H}_0 and $\{Y_{(j)}\}_{j \in \mathbf{N}}$ and we call $\{S^{(m)}\}_{m \in \mathbf{N}}$ with $S^{(m)} = \prod_{j=1}^m S_{(j)}$ the corresponding test martingale.

Henceforth, we omit the phrase ‘relative to \mathcal{H}_0 and $\{Y_{(j)}\}_{j \in \mathbf{N}}$ ’ whenever it is clear from the context. By the tower property of conditional expectation, one verifies that for any process of conditional E -variables $\{S_{(j)}\}_{j \in \mathbf{N}}$, we have for all m that the product $S^{(m)}$ is itself an ‘unconditional’ E -variable as in (1.1), i.e. $\mathbf{E}_P[S^{(m)}] \leq 1$ for all $P \in \mathcal{H}_0$. Definition 1 adapts and slightly modifies terminology from [Shafer et al., 2011]. As follows from that paper, in standard martingale terminology, what we call a test martingale is a non-negative supermartingale relative to the filtration induced by $\{Y_{(j)}\}_{j \in \mathbf{N}}$, with starting value 1.

Safety The interest in E -variables and test martingales derives from the fact that we have type-I error control irrespective of the stopping rule used: for any test martingale $\{S^{(j)}\}_{j \in \mathbf{N}}$, Ville’s inequality [Shafer, 2021] tells us that, for all $0 < \alpha \leq 1$, $P \in \mathcal{H}_0$,

$$P(\text{there exists } j \text{ such that } S^{(j)} \geq 1/\alpha) \leq \alpha. \quad (1.3)$$

Thus, if we measure evidence against the null hypothesis after observing j data units by $S^{(j)}$, and we reject the null hypothesis if $S^{(j)} \geq 1/\alpha$, then our type-I error will be bounded by α , no matter what stopping rule we used for determining j . We thus have type-I error control even if we use the most aggressive stopping rule compatible with this scenario, where we stop at the first j at which $S^{(j)} \geq 1/\alpha$ (or we run out of data, or money to generate new data). We also have type-I error control if the actual stopping rule is unknown to us, or determined by external factors independent of the data $Y_{(j)}$.

We will call any test based on $\{S^{(j)}\}_{j \in \mathbf{N}}$ and a (potentially unknown) stopping time τ that, after stopping, rejects iff $S^{(\tau)} \geq 1/\alpha$ a *level α -test that is safe under optional stopping*, or simply a *safe test*.

Example 1. Let P_0 and Q be any two distributions for the process $Y_{(1)}, Y_{(2)}, \dots$, and let $\mathcal{H}_0 = \{P_0\}$ represent a simple null. Let $S^{(m)}$ denote the likelihood ratio for m outcomes and $S_{(j)}$ its constituent factors, i.e.

$$S^{(m)} = \frac{q(Y^{(m)})}{p_0(Y^{(m)})} = \prod_{j=1}^m S_{(j)} \quad \text{with} \quad S_{(j)} = \frac{q(Y_{(j)} \mid Y^{(j-1)})}{p_0(Y_{(j)} \mid Y^{(j-1)})}$$

where $q(y_{(m)} \mid y^{(m-1)})$ denotes the conditional density corresponding to Q and $p_0(y_{(m)} \mid y^{(m-1)})$ the one corresponding to P_0 with respect to a common underlying measure. Then the likelihood ratio process $\{S^{(m)}\}_{m \in \mathbf{N}}$ constitutes a test martingale, and the process of conditional likelihoods $\{S_{(j)}\}$ is a conditional E -variable process relative to \mathcal{H}_0 . This can be immediately verified by directly calculating the conditional expectation of $S_{(j)}$ given $Y^{(j-1)}$, noticing that the densities $p_0(Y_{(j)} \mid Y^{(j-1)})$ cancel in the calculation.

GRO-Optimality, Simple \mathcal{H}_1 Just like for p -values, the definition of E -variables only requires explicit specification of \mathcal{H}_0 , not of an alternative hypothesis \mathcal{H}_1 . \mathcal{H}_1 becomes crucial once we distinguish between ‘good’ and ‘bad’ E -variables: E -variables have been designed to remain small, with high probability, under the null \mathcal{H}_0 . But if \mathcal{H}_1 rather than \mathcal{H}_0 is true, then ‘good’ E -variables should produce evidence (grow — because the larger the E -variable, the closer we are to rejecting the null) against \mathcal{H}_0 as fast as possible. To make this precise, first consider simple (singleton) $\mathcal{H}_1 = \{Q\}$. We start with the one-outcome setting of (1.1), i.e. we look at a single E -variable $S_{(j)}$ in isolation for a single outcome $Y_{(j)}$. Its optimality is measured in terms of

$$\mathbf{E}_Q[\log S_{(j)}], \tag{1.4}$$

and the E -variable which maximizes this quantity among all E -variables that can be written as functions of $Y_{(j)}$ (i.e. non-negative random variables satisfying (1.1)), assuming it exists, is called the *Growth Rate Optimal* E -variable for $Y_{(j)}$ relative to Q , or simply ‘ Q -GRO for $Y_{(j)}$ ’, and denoted as $S_{\text{GRO}(Q),(j)}$. More generally, E -variable $S^{(m)}$ is called *growth rate optimal* relative to Q for $Y^{(m)}$, or simply Q -GRO for $Y^{(m)}$, if, among all (unconditional) E -variables that can be written as a function of $Y^{(m)}$, it maximizes

$$\mathbf{E}_Q[\log S^{(m)}]. \tag{1.5}$$

We will denote this E -variable, if it exists, by $S_{\text{GRO}(Q)}^{(m)}$. The idea to maximize (1.5) goes back to Kelly [1956]; the GRO-terminology is from Grünwald et al. [2019]. The larger an E -variable or test martingale tends to be under the alternative, the better it scores in the GRO sense. Of course, the same would still hold if we were to replace the logarithm by another strictly increasing function. But there are various compelling reasons for why one should take a logarithm here — see Grünwald et al. [2019], Shafer [2021]. One interesting reason, not explicitly covered by these two papers, was already given by Breiman [1961] and is explained in detail by [ter Schure et al., 2021, Appendix B.1]: the Q -GRO test martingale, assuming it exists (i.e. Condition (1.6) below holds and data are i.i.d.), is also the test martingale which minimizes the expected number of data points needed before the null can be rejected if we use the safe test with the aggressive stopping rule described before (reject at the smallest j such that $S^{(j)} \geq 1/\alpha$). Thus, using the Q -GRO test martingale is quite analogous to employing a test that maximizes power — in Section 5 we provide some simulations to relate power to GRO. Note that we cannot directly use power in designing tests, since the notion of power requires a fixed sampling plan, which we will usually not have.

In ‘nice’ cases, the Q -GRO E -variable (1.5) for m outcomes can be obtained by multiplying the individual Q -GRO E -variables:

Proposition 1. *Let $\mathcal{H}_1 = \{Q\}$ be simple and \mathcal{H}_0 be potentially composite. Suppose the following condition holds (with p the density of P):*

$$\text{There exists a } P \in \mathcal{H}_0 \text{ such that } S_{(1)} = q(Y_{(1)})/p(Y_{(1)}) \text{ is an } E\text{-variable.} \tag{1.6}$$

Then $S_{(1)}$ is the only E-variable for $Y_{(1)}$ that can be written in this form (i.e. there is no $P' \neq P$ such that (1.6) holds), and $S_{(1)} = S_{\text{GRO}(Q),(1)}$ is the Q -GRO E-variable for $Y_{(1)}$. An E-variable of this form automatically exists if \mathcal{H}_0 is simple and, more generally, if \mathcal{H}_0 , restricted to a single outcome $Y_{(1)}$, is a convex set of distributions that is compact in the weak topology.

If we further assume that $Y_{(1)}, Y_{(2)}, \dots$ are i.i.d. according to all distributions in $\mathcal{H}_0 \cup \mathcal{H}_1$, then $S_{\text{GRO}(Q)}^{(m)} = \prod_{j=1}^m S_{\text{GRO}(Q),(j)}$, i.e. the Q -GRO optimal (unconditional) E-variable for $Y^{(m)}$ is the product of the individual Q -GRO optimal E-variables.

If Condition (1.6) holds and $Y_{(1)}, Y_{(2)}, \dots$ are i.i.d. according to all distributions in $\mathcal{H}_0 \cup \mathcal{H}_1$, it thus makes sense to define the Q -GRO test martingale to be the test martingale $(S_{\text{GRO}(Q)}^{(j)})_{j \in \mathbb{N}}$. We will then have that $S_{\text{GRO}(Q),(j)} = s_Q(Y_{(j)})$ for a fixed function $s_Q : \mathcal{Y} \rightarrow \mathbf{R}_0^+$.

Example 2. [Simple \mathcal{H}_1 and Simple \mathcal{H}_0] Consider $\mathcal{H}_1 = \{Q\}$ and simple $\mathcal{H}_0 = \{P_0\}$ and arbitrary Q' such that the $Y_{(j)}$ are i.i.d. according to P, Q and Q' . Then $S_{(j)} = q'(Y_{(j)})/p_0(Y_{(j)})$ is an E-variable for $Y_{(j)}$, irrespective of the definition of Q' , by the same argument as in Example 1. By the Proposition above, the Q -GRO E-variable for $Y_{(j)}$ is given by setting $q' = q$. Then $\mathbf{E}_Q[S_{\text{GRO}(Q),(j)}] = \mathbf{E}_{Y_{(j)} \sim Q}[\log q(Y_{(j)})/p_0(Y_{(j)})]$ also coincides with the KL divergence between Q and P_0 .

In Section 2 (Theorem 1) we develop functions s_Q (denoted $s(\cdot; n_a, n_b, \theta_a^*, \theta_b^*)$ there) for simple $\mathcal{H}_1 = \{Q\}$ so that $S_{Q,(1)} = s_Q(Y_{(1)})$ is an E-variable even though \mathcal{H}_0 is composite and not convex, so that Proposition 1 does not apply. Since we invariably assume the $Y_{(j)}$ are i.i.d., $S_{Q,(j)} := s_Q(Y_{(j)})$ is an E-variable as well and with $S_Q^{(m)} := \prod_{j=1}^m S_{Q,(j)}$, $(S_Q^{(m)})_{m \in \mathbb{N}}$ is a test martingale. The construction works for the general setting of two data streams discussed in the introduction, and for some special \mathcal{H}_0 (even though composite and nonconvex), the $S_{Q,(j)}$ will in fact be Q -GRO and $(S_Q^{(m)})_{m \in \mathbb{N}}$ will be the Q -GRO test martingale. These include the \mathcal{H}_0 that arise in the 2×2 setting, our main application. For other \mathcal{H}_0 , the E-variables $S_{Q,(j)}$ will not necessarily have the Q -GRO-property; they are designed to have (1.5) large, but it may be even larger for other E-variables.

GRO and Composite \mathcal{H}_1 In case \mathcal{H}_1 is composite, no direct analogue of the GRO-criterion for designing E-variables exists, since it is not clear under what distribution $Q \in \mathcal{H}_1$ we should maximize (1.5). In this paper, we deal with this situation by *learning* Q from the data in a Bayesian fashion. It is now convenient to write $\mathcal{H}_1 = \{P_\theta : \theta \in \Theta_1\}$ in a parameterized manner (accordingly, henceforth we shall write θ_1 -GRO E-variable instead of P_{θ_1} -GRO E-variable and $S_{\text{GRO}(\theta),(j)}$ instead of $S_{\text{GRO}(P_\theta),(j)}$). We will assume i.i.d. data, thus, if \mathcal{H}_1 were true, then data would be i.i.d. $\sim P_{\theta_1^*}$ for some $\theta_1^* \in \Theta_1$. Starting with a distribution W on Θ_1 , i.e. a prior, at each point in time j , we determine the Bayesian posterior $W | Y^{(j-1)}$ and use the Bayes predictive $P_{W|Y^{(j-1)}} := \int_{\Theta_1} P_\theta dW(\theta | Y^{(j-1)})$ as an estimate for the ‘true’ $P_{\theta_1^*}$. As is well-known, under conditions on W and \mathcal{H}_1 (which, if \mathcal{H}_1 is finite-dimensional parametric, are very mild), the posterior will concentrate around θ^* and hence $P_{W|Y^{(j-1)}}$ will resemble $P_{\theta_1^*}$ more and more, with very high probability, as more data becomes available.

At each point in time j , we use our current estimate $P_{W|Y^{(j-1)}}$ to design a conditional E-variable $S_{(j)}$. On an informal level, as long as $P_{W|Y^{(j-1)}}$ converges to the ‘true’ $P_{\theta_1^*}$, the $S_{(j)}$ will in fact also start to more and more resemble the E-variables $S_{\text{GRO}(\theta_1^*), (j)}$ we designed for $\mathcal{H}_1 = \{P_{\theta_1^*}\}$ and which were designed to have a large expected growth under the ‘true’ $P_{\theta_1^*}$. Assuming the convergence happens fast, we have that

$$\mathbf{E}_{Y^{(m)} \sim P_{\theta_1^*}} \left[\log S_{\text{GRO}(\theta_1^*)}^{(m)} - \log \prod_{j=1}^m S_{(j)} \right] \quad (1.7)$$

is small, i.e. we may expect that the test martingale $\prod_{j=1}^m S_{(j)}$ grows not much slower than $S_{\text{GRO}(\theta_1^*)}^{(m)} = \prod_{j=1}^m S_{\text{GRO}(\theta_1^*), (j)}$, the best test martingale (maximizing $\mathbf{E}_{Y^{(m)} \sim P_{\theta_1^*}} [\log S]$ over all E-variables S for $Y^{(m)}$) we could have used if we had known the true $P_{\theta_1^*}$ all along.

2 Two-Stream Safe Tests

Consider the two-stream setting introduced in the beginning of the paper. To formalize it further, we introduce *calendar time* $t = 1, 2, \dots$ and corresponding random variables V_t and G_t : at each t , we obtain an outcome V_t in \mathcal{Y} in group $G_t \in \{a, b\}$. Importantly though, at this point we make no assumptions about the relative ordering of outcomes from the two groups. At time t , we have that t_a , the number of a ’s that are observed so far, and t_b , the number of b ’s observed so far, satisfy $t_a + t_b = t$, but subject to this constraint we allow them coming in any order, e.g. first all a ’s, or first all b ’s, or interleaved. For example, with $t_a = 3$ and $t_b = 2$, we might have $V_1 = Y_{1,a}, V_2 = Y_{2,a}, V_3 = Y_{3,a}, V_4 = Y_{1,b}, V_5 = Y_{2,b}$ (all a s come first, $G_1 = G_2 = G_3 = a, G_4 = G_5 = b$) but also, for example $V_1 = Y_{1,a}, V_2 = Y_{1,b}, V_3 = Y_{2,a}, V_4 = Y_{3,a}, V_5 = Y_{2,b}$.

We thus have that the (marginal) probability of the first $t = t_a + t_b$ outcomes, given that t_a of these are in group a and t_b in group b , and writing $y^t = (y_1, \dots, y_t)$, is given by the probability density (or mass function)

$$p_{\theta_a, \theta_b}(y_a^{t_a}, y_b^{t_b}) := p_{\theta_a}(y_a^{t_a}) p_{\theta_b}(y_b^{t_b}) = \prod_{t=1}^{t_a} p_{\theta_a}(y_{t,a}) \prod_{t=1}^{t_b} p_{\theta_b}(y_{t,b}). \quad (2.1)$$

To indicate that random vector $(Y_a^{t_a}, Y_b^{t_b}) := (Y_{1,a}, \dots, Y_{t_a,a}, Y_{1,b}, \dots, Y_{t_b,b})$ has a distribution represented by (3.1) we write ‘ $Y_a^{t_a}, Y_b^{t_b} \sim P_{\theta_a^*, \theta_b^*}$ ’.

According to the *null hypothesis* $\mathcal{H}_0 = \{P_{\theta_a, \theta_b} : (\theta_a, \theta_b) \in \Theta_0\}$, $\Theta_0 = \{(\theta, \theta) : \theta \in \Theta\}$, both processes coincide. Thus, we have that $\theta_a^* = \theta_b^* = \theta_0$ for some $\theta_0 \in \Theta$ and then the density of data $y_a^{t_a}, y_b^{t_b}$ is given by $p_{\theta_0}(y_{1,a}, \dots, y_{t_a,a}, y_{1,b}, \dots, y_{t_b,b})$.

2.1 The simple E-variable for 2-stream–blocks

We first consider the case in which the alternative hypothesis is simple: $\Theta_1 = \{\theta_1\}$ for some fixed $\theta_1 = (\theta_a^*, \theta_b^*) \in \Theta^2$. Consider a fixed sample size of size n , and assume that we will

observe a block of n_a outcomes in group a and n_b outcomes in group b . In this case, we can define an E-variable as the likelihood ratio between $p_{\theta_a^*, \theta_b^*}$ and a carefully chosen distribution that is a product of mixtures of distributions from Θ_0 : for $n_a, n_b \in \mathbf{N}$, $n := n_a + n_b$ and $y_a^{n_a} = (y_{1,a}, \dots, y_{n_a,a}) \in \{0, 1\}^{n_a}$ and $y_b^{n_b} = (y_{1,b}, \dots, y_{n_b,b}) \in \{0, 1\}^{n_b}$, we define:

$$s(y_a^{n_a}, y_b^{n_b}; n_a, n_b, \theta_a^*, \theta_b^*) := \frac{p_{\theta_a^*}(y_a^{n_a})}{\prod_{i=1}^{n_a} \left(\frac{n_a}{n} p_{\theta_a^*}(y_{i,a}) + \frac{n_b}{n} p_{\theta_b^*}(y_{i,a}) \right)} \cdot \frac{p_{\theta_b^*}(y_b^{n_b})}{\prod_{i=1}^{n_b} \left(\frac{n_a}{n} p_{\theta_a^*}(y_{i,b}) + \frac{n_b}{n} p_{\theta_b^*}(y_{i,b}) \right)}. \quad (2.2)$$

Theorem 1. *The random variable $S_{[n_a, n_b, \theta_a^*, \theta_b^*]} := s(Y_a^{n_a}, Y_b^{n_b}; n_a, n_b, \theta_a^*, \theta_b^*)$ is an E-variable, i.e. with $s'(\cdot) = s(\cdot; n_a, n_b, \theta_a^*, \theta_b^*)$, we have:*

$$\sup_{\theta \in \Theta} \mathbf{E}_{V^n \sim P_\theta} [s(V^n; n_a, n_b, \theta_a^*, \theta_b^*)] \leq 1. \quad (2.3)$$

Moreover, if $\{P_\theta : \theta \in \Theta\}$ is a convex set of distributions, then $S_{[n_a, n_b, \theta_a^*, \theta_b^*]}$ is the (θ_a^*, θ_b^*) -GRO E-variable: for any non-negative function s' on $\mathcal{Y}^{n_a+n_b}$ satisfying $\sup_{\theta \in \Theta} \mathbf{E}_{V^n \sim P_\theta} [s'(V^n)] \leq 1$, we have:

$$\mathbf{E}_{Y_a^{n_a}, Y_b^{n_b} \sim P_{\theta_a^*, \theta_b^*}} [\log s(Y_a^{n_a}, Y_b^{n_b}; n_a, n_b, \theta_a^*, \theta_b^*)] \geq \mathbf{E}_{Y_a^{n_a}, Y_b^{n_b} \sim P_{\theta_a^*, \theta_b^*}} [\log s'(Y_a^{n_a}, Y_b^{n_b})].$$

Crucially, in the second part of the theorem, we do not require convexity of \mathcal{H}_0 , a set of distributions over $\mathcal{Y}^{n_a+n_b}$ (if \mathcal{H}_0 were convex, the GRO property would already follow by Proposition 2), but instead of $\{P_\theta : \theta \in \Theta\}$, a set of distributions on \mathcal{Y} . In the 2×2 case \mathcal{H}_0 is not convex, since the set of i.i.d. Bernoulli distributions over $n_a + n_b > 1$ outcomes is not convex; but $\{P_\theta : \theta \in \Theta\}$ is just the Bernoulli model on one outcome, which is convex.

To illustrate, consider the basic case in which data comes in in fixed batches $Y_{(1)}, Y_{(2)}, \dots$, with each batch $Y_{(j)} = ((Y_{(j-1)n_a+1,a}, Y_{(j-1)n_a+2,a}, \dots, Y_{jn_a,a}), (Y_{(j-1)n_b+1,b}, Y_{(j-1)n_b+2,b}, \dots, Y_{jn_b,b}))$, having exactly n_a outcomes in group a and n_b outcomes in group b , and let $n = n_a + n_b$. This case would obtain, for example, in a sequential clinical trial in which patients come in one by one, each odd patient is given the treatment and each even patient is given the placebo. Then $n = 2$, $n_a = n_b = 1$. We may then measure the evidence against the null hypothesis by the product E-value

$$S_{[n_a, n_b, \theta_a^*, \theta_b^*]}^{(m)} := \prod_{j=1}^m S_{(j), [n_a, n_b, \theta_a^*, \theta_b^*]} \quad ; \quad S_{(j), [n_a, n_b, \theta_a^*, \theta_b^*]} := s(Y_{(j)}; n_a, n_b, \theta_a^*, \theta_b^*). \quad (2.4)$$

By Ville's inequality (1.3), the probability under any distribution in the null that there is an m with $S_{[n_a, n_b, \theta_a^*, \theta_b^*]}^{(m)}$ larger than $1/\alpha$, is bounded by α , hence, type-I error guarantees are preserved under optional stopping if we perform the test based on $\{S_{[n_a, n_b, \theta_a^*, \theta_b^*]}^{(m)}\}_{m \in \mathbf{N}}$ as defined underneath (1.3), as long as we stop between and not 'within' batches (if we stop within a batch, the E-variable $S_{[n_a, n_b, \theta_a^*, \theta_b^*]}^{(m)}$ is undefined).

If the data do not come in batches of equal size, we may proceed as follows. First, we need to fix some $n_a \geq 1$ and $n_b \geq 1$ of our own choice. The treatment below will give valid

E -variables irrespective of our choice of n_a and n_b , but it will be seen that some choices are much more reasonable (will lead to much more evidence against the null, if the null is false) than others.

Thus, fix n_a and n_b , set $n = n_a + n_b$. At each time t , we will have observed, so far, some number t_a of outcomes in group a , and t_b in group b . Now let m_t be the largest m such that $mn_a \leq t_a$ and $mn_b \leq t_b$. Now, for $m = 1, 2, \dots$, define $Y_{(m)}$ as above. At any given time t , $Y_{(1)}, Y_{(2)}, \dots, Y_{(m_t)}$ will have been observed, and there may be a number n'_j remaining observations in group $j \in \{a, b\}$ so that either $n'_a < n_a$ or $n'_b < n_b$ or both. Since the $\{Y_{(j)}\}_{j \in \mathbf{N}}$ determine a test martingale in the sense of Definition 1, optional stopping while preserving type-I error guarantees is then possible at any point in time t , as long as the E -variable is calculated as (2.4) above for $m = m_t$, thus ignoring the final $n'_a + n'_b$ outcomes.

How should n_a and n_b be chosen in practice? For example, consider a variation of the clinical trial setting above in which the treatment-control assignment is randomized: for each incoming patient, a fair coin is flipped to decide treatment (a) or placebo (b). Then at any given time the number of patients in group a and b will not be precisely equal, but if we choose $n_a = n_b = 1$ as above it is highly unlikely that the amount of data we have to ignore at any given time t is very large. Similarly, if G_t , the group membership of the t -th observation is itself i.i.d. according to some distribution P^* , we might have some idea of the probability $p^*(a)$ assigned to group a ; if $p^*(a) = 2/5$ (say), we would choose $n_a = 2, n_b = 3$.

We can add a significant amount of extra flexibility by allowing for variable group sizes, i.e., the chosen n_a and n_b may depend on the past. For this, we introduce a function $f : \bigcup_{t \geq 0} \mathcal{Y}^t \times \{0, 1\}^t \rightarrow \{\text{STOP-BLOCK}, \text{CONTINUE}\}$ that, at each point in time t , decides whether the current block should end ($f(V^t, G^t) = \text{STOP-BLOCK}$) or not ($f(V^t, G^t) = \text{CONTINUE}$). As long as the value of this function does not depend on the actual outcomes V_t observed after the last block that was completed, all requirements for having a test martingale and thus for safe optional stopping are met. For example, suppose that on data $V_1, G_1, V_2, G_2, \dots, V_t, G_t$ observed so-far, f has output STOP-BLOCK at m occasions, the last time at $t' = t - k$ for some $k > 0$. Then $f(t)$ is allowed to depend on $Y^{(m)}$ and G^t , but for any fixed $Y^{(m)} = y^{(m)}, G^t = g^t$, for all $y^k, y'^k \in \mathcal{Y}^k$, we must have $f((y^{(m)}, y^k), g^t) = f((y^{(m)}, y'^k), g^t)$. In this way, one can in principle *learn* $p^*(a)$ from the data, changing group sizes n_a and n_b flexibly as data come in. For simplicity, we have not followed this approach here, but all our results readily extend to this case.

2.2 The simple E-variable with Bayesian alternative

Now fix some prior W_1 with density w_1 on the alternative $\Theta_1 \subseteq \Theta^2$. We can trivially extend the definition of simple E-variable relative to singleton (θ_a^*, θ_b^*) to a simple E-variable relative to arbitrary prior W_1 on (θ_a^*, θ_b^*) : define $p_{W_1, a}(y) := \int p_{\theta_a}(y) dW_1(\theta_a)$, the integration being over the marginal prior distribution over θ_a , and similarly, $p_{W_1, b}(y) := \int p_{\theta_b}(y) dW_1(\theta_b)$.

Then, as a corollary of Theorem 1,

$$s(y_a^{n_a}, y_b^{n_b}; n_a, n_b, W_1) := \frac{\prod_{i=1}^{n_a} p_{W_{1,a}}(y_{i,a})}{\prod_{i=1}^{n_a} \left(\frac{n_a}{n} p_{W_{1,a}}(y_{i,a}) + \frac{n_b}{n} p_{W_{1,b}}(y_{i,a}) \right)} \cdot \frac{\prod_{i=1}^{n_b} p_{W_{1,b}}(y_{i,b})}{\prod_{i=1}^{n_b} \left(\frac{n_a}{n} p_{W_{1,a}}(y_{i,b}) + \frac{n_b}{n} p_{W_{1,b}}(y_{i,b}) \right)}. \quad (2.5)$$

is itself also an E-variable, as follows from applying Theorem 1 with a ‘meta’-set of distributions, which is possible since we made no assumptions at all on the set Θ in Theorem 1: we replace Θ by $\mathcal{W}(\Theta)$, the set of distributions on Θ ; we replace the background set of distributions $\{p_\theta : \theta \in \Theta\}$ by the set of distributions $\{p_W : W \in \mathcal{W}(\Theta)\}$; we replace the simple $\mathcal{H}_1 = \{P_{\theta_a^*, \theta_b^*}\}$ by a ‘simple’ $\mathcal{H}'_1 = \{P_{W_a, W_b}\}$ for some distributions W_a and W_b on Θ . Such W_1 -based simple E-variables can be used to *learn* the parameters θ_a^*, θ_b^* as more data in both streams come in, and this is how we will use them in a sequential context with optional stopping. Thus, assume again that data comes in batches $Y_{(1)}, Y_{(2)}, \dots$ with each $Y_{(j)}$ consisting of n_a outcomes in group a and n_b outcomes in group b (generalization to flexible group sizes changing in time and depending on the past as described at the end of Section 2.1 is straightforward). We start with some prior W_1 for the first batch $Y_{(1)}$ but we now use, for the j -th batch $Y_{(j)}$, the *Bayesian posterior* $W_1 \mid Y^{(j-1)}$ as prior to define the j -th E-variable with:

$$S_{[n_a, n_b, W_1]}^{(m)} := \prod_{j=1}^m S_{(j), [n_a, n_b, W_1]} \quad ; \quad S_{(j), [n_a, n_b, W_1]} := s(Y_{(j)}; n_a, n_b, W_1 \mid Y^{(j-1)}). \quad (2.6)$$

Again, $\{S_{(j), [n_a, n_b, W_1]}\}_{j \in \mathbb{N}}$ is a conditional E-variable process, so testing based on the corresponding test martingale is safe under optional stopping by (1.3). If data are sampled from some alternative hypothesis (θ_a^*, θ_b^*) , then as data accumulates, the posterior W_1 will, with high probability, concentrate narrowly around (θ_a^*, θ_b^*) and so $S_{(j), [n_a, n_b, W_1]}$ will behave more and more similarly to the ‘best’ (θ_a^*, θ_b^*) E-variable. Still, with the exception of a special case we indicate below, in general we cannot expect it to be the W_1 -GRO E-variable. But we are not particularly concerned by this: our experiments in Section 5 indicate that, at least in the 2×2 table setting, it behaves quite well in terms of power, which is often the main practical interest.

Simplification when $\{P_\theta : \theta \in \Theta\}$ is Convex Denoting $W_{1,g} \mid Y^{(m)}$ as the marginal posterior for θ_g , for $g \in \{a, b\}$, we can rewrite (2.6) as

$$S_{[n_a, n_b, W_1]}^{(m)} = \prod_{j=1}^m \frac{\prod_{i=1}^{n_a} p_{W_{1,a} \mid Y^{(j-1)}}(Y_{(j-1)n_a+i,a}) \prod_{i=1}^{n_b} p_{W_{1,b} \mid Y^{(j-1)}}(Y_{(j-1)n_b+i,b})}{\prod_{g \in \{a,b\}} \prod_{i=1}^{n_g} \frac{n_a}{n} \left(p_{W_{1,a} \mid Y^{(j-1)}}(Y_{(j-1)n_g+i,g}) + \frac{n_b}{n} p_{W_{1,b} \mid Y^{(j-1)}}(Y_{(j-1)n_g+i,g}) \right)}$$

if $\{P_\theta : \theta \in \Theta_0\}$ convex $\stackrel{=}{=} \prod_{j=1}^m \prod_{i=1}^{n_a} \frac{p_{W_{1,a} \mid Y^{(j-1)}}(Y_{(j-1)n_a+i,a})}{p_{\theta_0 \mid Y^{(j-1)}}(Y_{(j-1)n_a+i,a})} \prod_{i=1}^{n_b} \frac{p_{W_{1,b} \mid Y^{(j-1)}}(Y_{(j-1)n_b+i,b})}{p_{\theta_0 \mid Y^{(j-1)}}(Y_{(j-1)n_b+i,b})} \quad (2.7)$

with $\check{\theta}_0|Y^{(j-1)} \in \Theta$ s.t. $p_{\check{\theta}_0|Y^{(j-1)}} = (n_a/n)p_{W_{1,a}|Y^{(j-1)}} + (n_b/n)p_{W_{1,b}|Y^{(j-1)}}$, the existence of $\check{\theta}_0|Y^{(j-1)}$ being immediate if $\{P_\theta : \theta \in \Theta\}$ is convex. This rewrite will enable several additional results for such Θ .

Connection to Bayes Factors Consider W_1 such that θ_a and θ_b are independent under W_1 with marginal distributions W_a and W_b , and now further take $n_a = n_b = 1$. By basic telescoping, we can then further rewrite (2.6) as

$$\frac{\int p_{\theta_a}(Y_a^m)dW_a(\theta_a) \int p_{\theta_b}(Y_b^m)dW_b(\theta_b)}{\prod_{j=1}^m \prod_{g \in \{a,b\}} \left(\frac{1}{2}p_{W_{1,a}|Y^{(j-1)}}(Y_{j,g}) + \frac{1}{2}p_{W_{1,b}|Y^{(j-1)}}(Y_{j,g}) \right)} \stackrel{\text{if } \mathcal{H}_0 \text{ convex}}{=} \frac{\int p_{\theta_a}(Y_a^m)dW_a(\theta_a) \int p_{\theta_b}(Y_b^m)dW_b(\theta_b)}{\prod_{j=1}^m \prod_{g \in \{a,b\}} p_{\check{\theta}_0|Y^{(j-1)}}(Y_{j,g})} = \prod_{j=1}^m \frac{p_{W_a|Y^{(j-1)}}(Y_{j,a}) \cdot p_{W_b|Y^{(j-1)}}(Y_{j,b})}{\prod_{g \in \{a,b\}} p_{\check{\theta}_0|Y^{(j-1)}}(Y_{j,g})}. \quad (2.8)$$

The numerator of the simple product E-value is now equal to the Bayesian marginal likelihood of the data based on prior W_1 . In general, this rewrite of the numerator breaks down if n_a or n_b is larger than one and/or θ_a and θ_b are dependent under the prior. The reason is that according to (2.6) (see the line above (2.7)), the outcomes in each group within each of the m blocks are treated as independent (they have the same density). In contrast, while a Bayes factor for m blocks of n data points can also be rewritten as a product of $n \cdot m$ Bayes predictive densities, in general it makes every outcome dependent on every previous outcome — the posterior, and hence the posterior predictive, is updated also within each block. If θ_a and θ_b are independent under the prior though, then they are also independent under the posterior, and if $n_a = n_b = 1$ then the Bayes predictive densities for the two outcomes within each block will also be independent, and we get (2.8).

Thus, in this special case, if the denominator could also be written as a Bayes marginal likelihood, then our E-variable would really be a Bayes factor. Yet, even if $\{P_\theta : \theta \in \Theta\}$ is convex, it cannot be written in this way, though it is very ‘close’: each of the m factors in the denominator in (2.8) is the product density function of two identical distributions for one outcome, and Proposition 2 shows that, in the special case of the 2×2 model with W_a and W_b independent beta priors, this distribution may itself be the Bayes predictive distribution obtained by equipping Θ_0 with another beta prior. Still, for a real Bayes factor corresponding to \mathcal{H}_0 , for each j , the two outcomes $Y_{j,a}, Y_{j,b}$ in the j -th block would again not be independent given $Y^{(j-1)}$, whereas in (2.8) they are.

3 Safe tests for two proportions

We assume the setting above, but now we further assume that both streams are Bernoulli. This will substantially simplify the formulae. Thus, $\Theta = [0, 1]$ and (2.1) now specializes to

$$p_{\theta_a, \theta_b}(y_a^{t_a}, y_b^{t_b}) := p_{\theta_a}(y_{1,a}, \dots, y_{t_a,a})p_{\theta_b}(y_{1,b}, \dots, y_{t_b,b}) = \theta_a^{t_a}(1 - \theta_a)^{t_a - t_a} \theta_b^{t_b}(1 - \theta_b)^{t_b - t_b}. \quad (3.1)$$

with t_{a1} the number of outcomes 1 in stream a among the first t_a ones, and t_{b1} the number of outcomes 1 in stream b among the first t_b ones. According to the null hypothesis, we have that $\theta_a^* = \theta_b^* = \theta_0$ for some $\theta_0 \in \Theta = [0, 1]$. (3.1) now simplifies to:

$$p_{\theta_0}(y_a^{t_a}, y_b^{t_b}) := \theta_0^{t_1} (1 - \theta_0)^{t_0}, \quad (3.2)$$

with t_1 the number of ones in the sequence $y^{t_a+t_b} = y_1, \dots, y_{t_a+t_b}$, and similarly for t_0 .

We now run through the results of the previous section for this instantiation of our test. Again, we start with the case of a simple $\mathcal{H}_1 = \{P_{\theta_a^*, \theta_b^*}\}$. (2.2) considerably simplifies and can be written as:

$$s(y_a^{n_a}, y_b^{n_b}; n_a, n_b, \theta_a^*, \theta_b^*) := \frac{p_{\theta_a^*}(y_a^{n_a})}{p_{\theta_0}(y_a^{n_a})} \cdot \frac{p_{\theta_b^*}(y_b^{n_b})}{p_{\theta_0}(y_b^{n_b})}, \quad \text{where } \theta_0 = \frac{n_a}{n} \theta_a^* + \frac{n_b}{n} \theta_b^*. \quad (3.3)$$

Theorem 1 tells us that this is an E-variable. Since $\{P_\theta : \theta \in \Theta\}$, the Bernoulli model, is convex, the theorem also tells us that in this case the simple E-variable with simple alternative is always (θ_a^*, θ_b^*) -GRO.

We now turn to the simple E-variable relative to arbitrary prior W_1 . For the Bernoulli model the Bayes posterior predictive distribution is itself a Bernoulli distribution, with its parameter equal to the posterior mean. Therefore, we again get a considerable simplification: the simple E-variable relative to prior W_1 is still given by (2.5), but this now simplifies to:

$$s(y_a^{n_a}, y_b^{n_b}; n_a, n_b, W_1) := s(y_a^{n_a}, y_b^{n_b}; n_a, n_b, \theta_a^*, \theta_b^*) \text{ for } \theta_g^* = \mathbf{E}_{\theta_g \sim W_1}[\theta_g], g \in \{a, b\}. \quad (3.4)$$

Combining this with (2.7) we infer that

$$S_{[n_a, n_b, W_1]}^{(m)} = \prod_{j=1}^m \prod_{i=1}^{n_a} \frac{p_{\check{\theta}_a | Y^{(j-1)}}(Y_{(j-1)n_a+i, a})}{p_{\check{\theta}_0 | Y^{(j-1)}}(Y_{(j-1)n_a+i, a})} \prod_{i=1}^{n_b} \frac{p_{\check{\theta}_b | Y^{(j-1)}}(Y_{(j-1)n_b+i, b})}{p_{\check{\theta}_0 | Y^{(j-1)}}(Y_{(j-1)n_b+i, b})} \quad (3.5)$$

where $\check{\theta}_a | Y^{(j-1)} = \mathbf{E}_{\theta_a \sim W | Y^{(j-1)}}[\theta_a]$ and $\check{\theta}_b | Y^{(j-1)} = \mathbf{E}_{\theta_b \sim W | Y^{(j-1)}}[\theta_b]$ and $\check{\theta}_0 | Y^{(j-1)} = (n_a/n)\check{\theta}_a | Y^{(j-1)} + (n_b/n)\check{\theta}_b | Y^{(j-1)}$.

Simplified Calculations with Independent Beta Priors Now take the special case in which θ_a and θ_b are independent under the prior W_1 with marginals W_a and W_b . In this case, θ_a and θ_b are also independent under the posterior, and we can simplify $\check{\theta}_a | Y^{(j-1)} = \mathbf{E}_{\theta_a \sim W_a | Y_a^{(j-1)n_a}}[\theta_a]$, the expectation of θ_a under the posterior W_a given all data so far in group a , and similarly for group b . Using beta priors, this expectation is easy to calculate and we get:

Proposition 2. *Let θ_a, θ_b be independent under W_1 , with marginals W_a and W_b respectively. Suppose that these are beta priors with parameters (α_a, β_a) and (α_b, β_b) respectively. Then, upon defining $U_a = \sum_{i=1}^{(j-1)n_a} Y_{i,a}$, $U_b = \sum_{i=1}^{(j-1)n_b} Y_{i,b}$, $U = \sum_{i=1}^{(j-1)n} (Y_{i,a} + Y_{i,b})$ we have that $\check{\theta}_a, \check{\theta}_b, \check{\theta}_0$ as above satisfy: $\check{\theta}_a | Y^{(j-1)} = (U_a + \alpha_a) / ((j-1)n_a + \alpha_a + \beta_a)$, $\check{\theta}_b | Y^{(j-1)} = (U_b + \alpha_b) / ((j-1)n_b + \alpha_b + \beta_b)$ respectively, and $\check{\theta}_0 | Y^{(j-1)}$ is as further above. In the special case that we fix the prior parameters in the groups proportional to the group size fraction $\kappa := n_b/n_a$, i.e we fix $\alpha_b = \kappa\alpha_a$, $\beta_b = \kappa\beta_a$, the expression for $\check{\theta}_0$ simplifies to $\check{\theta}_0 | Y^{(j-1)} = (U + (1 + \kappa)\alpha_a) / ((j-1)n + (1 + \kappa)\alpha_a + (1 + \kappa)\beta_a)$.*

4 (Un)Restricted Composite \mathcal{H}_1 in the 2×2 setting

In this section we describe the main instantiations of the 2×2 stream testing scenario that are relevant in practice. These differ in the choice of \mathcal{H}_1 : the choice can be fully unrestricted (we simply want to find whether there is any discrepancy from \mathcal{H}_0 at all); restricted in terms of effect size; or restricted because we have prior knowledge about either θ_a^* or θ_b^* . We consider each in turn, the second and third scenario in a separate subsection. Section 5 provides extensive numerical simulations for all three scenarios.

In the first scenario, a researcher wants to perform a *two-sided test*; they simply aim to find any discrepancy from \mathcal{H}_0 if it exists, with no restrictions are placed on \mathcal{H}_1 . In this case, if we choose W_1 as independent beta priors on θ_a and θ_b , we can simply proceed as described in Proposition 2 above, taking a beta prior for simplicity. We will develop a reasonable ‘default’ choice for the hyper parameters by experiment in Section 5.

4.1 Dealing with Effect Sizes

In the second scenario that we will put to the test, we really want to test \mathcal{H}_0 against a restricted \mathcal{H}_1 consisting of those hypotheses that have a certain minimal *effect size* δ . This would then be a one-sided test. For example, a researcher might know that a new treatment must cure at least a certain number of patients more compared to a control treatment to provide a *clinically relevant treatment effect* δ . In this case, \mathcal{H}_1 could be restricted to either of the sets $\Theta(\delta)$ or $\Theta^+(\delta)$, where

$$\Theta(\delta) = \{\theta \in [0, 1]^2 : d(\theta) = \delta\} \quad ; \quad \Theta^+(\delta) = \begin{cases} \{\theta \in [0, 1]^2 : d(\theta) \geq \delta\} & \text{if } \delta > 0 \\ \{\theta \in [0, 1]^2 : d(\theta) \leq \delta\} & \text{if } \delta < 0, \end{cases} \quad (4.1)$$

where we set $d((\theta_a, \theta_b)) = \theta_b - \theta_a$. A second notion of effect size that often will be applicable in this sort of research is the *log odds ratio* between θ_b and θ_a , with restricted parameter space again given by (4.1) but d set to $d((\theta_b, \theta_a)) = \log [(\theta_b/(1 - \theta_b))((1 - \theta_a)/\theta_a)]$. These are the two effect size notions that will feature in our experiments. An illustration of both divergence measures and the resulting restricted parameter spaces is given in Figure 1.

A third popular notion of effect size, the relative risk, behaves, for small θ_a and $\delta > 0$, very similarly to the odds ratio, and will therefore not be separately considered in our experiments. We will also not consider two-sided testing against a restricted \mathcal{H}_1 (i.e., one wants to restrict the alternative hypothesis to a treatment being either substantially better *or* substantially worse than a control) since this is not a very common scenario. Such E-variables for the 2×2 setting and stream data could however be constructed with a method analogous to the one we describe below, by combining two ‘simple’ E-variables [Turner, 2019].

If we pick \mathcal{H}_1 restrict to $\Theta(\delta')$, then we could simply use the beta prior mentioned before with support conditioned on this set. What about the more realistic case of a \mathcal{H}_1 with $\delta \in \Theta^+(\delta')$? A first, intuitive (and certainly defensible) approach would be to use a prior W'_1 that is spread out over $\Theta^+(\delta')$, e.g. (if $\delta' > 0$) the beta prior as above conditioned

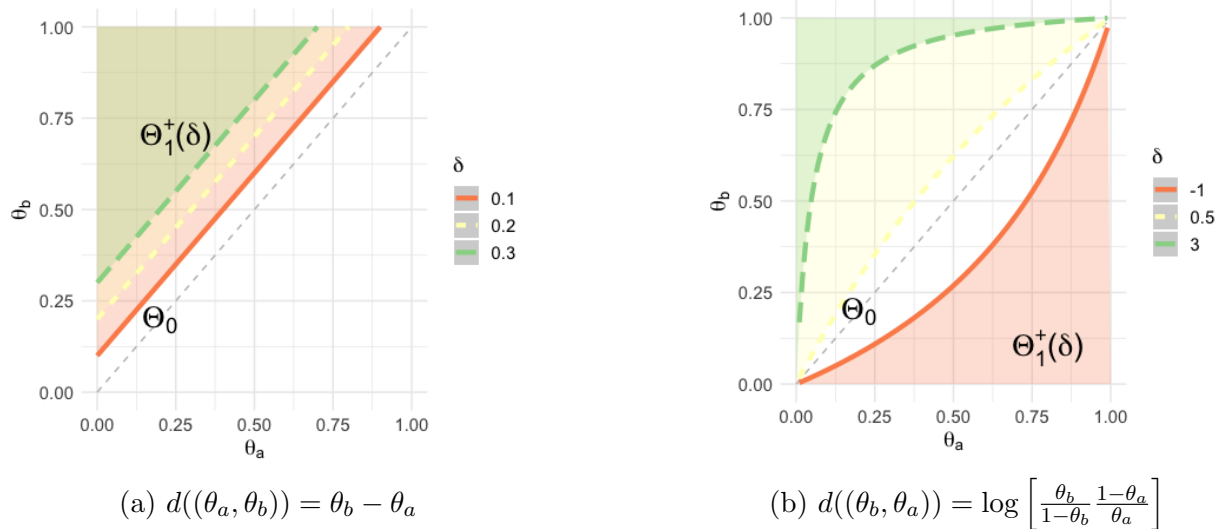


Figure 1: Examples of restricted alternative hypothesis parameter spaces for several values of two divergence measures; the difference between group means and the log odds ratio. Θ_0 denotes the null hypothesis parameter space; $\Theta_1^+(\delta)$ the restricted alternative hypothesis parameter space.

on $\delta \geq \delta'$. However, in terms of the GRO criterion, there are good reasons to still use a prior W_1^* that puts all prior mass on $\Theta(\delta')$, the boundary of the real parameter space $\Theta(\delta')$. Namely, for the resulting E-variable process $S_{[n_a, n_b, W_1^*]}^{(1)}, S_{[n_a, n_b, W_1^*]}^{(2)}, \dots$, it holds for every m that

$$\text{for all } (\theta_a, \theta_b) \text{ with } d((\theta_a, \theta_b)) > \delta', \mathbf{E}_{Y^{(m)} \sim P_{(\theta_a, \theta_b)}}[\log S_{[n_a, n_b, W_1^*]}^{(m)}] \geq \min_{\theta \in \Theta(\delta')} \mathbf{E}_{Y^{(m)} \sim P_\theta}[\log S_{[n_a, n_b, W_1^*]}^{(m)}]. \quad (4.2)$$

Thus, we might want to use the prior W_1^* also if δ can be more extreme than δ' , since if δ is actually more extreme, the expected evidence against \mathcal{H}_0 using W_1^* (even though designed for δ') will actually get larger anyway.

The advantage of the first approach is that it will lead to much higher GROwth ($\mathbf{E}_{P_{(\theta_a, \theta_b)}}[\log S_{[n_a, n_b, W_1']}]$ much larger than $\mathbf{E}_{P_{(\theta_a, \theta_b)}}[\log S_{[n_a, n_b, W_1^*]}^{(m)}]$) if we are ‘lucky’ and $|d(\theta_a, \theta_b)| \gg |\delta'|$. The price to pay is that it will lead to somewhat smaller growth if $d((\theta_a, \theta_b))$ is close to δ' (experiments omitted). It is easy to see why this is the case: the prior W_1' must spread out its mass over a much larger subset of $[0, 1]^2$ than W_1^* . Therefore, the E-variables based on W_1' will perform somewhat worse than those based on W_1^* if the data are sampled from a point (θ_a^*, θ_b^*) in the support of W_1^* , simply because W_1^* gives much larger prior support in a neighborhood of (θ_a^*, θ_b^*) . For this reason, and also because it is computationally a lot simpler, we decided to focus our experiments on the second approach rather than the first.

Calculating the prior and posterior for restricted \mathcal{H}_1 For both notions of effect size, θ_a and θ_b can no longer be independent for any prior on $\Theta(\delta)$. Hence, the prior and posterior do not longer admit the composition in terms of beta densities as in Proposition 2. For example, when putting a prior on $\Theta(\delta)$ with the additive effect size notion, we know the new domain of θ_a would be $[0, 1 - \delta]$. θ_b is completely determined by θ_a and δ in this case. We will still use a beta prior on $\Theta(\delta)$ and calculate posteriors by a numerical approach, explained in Appendix B Supplementary Material).

4.2 Working with Restrictions on \mathcal{H}_1

In practice, researchers often already have estimates of the occurrence rate of events in the *control group* in their experiments; for example, estimates of the proportion of patients that recover from a disease under standard care are known, and researchers investigate whether the proportion of recovered patients is higher in a group receiving an experimental treatment. This restriction on θ_a can be incorporated in the E-variable. This incorporation becomes especially easy if \mathcal{H}_1 is already restricted to a set $\Theta^+(\delta')$ with minimal relevant effect size δ' . For then $\Theta(\delta')$ contains just one point (θ_a^*, θ_b^*) (in the case of the linear effect size, this is $(\theta_a, \theta_a + \delta)$), and the E-variable constructed according to the guidelines of the previous subsection, which puts all its mass on δ' even though we allow $\delta \geq \delta'$, would be the simple E-variable corresponding to putting prior mass 1 on (θ_a^*, θ_b^*) .

5 Illustration via Simulated Data

In this section, we illustrate properties of our E-variables for 2×2 application through simulated data.¹ First, we determine a reasonable choice of beta prior hyper-parameter to use in (3.5) in terms of our GRO-criterion. Second, we show that in the optional continuation setting, as predicted by theory, type-I error control is achieved irrespective of the restriction on \mathcal{H}_1 used. Thereafter, we show by more simulations that our proposal for the beta prior hyper-parameter based on GRO also performs well in terms of power. Finally, we compare the power of our E-variable with this default prior choice and different restrictions on \mathcal{H}_1 to Fisher’s exact test.

REGROW For simplicity, in all our experiments we will invariably set the beta prior hyperparameters to $\alpha_a = \alpha_b = \beta_a = \beta_b = \gamma$ for some $\gamma > 0$ (recall that any such choice leads to a valid E-variable). We will aim for the γ that minimizes (1.7) in the worst-case over all $\theta_1^* \in [0, 1]^2$, thereby following the REGROW (*relative growth-rate optimality in worst-case*) criterion of Grünwald et al. [2019], who give a minimax regret motivation for this choice. In essence, the prior minimizing, among all distributions over $[0, 1]^2$, the maximum of (1.7) over all θ_1^* can be viewed as the prior that allows us to learn θ_1^* as fast as possible in the worst-case. Here we are contented to adopt a sub-optimal but computationally convenient

¹R code will be made publicly available; for now it is available on request.

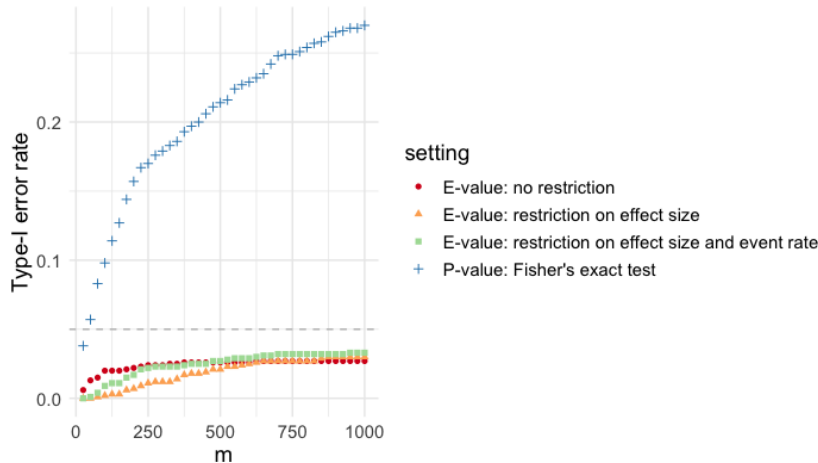


Figure 2: Type-I error rates for various E-variables and Fisher’s exact test under optional stopping estimated with 1000 simulations of two Bernoulli(0.1) data streams of length 1000, with $n_a = n_b = 1$. Significance level $\alpha = 0.05$ was used (grey dashed line). For the safe tests, beta prior parameter values used were $\gamma = \alpha_a = \beta_a = \alpha_b = \beta_b = 1/2$ ($\gamma = 0.18$ gave comparable results). For the E-variables with restrictions on \mathcal{H}_1 , we used $\delta = 0.05$ and $\theta_a = 0.1$.

prior by restrict the minimum to be over a 1-dimensional family of beta priors with hyper parameter γ . We find the minimizing γ by experiment (details omitted). It depends on m , which is unknown in advance, but for large m , in the setting with $n_a = n_b = 1$, it converges to $\gamma \approx 0.18$, and this is the value we will take as our default choice — our experiments below indicate that it remains a good choice, also when our main concern is power, and also under restrictions on \mathcal{H}_1 .

Type-I Error In Figure 2 type-I error rates of several E-variables and Fisher’s exact test estimated through a simulation experiment are depicted. 2000 samples of length 1000 were drawn according to a Bernoulli(0.1) distribution to represent 1000 data streams in two groups. After each complete block $m \in \{1, \dots, 1000\}$ an E-value or p-value was calculated and the proportion of rejected experiments up until m with each test type was recorded. As the stream lengths increase, the type-I error rate under (incorrectly applied) optional stopping with Fisher’s exact test increases quickly. The type-I error rate of the E-variables remains bounded.

Power Whereas GROwth is the natural performance measure in experiments that may always be continued at some point in the future, traditionally oriented researchers may be more interested in power. The question is then whether the optimal asymptotic choice $\gamma \approx 0.18$ in terms of the relative GRO property for unrestricted \mathcal{H}_1 is also the optimal choice in terms of power (which is usually considered in combination with some minimal effect size, i.e. a restricted \mathcal{H}_1). The following experiment shows that by and large it is.

For simplicity we only illustrate the case $n_a = n_b = 1$ and a desired power of 0.8. For various effect sizes δ , and various values of γ , we first determined the smallest sample size (number of blocks) m such that, under optional stopping up until and including m , the power is ≥ 0.8 in the worst case over all (θ_a, θ_b) with $\delta = \theta_b - \theta_a$. Here by ‘optional stopping up until and including m ’, we mean ‘we stop and reject the null iff $S_{[n_a, n_b, W_{[\gamma]}]}^{(m')} > \alpha^{-1}$ for some $m' \in \{1, 2, \dots, m\}$, and we stop and accept the null if this is not the case (so m is the maximal sample size we consider)’. We call this m the *worst-case* sample size needed for 80% power at effect size δ with prior parameter γ . The reason for calling it worst-case is that in practice, by engaging in optional stopping with a fixed maximal sample size, the *expected sample size* of this procedure is smaller: if, for $m' < m$, we already have $S_{[n_a, n_b, W_{[\gamma]}]}^{(m')} > \alpha^{-1}$ then we stop and reject early; if not, we go on until we have seen m blocks and then stop (and reject iff $S_{[n_a, n_b, W_{[\gamma]}]}^{(m)} > \alpha^{-1}$). We thus performed two simulation experiments: first, to estimate the worst-case sample size (at $\alpha = 0.05$), and second, to estimate the expected sample size. Again, the estimates were obtained by re-simulating a sequence of data blocks K times for a large number of K , making sure the bias and variance of the estimates were sufficiently small.

In Figure 3 results of these experiments are depicted. We make two observations: first, almost no difference in sample sizes to plan for between $\gamma = 0.18$ and $\gamma = 0.05$ was observed for distributions with small expected sample sizes (represented by the triangles and the dots, which overlap for most data points), and other values of γ obtained smaller power, indicating that the relative growth-optimal $\gamma = 0.18$ could in practice be used as a default setting for our E-variable — and as a consequence, we recommend it as such. Second, in the rightmost panel we see that for distributions with *very* small relative differences between θ_a and θ_b , e.g. $P_{0.5, 0.58}$, values of γ higher than 0.18 yielded a higher power, whereas for such δ , the relative GROW criterion was optimized for $\gamma = 0.18$ for the corresponding (very large) stopping times in our simulation experiments. This is not surprising given what is known for simple $\mathcal{H}_0 = \{P_{\theta_0}\}$: when testing a point null θ_0 with a 1-dimensional exponential family alternative, safe tests based on Bayes factors with standard Bayesian (e.g. Gaussian or conjugate) priors do not obtain optimal power in an asymptotic sense: they reject if $|\hat{\theta} - \theta_0|^2 \gtrsim (\log n)/n$ (with $\hat{\theta}$ denoting the MLE; see the example on Z-tests by Grünwald et al. [2019]) whereas based on nonstandard ‘switching’ [van der Pas and Grünwald, 2018] or ‘stitching’ methods [Howard et al., 2021], corresponding to special priors with densities going to infinity as effect size goes to 0, one can get rejection if $|\hat{\theta} - \theta_0|^2 \gtrsim (\log \log n)/n$. However, there is a significant price to pay in terms of the constants hidden in the asymptotics, and in practice, ‘standard’ priors may very well perform better at all but very large sample sizes [Maillard, 2019]. Given that the higher γ , the more the beta prior behaves like a switch prior, we conjecture that what we see in Figure 3(b) at very small δ is a version of the switching/stitching phenomenon with a composite null; since it only kicks in at very large sample sizes, we prefer $\gamma = 0.18$ as the default choice after all.

Finally, we compared the performance of our E-variables with the “default” beta priors with $\gamma = 0.18$ with their classical counterpart, Fisher’s exact test. Worst-case and expected stopping times of the E-variables with- and without restrictions on \mathcal{H}_1 were compared for

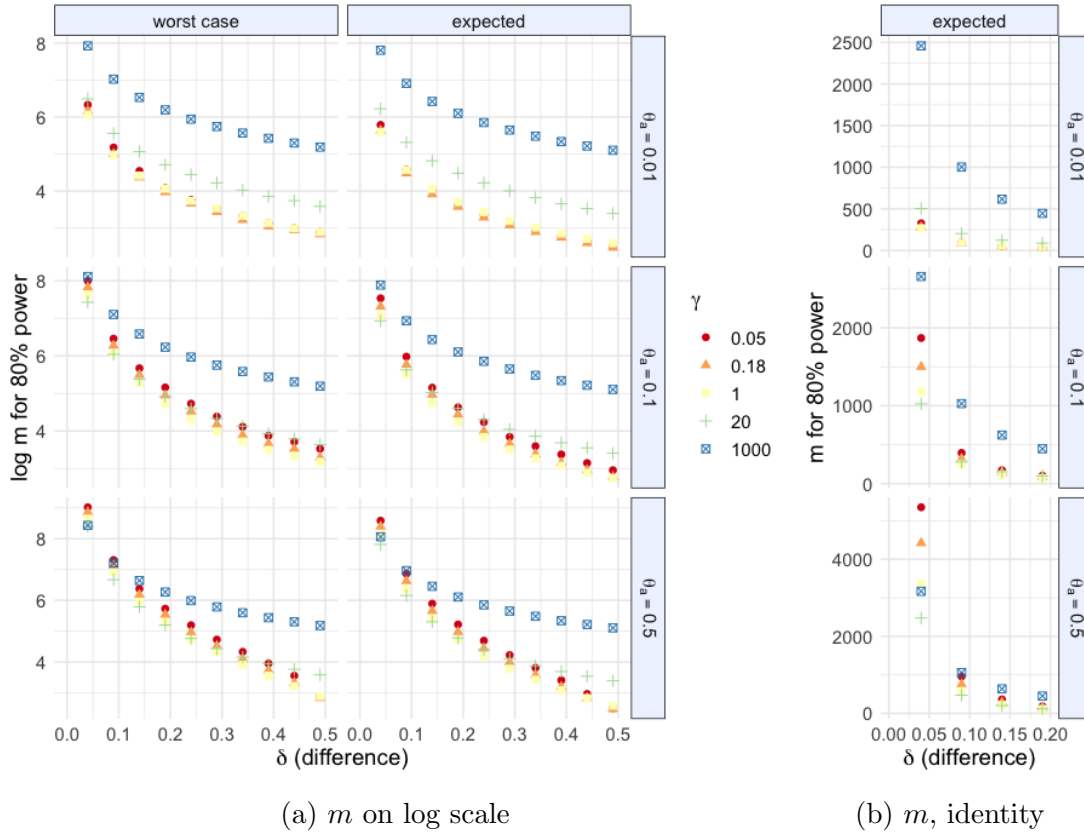


Figure 3: In 2000 simulations, the (natural logarithm of) the number of data blocks m (“sample sizes”) needed for achieving 80% power while testing at $\alpha = 0.05$ for distributions with varying group means and varying differences between group means were estimated for different beta prior parameter values.

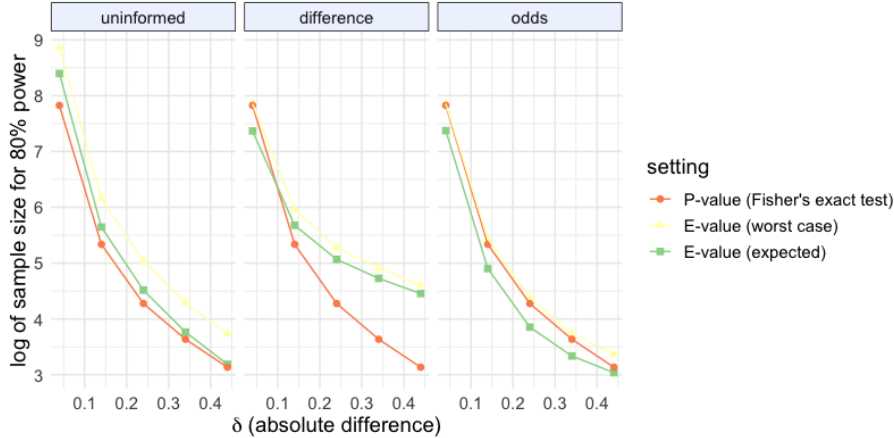


Figure 4: Estimates from 1000 simulations of worst-case and expected sample sizes for achieving 80% power estimated for three types of E-variables with different restrictions on \mathcal{H}_1 , and the sample size to plan for with Fisher’s exact test. Hypothesized effect sizes were 0.04 for the E-variables with prior information on the absolute difference and were converted equivalently for the log odds ratio prior information case, and we set $\gamma = 0.18$ for the beta priors.

sample sizes one would need to plan for when analyzing experiment results with Fisher’s exact test; see Figure 4. We noticed that the expected sample sizes achieved under optional stopping with the E-variable with unrestricted \mathcal{H}_1 were very similar to the sample sizes needed to plan for with Fisher’s exact test. When using a correctly specified restriction on \mathcal{H}_1 (the leftmost data points in the second and third subfigures), this expected number of samples is even considerably lower than the sample size to plan for with Fisher’s exact test. However, under misspecification, when the difference or log odds ratio used in the design of the E-variable turns out to be a lot smaller than the real difference present in the data generating machinery, one should expect to collect more samples (the data points towards the right in the second and third subfigures). This effect would disappear if were to put a prior on the full $\Theta^+(\delta)$ rather than the boundary $\Theta(\delta)$, at the price of slightly worse behaviour in the well-specified case when data is sampled from $\Theta(\delta)$.

6 Illustration via Real World Data

We will now demonstrate the approach through a real-world example: the SWEPIS study on labor induction [Wennerholm et al., 2019]. Wagenmakers and Ly have used this example before to illustrate how using single p-values to make decisions can hide valuable information in research data [Wagenmakers and Ly, 2020].

In the SWEPIS study, two groups of pregnant women were followed. In the first group labor was induced at 41 weeks, and in the second labor was induced after 42 weeks. The study was stopped early, as 6 cases of stillbirth were observed in the 42-weeks group (at

$n_b = 1379$), as compared to 0 in the 41-weeks group (at $n_a = 1381$). These data yield a significant Fisher’s exact test, $P \approx 0.015$, for testing that the number of stillbirths in the 42-weeks group is higher, when (wrongly) assuming that n_a and n_b were fixed in advance to the above values.

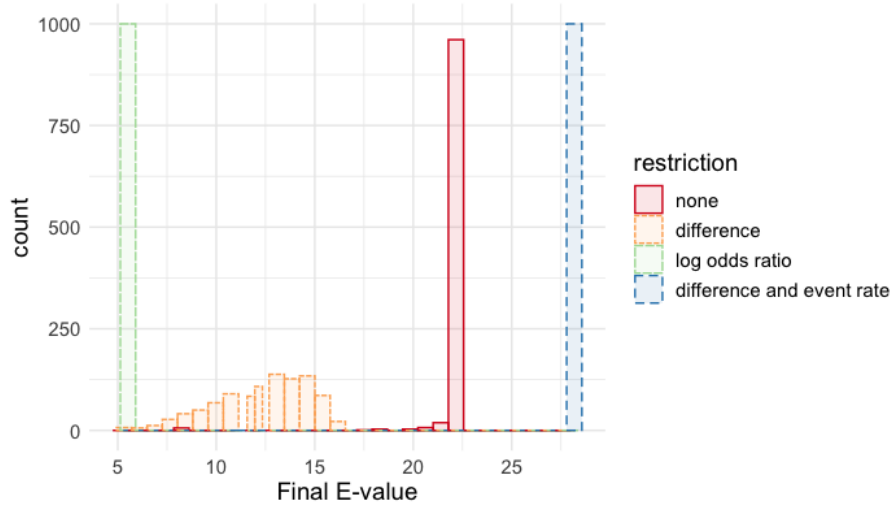
If we had used E-variables for continuously analyzing this data, would we then have found evidence for superiority of the 41 weeks approach, and would we have stopped the study earlier? As the E-variables we propose are not exchangeable, i.e. their values change under permutations of the data sequences, a direct comparison to the results of the SWEPIs study is not possible as the exact data stream is not available. To simulate a “real-time” scenario equivalent to the SWEPIs study, we assume we collect a total of 1380 data blocks, with $n_a = n_b = 1$, with a total of 2760 observations. We already know that in group a, 0 events are observed. In group b, 6 events are observed, of which we know that the last event was observed in data block 1380, directly before the study was stopped. Hence, we can simulate the “real-time” data by permuting the indices of the observations in group b in the 1379 first data blocks.

Four different approaches for analyzing the data with E-variables were explored: without any restriction on \mathcal{H}_1 , with a restriction based on the additive divergence measure (the minimal difference between the groups), with a restriction based on the log odds ratio, and with a restriction on the event rate in the control group *and* on the minimal difference. The minimal difference, log odds ratio and event rate used were chosen based on a large recent meta-analysis on stillbirths [Muglu et al., 2019]; we used $\delta = 0.00318$ as a restriction on the difference between the groups, $\log(2)$ for the log odds ratio and 0.0001 as the event rate. For all E-variables, the default beta prior hyperparameters with $\gamma = 0.18$ as earlier were used.

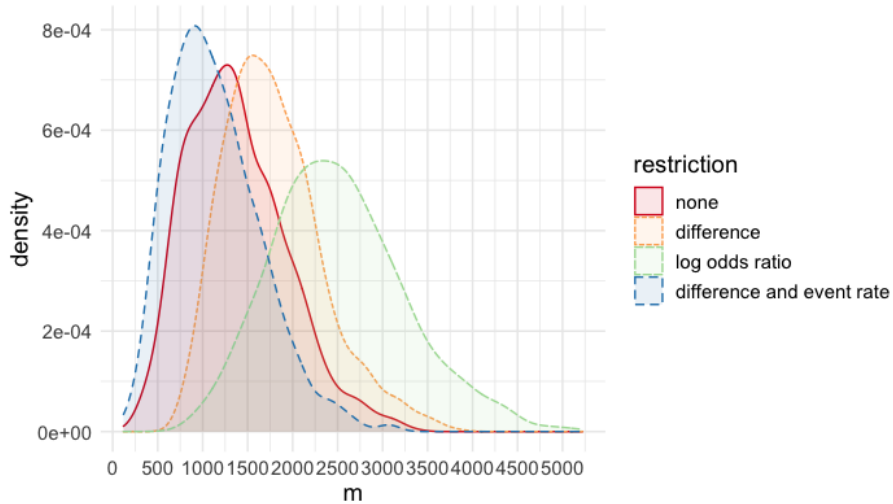
In Figure 5 the spread of the evidence collected with the four types of E-variables in 1000 simulations analogous to the SWEPIs setting is depicted. Because the observed effect size was higher than expected, E-values obtained with the (too low) restriction on the effect size were lower than the E-values obtained with the E-variable without restrictions. Adding the restriction on the event rate increased the E-values, and in all 1000 simulations, the SWEPIs study would have been stopped before the occurrence of the sixth stillbirth. Figure 5 also depicts results of a second simulation experiment, where we sampled 1000 data streams from $P_{0,6/1380}$ and recorded the stopping times while analyzing the streams with the four E-variables with different restrictions on \mathcal{H}_1 . With the E-variables without restriction, or with a restriction on the event rate and difference between the groups, we would have often stopped data collection earlier than in the SWEPIs setting.

We can thus conclude that, would the monitoring of the study have been performed with E-variables instead of p-values, first of all we would have collected *correct* evidence for a higher proportion of stillbirths in the 42-weeks group, and second, the degree of evidence is quite similar to that collected with the (incorrectly determined) p-value: both are significant at the 0.05 level. Wagemakers and Ly with their method also found evidence for the existence of a difference between the two groups, but not nearly of the same degree: they reported Bayes factors that varied, depending on the choice of the prior, between 1

and 5.4 (note that whenever we reject, our product of E-values, which like a Bayes factor can be thought of as a prequential likelihood ratio, must be ≥ 20). A possible explanation for this difference could be that the Bayes factors used for collecting evidence in their study are not designed for analyzing stream data. As we also saw in our experiments, choosing the wrong prior or restriction on \mathcal{H}_1 can make a large difference for the evidence collected. These results show that when planning a prospective study, using E-variables for analysis could, through their flexibility, contribute to earlier evidence collection compared to existing methods.



(a) Simulated E-values in SWEPIs setting, stopping at $m = 1380$ or when $E \geq 20$



(b) Simulated stopping times in setting with continuing until $E \geq 20$

Figure 5: Spread of E-values and stopping times observed with safe analysis of 1000 simulations of data streams analogous to the SWEPIs scenario, with four different types of restrictions on \mathcal{H}_1 .

7 Conclusion

We have established E-variables and test martingales for the general two-i.i.d.-data streams problem. We have demonstrated, using theory, simulations and a real-world example that, for tests of two proportions, by choosing an appropriate prior on Θ_1 , the method can be made competitive with classical methods that do not allow for optional stopping.

Whereas in this paper, we have focused on testing, our E-variables can also be extended to get *anytime-valid confidence sequences* [Howard et al., 2021, Lai, 1976], i.e. confidence sequences for effect sizes that are valid even under optional stopping. This requires us to first extend the testing to scenarios with $\delta \geq \delta_1$ vs. $\delta \leq \delta_0$ for $\delta_0 \neq 0$, that is, null hypotheses with $\theta_a \neq \theta_b$. We will report on this extension elsewhere. Our work also suggests a question for future work that is practically relevant, easy to state but hard to answer: to what extent do our findings generalize to logistic regression?

Appendices

Appendix A Proofs

A.1 Proof of Propositions

Proof of Proposition 1 We will actually prove Proposition 3 below, a generalization of Proposition 1 that will be useful when proving Theorem 1. Here we use the notation adopted later in the paper: for $\mathcal{H}_0 = \{P_\theta : \theta \in \Theta_0\}$ and, for W a distribution on Θ_0 , we write $P_W = \int P_\theta dW(\theta)$.

Proposition 3. *Let $\mathcal{H}_1 = \{Q\}$ be a singleton and let $\mathcal{H}_0 = \{P_\theta : \theta \in \Theta_0\}$. Define $S_{\theta,(j)} := q(Y_{(j)})/p_\theta(Y_{(j)})$ and $S_{W,(j)} = q(Y_{(j)})/p_W(Y_{(j)})$. Suppose there exists a distribution W on Θ_0 with finite support such that $S_{W,(1)}$ is an \mathbf{E} -variable. Then:*

1. $S_{W,(1)}$ is the Q -GRO \mathbf{E} -variable for $Y_{(1)}$.
2. A valid \mathbf{E} -variable of this form, with W putting mass 1 on a single $\theta^\circ \in \Theta_0$ so that $S_{W,(1)} = S_{\theta^\circ,(1)}$, automatically exists if \mathcal{H}_0 is a convex set of distributions that is compact in the weak topology.
3. If, for some $\theta^\circ \in \Theta_0$, $S_{\theta^\circ,(1)}$ is an \mathbf{E} -variable and we further assume that $Y_{(1)}, Y_{(2)}, \dots$ are i.i.d. according to all distributions in $\mathcal{H}_0 \cup \mathcal{H}_1$, then $S_{\text{GRO}(Q)}^{(m)} = \prod_{j=1}^m S_{\theta^\circ,(j)}$, i.e. the Q -GRO optimal (unconditional) \mathbf{E} -variable for $Y^{(m)}$ is the product of the individual Q -GRO optimal \mathbf{E} -variables.

Proof. Part 1 [Grünwald et al., 2019, Theorem 1] implies, for each $m \geq 1$, that (a) there can be at most one distribution W such that $S_{W,(1)}$ is an \mathbf{E} -variable, and (b), if such a W exists, then $S_{W,(1)}$ must be the Q -GRO \mathbf{E} -variable for $Y_{(1)}$. This implies the statement.

Part 2 Assume \mathcal{H}_0 is convex. Then for every distribution W on Θ_0 with finite support, we must have $S_{W,(1)} = S_{\theta,(1)}$ for some $\theta \in \Theta_0$ such that $p_\theta = \sum_{i=1}^k w(\theta_k) p_{\theta_k}$ where $\{\theta_1, \dots, \theta_k\}$ is the support of W and w is the corresponding probability mass function.

Let $D(Q(Y_{(1)}) \| P_\theta(Y_{(1)}))$ be the KL divergence between the marginal distributions for $Y_{(1)}$ according to Q and P_θ . We claim that there exists $\theta^\circ \in \Theta_0$ such that

$$D(Q(Y_{(1)}) \| P_{\theta^\circ}(Y_{(1)})) = \inf_{\theta \in \Theta_0} D(Q(Y_{(1)}) \| P_\theta(Y_{(1)})) = \inf_W D(Q(Y_{(1)}) \| P_W(Y_{(1)})), \quad (\text{A.1})$$

where the second infimum is over all distributions W with finite support. Here the first equality (stating that the minimum is achieved) follows from Posner [1975] who showed that the KL divergence is lower semi-continuous in its second argument, together with compactness of \mathcal{H}_0 . The second equality follows by convexity of \mathcal{H}_0 . But (A.1) expresses that P_{θ° is the *reverse information projection (RIPr)* [Li, 1999, Li and Barron, 2000, Grünwald et al., 2019] of Q onto the convex hull of \mathcal{H}_0 (restricted to single outcomes

$Y_{(1)}$). [Grünwald et al., 2019, Theorem 1] then immediately gives that $q(Y_{(1)})/p_{\theta^\circ}(Y_{(1)})$ is an \mathbf{E} -variable.

Part 3 The assumption implies that $S_{\theta^\circ, (1)}$ is an \mathbf{E} -variable. Moreover, the i.i.d. assumption implies that $S_{\theta^\circ}^{(m)} := \prod_{j=1}^m S_{\theta^\circ, (j)} = \prod q(Y_{(j)})/p_{\theta^\circ}(Y_{(j)})$ is also an \mathbf{E} -variable. But [Grünwald et al., 2019, Theorem 1] implies, for \mathcal{H}_0 for which data are i.i.d., for each $m \geq 1$, that (a) if a $\theta \in \Theta_0$ exists such that $S_\theta^{(m)}$ is an \mathbf{E} -variable, then $S_\theta^{(m)}$ must be the Q -GRO \mathbf{E} -variable for $Y^{(m)}$. This proves the statement. \square

Proof of Proposition 2 The formulae for $\check{\theta}_a|Y^{(j-1)}$ and $\check{\theta}_b|Y^{(j-1)}$ are standard expressions for the Bayes predictive distribution based on the given beta priors; we omit further details. As to the expression for $\check{\theta}_0|Y^{(j-1)}$ in terms of $\kappa = n_b/n_a$: Straightforward rewriting gives, for general $\alpha_a, \alpha_b, \beta_a, \beta_b$:

$$\check{\theta}_0|Y^{(j-1)} = \frac{1}{1 + \kappa} \check{\theta}_a|Y^{(j-1)} + \frac{\kappa}{1 + \kappa} \check{\theta}_b|Y^{(j-1)}.$$

If we plug in the expressions for $\check{\theta}_a|Y^{(j-1)}, \check{\theta}_b|Y^{(j-1)}$ and we instantiate to $\alpha_b = \kappa\alpha_a$, and $\beta_b = \kappa\beta_a$, this becomes

$$\begin{aligned} \check{\theta}_0|Y^{(j-1)} &= \frac{1}{1 + \kappa} \frac{U_a + \alpha_a}{n_a(j-1) + \alpha_a + \beta_a} + \frac{\kappa}{1 + \kappa} \frac{U_b + \alpha_b}{\kappa(n_a(j-1) + \alpha_a + \beta_a)} \\ &= \frac{1}{1 + \kappa} \frac{U_a + U_b + (1 + \kappa)\alpha_a}{n_a(j-1) + \alpha_a + \beta_a} = \frac{U + (1 + \kappa)\alpha_a}{n(j-1) + (1 + \kappa)\alpha_a + (1 + \kappa)\beta_a}, \end{aligned}$$

which is what we had to prove.

A.2 Proof of Theorem 1

The following fact plays a central role in the proof:

Fact 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 $u^{n_a} v^{n_b} \leq 1$.

This result 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 .

Proof of Theorem 1 *Part 1* For $y \in \mathcal{Y}$, set $p^\circ(y) := (n_a/n)p_{\theta_a^*}(y) + (n_b/n)p_{\theta_b^*}(y)$ and $p^\circ(y^k) = \prod_{i=1}^k p^\circ(y_i)$. For all $\theta \in \Theta$ we have:

$$\begin{aligned} \mathbf{E}_{Y_a^{n_a}, Y_b^{n_b} \text{ i.i.d. } \sim P_\theta} [s(V^n; n_a, n_b, \theta_a^*, \theta_b^*)] &= \mathbf{E}_{Y_a^{n_a} \sim P_\theta} \left[\frac{p_{\theta_a^*}(Y_a^{n_a})}{p^\circ(Y_a^{n_a})} \right] \cdot \mathbf{E}_{Y_b^{n_b} \sim P_\theta} \left[\frac{p_{\theta_b^*}(Y_b^{n_b})}{p^\circ(Y_b^{n_b})} \right] = \\ &= \left(\mathbf{E}_{Y \sim P_\theta} \left[\frac{p_{\theta_a^*}(Y)}{p^\circ(Y)} \right] \right)^{n_a} \cdot \left(\mathbf{E}_{Y \sim P_\theta} \left[\frac{p_{\theta_b^*}(Y)}{p^\circ(Y)} \right] \right)^{n_b}. \end{aligned} \quad (\text{A.2})$$

We also have

$$\begin{aligned} & \frac{n_a}{n} \mathbf{E}_{Y \sim P_\theta} \left[\frac{p_{\theta_a^*}(Y)}{p^\circ(Y)} \right] + \frac{n_b}{n} \mathbf{E}_{Y \sim P_\theta} \left[\frac{p_{\theta_b^*}(Y)}{p^\circ(Y)} \right] \\ &= \mathbf{E}_{Y \sim P_\theta} \left[\frac{n_a}{n} \cdot \frac{p_{\theta_a^*}(Y)}{\frac{n_a}{n} p_{\theta_a^*}(Y) + \frac{n_b}{n} p_{\theta_b^*}(Y)} + \frac{n_b}{n} \cdot \frac{p_{\theta_b^*}(Y)}{\frac{n_b}{n} p_{\theta_a^*}(Y) + \frac{n_b}{n} p_{\theta_b^*}(Y)} \right] = 1. \end{aligned} \quad (\text{A.3})$$

The result now follows by combining (A.2) with (A.3) using the Fact further above.

Part 2 By convexity of $\{P_\theta : \theta \in \Theta\}$, there exists $\theta^\circ \in \Theta$ such that $p_{\theta^\circ} = (n_a/n)p_{\theta_a^*} + (n_b/n)p_{\theta_b^*}$ and then the numerator in (2.2) can be rewritten as $p_{\theta^\circ}(y_a^{n_a}, y_b^{n_b})$. The GRO-property is now an immediate consequence of Proposition 3, Part 3.

Appendix B Numerical approach to calculating E-variables for restricted \mathcal{H}_1

In this subsection we describe how we propose to approximate the beta prior and posterior on the restricted \mathcal{H}_1 with parameter space $\Theta(\delta)$, as defined in (4.1). Note that we limit ourselves to $\delta > 0$ in this detailed description; for $\delta < 0$ one can apply an entirely equivalent approach, with an extra term in the reparameterization. We define

$$\zeta = \begin{cases} \delta & \text{if } d((\theta_a, \theta_b)) = \theta_b - \theta_a, \\ 0 & \text{if } d((\theta_a, \theta_b)) = \log\text{-odds-ratio}(\theta_a, \theta_b), \end{cases} \quad (\text{B.1})$$

such that we have $\theta_a \in (0, 1 - \zeta)$ and in both cases, θ_b is completely determined by θ_a : $\theta_b = d^{-1}(\delta; \theta_a)$. Hence, our density estimation problem now becomes one-dimensional, which enables us to put a discretized prior on the restricted parameter space.

First, we discretize the parameter space Θ_a to a grid (a vector) with precision K , $K \in (0, 1 - \zeta)$ and $1/K \in \mathbb{N}^+$: $\bar{\theta}_a = (K, 2K, 3K, \dots, 1 - \zeta)$. Then, we reparameterize $\theta_a = (1 - \zeta)\rho$, with $\rho \in (0, 1)$. Then, we have $\bar{\rho} = (K/(1 - \zeta), 2K/(1 - \zeta), \dots, 1)$. For the discretized grid $\bar{\rho}$, we compute the prior $W = \text{Beta}(\alpha, \beta)$ densities and normalize them, which also gives us the discretized densities for each $\theta_a^i \in \bar{\theta}_a$ (with $i \in (1, 2, \dots, 1/K)$):

$$\pi_{\alpha, \beta, \zeta}(\theta_a^i) = \frac{\text{Beta}(\frac{\theta_a^i}{1 - \zeta}; \alpha, \beta)}{\sum_{k=1}^{\frac{1}{K}} \text{Beta}(\frac{\theta_a^k}{1 - \zeta}; \alpha, \beta)}.$$

For all elements of $\bar{\theta}_a$, the corresponding θ_b is retrieved and the likelihood of incoming data points $p_{\theta_a, \theta_b}(Y^{(j-1)})$ is calculated. We can then estimate the posterior density of $\theta_a^i \in \bar{\theta}_a$:

$$p(\theta_a^i | Y^{(j-1)}) = \frac{\pi_{\alpha, \beta, \zeta}(\theta_a^i) p_{\theta_a^i, \theta_b^i}(Y^{(j-1)})}{\sum_{k=1}^{\frac{1}{K}} \pi_{\alpha, \beta, \zeta}(\theta_a^k) p_{\theta_a^k, \theta_b^k}(Y^{(j-1)})}.$$

We can then estimate $\check{\theta}_a | Y^{(j-1)} = \mathbf{E}_{\theta_a \sim W | Y^{(j-1)}}[\theta_a]$ as $\sum_{i=1}^{\frac{1}{K}} p(\theta_a^i | Y^{(j-1)}) \theta_a^i$, and $\check{\theta}_b | Y^{(j-1)} = d^{-1}(\delta; \theta_a | Y^{(j-1)})$.

Appendix C The Gunel-Dickey Bayes Factors do not give rise to E-variables

Sampling scheme	Fixed parameters	Bayes factor (10) for 2x2 table
Poisson	none	$\frac{8(n+1)(n_1+1)}{(n+4)(n+2)} \left[\frac{n_{a1}!n_{b1}!n_{a0}!n_{b0}!n!}{(n_1+1)!n_0!n_a!n_b!} \right]$
Joint multinomial	n	$\frac{6(n+1)(n_1+1)}{(n+3)(n+2)} \left[\frac{n_{a1}!n_{b1}!n_{a0}!n_{b0}!n!}{(n_1+1)!n_0!n_a!n_b!} \right]$
Independent multinomial	n_a, n_b	$\frac{\binom{n}{n_1} (n+1)}{\binom{n_a}{n_{a1}} \binom{n_b}{n_{b1}} (n_a+1)(n_b+1)}$
Hypergeometric	n_a, n_b, n_1	$\frac{n_{a1}!n_{b1}!n_{a0}!n_{b0}!n!}{\prod_{i \in \{a,b,0,1\}} (n_i + \mathbb{1}_{n_i = \min(n_a, n_b, n_0, n_1)})!}$

Table 1: Overview of (objective) Bayes factors for contingency table testing provided by Gunel and Dickey [1974] and Jamil et al. [2017].

We will not consider the hypergeometric and joint multinomial scenarios for this paper, where the number of successes n_1 is fixed, as they do not match the block-wise data design in this paper. The Bayes factor for the Poisson sampling scheme is not an E-variable, as the expectation under the null hypothesis with Poisson distributions on individual cell counts exceeds 1 for rates $\lambda \geq 1$:

$$\begin{aligned} & \mathbb{E}_{n_{rc} \sim \text{Poisson}(\lambda_{rc})} [BF_{10}(N_{a1}, N_{b1}, N_{a0}, N_{b0})] = \\ & \sum_{n_{a1}=0}^{\infty} \dots \sum_{n_{b0}=0}^{\infty} \pi_{\lambda_{a1}}(n_{a1}) \dots \pi_{\lambda_{b0}}(n_{b0}) BF_{10}(n_{a1}, n_{b1}, n_{a0}, n_{b0}) = \\ & \frac{8}{\exp(\lambda_{a1} + \dots + \lambda_{b0})} \sum_{n_{a1}=0}^{\infty} \dots \sum_{n_{b0}=0}^{\infty} \lambda_{a1}^{n_{a1}} \dots \lambda_{b0}^{n_{b0}} \frac{(n+1)(n_1+1)}{(n+4)(n+2)} \frac{n!}{(n_1+1)!n_0!n_a!n_b!}, \end{aligned}$$

as illustrated numerically in Figure 6 for increasing limits for the sums $\sum_{n_{rc}=1}^{\max n_{rc}}$. For the independent multinomial sampling scheme, let, without loss of generality, $n_a < n_b$. We get, with $n_0 = n - n_1$,

$$\begin{aligned} & \mathbb{E}_{N_{a1}, N_{b1} \sim \text{Binomial}(\theta)} [BF_{10}(N_{a1}, N_{b1} | n_a, n_b)] = \\ & \sum_{n_{a1}=0}^{n_a} \sum_{n_{b1}=0}^{n_b} \binom{n_a}{n_{a1}} \binom{n_b}{n_{b1}} \theta^{n_1} (1-\theta)^{n_0} \frac{\binom{n}{n_1}}{\binom{n_a}{n_{a1}} \binom{n_b}{n_{b1}}} \frac{(n+1)}{(n_a+1)(n_b+1)} = \\ & \frac{(n+1)}{(n_a+1)(n_b+1)} \sum_{n_{a1}=0}^{n_a} \sum_{n_{b1}=0}^{n_b} \binom{n}{n_1} \theta^{n_1} (1-\theta)^{n_0} \end{aligned}$$

Numerical simulations show that, for a range of choices for n, n_a and θ this exceeds 1; see Figure 7.

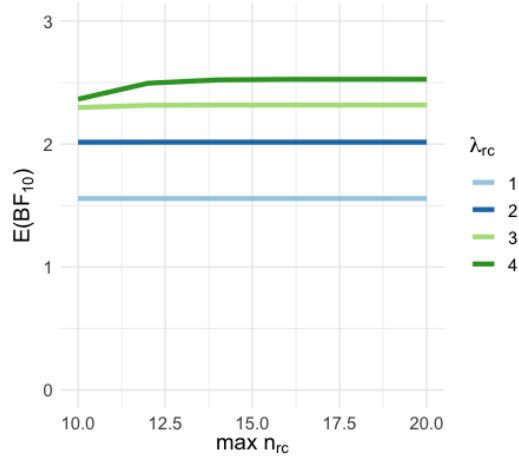


Figure 6: The Gunel-Dickey Bayes factor for the Poisson sampling scheme is not an E-variable: $\sum_{n_{a1}=0}^{\max n_{rc}} \dots \sum_{n_{b0}=0}^{\max n_{rc}} \pi_{\lambda_{a1}}(n_{a1}) \dots \pi_{\lambda_{b0}}(n_{b0}) BF_{10}(n_{a1}, n_{b1}, n_{a0}, n_{b0})$ for various $\max n_{rc}$ and λ_{rc} .

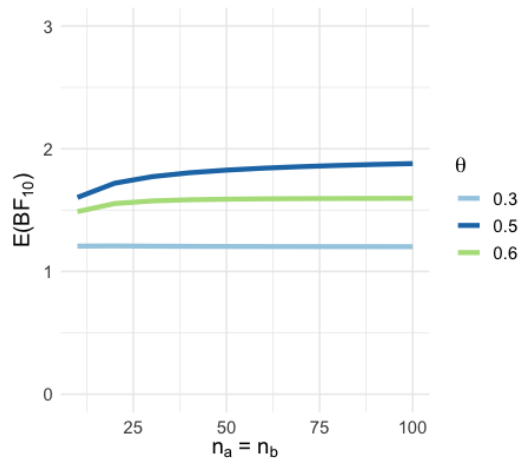


Figure 7: The Gunel-Dickey Bayes factor for the independent multinomial sampling scheme is not an E-variable: $\mathbb{E}_{N_{a1}, N_{b1} \sim \text{Binomial}(\theta)} [BF_{10}(N_{a1}, N_{b1} | n_a, n_b)]$ for various choices of θ and n_g .

References

- George A Barnard. Sequential tests in industrial statistics. *Supplement to the Journal of the Royal Statistical Society*, 8(1):1–21, 1946.
- Leo Breiman. Optimal gambling systems for favorable games. *Fourth Berkeley Symposium*, 1961.
- D.A. Darling and H. Robbins. Confidence sequences for mean, variance, and median. *Proceedings of the National Academy of Sciences of the United States of America*, 58(1):66, 1967.
- A.P. Dawid. Present position and potential developments: Some personal views, statistical theory, the prequential approach. *Journal of the Royal Statistical Society, Series A*, 147(2):278–292, 1984.
- Peter Grünwald, Rianne de Heide, and Wouter Koolen. Safe testing. *arXiv preprint arXiv:1906.07801*, 2019.
- E. Gunel and J. Dickey. Bayes factors for independence in contingency tables. *Biometrika*, 61(3):545–557, 1974.
- Alexander Henzi and Johanna F. Ziegel. Valid sequential inference on probability forecast performance. *arXiv preprint arXiv:2103.08402*, 2021.
- Steven R Howard, Aaditya Ramdas, Jon McAuliffe, and Jasjeet Sekhon. Uniform, nonparametric, non-asymptotic confidence sequences. *Annals of Statistics*, 2021.
- Tahira Jamil, Alexander Ly, Richard D Morey, Jonathon Love, Maarten Marsman, and Eric-Jan Wagenmakers. Default “Gunel and Dickey” Bayes factors for contingency tables. *Behavior Research Methods*, 49(2):638–652, 2017.
- Leslie K John, George Loewenstein, and Drazen Prelec. Measuring the prevalence of questionable research practices with incentives for truth telling. *Psychological science*, 23(5):524–532, 2012.
- Robert E Kass and Suresh K Vaidyanathan. Approximate Bayes factors and orthogonal parameters, with application to testing equality of two binomial proportions. *Journal of the Royal Statistical Society: Series B (Methodological)*, 54(1):129–144, 1992.
- J.L. Kelly. A new interpretation of information rate. *Bell System Technical Journal*, pages 917–926, 1956.
- Tze Leung Lai. On confidence sequences. *The Annals of Statistics*, 4(2):265–280, 1976.
- Alix Lhéritier and Frédéric Cazals. A sequential non-parametric multivariate two-sample test. *IEEE Transactions on Information Theory*, 64(5):3361–3370, 2018.

- J.Q. Li. *Estimation of Mixture Models*. PhD thesis, Yale University, New Haven, CT, 1999.
- J.Q. Li and A.R. Barron. Mixture density estimation. In S.A. Solla, T.K. Leen, and K-R. Müller, editors, *Advances in Neural Information Processing Systems*, volume 12, pages 279–285, Cambridge, MA, 2000. MIT Press.
- Odalric-Ambrym Maillard. Mathematics of statistical sequential decision making, 2019. Thèse de Habilitation.
- Tudor Manole and Aaditya Ramdas. Sequential estimation of convex divergences using reverse submartingales and exchangeable filtrations. *arXiv preprint arXiv:2103.09267*, 2021.
- Javaid Muglu, Henna Rather, David Arroyo-Manzano, Sohinee Bhattacharya, Imelda Balchin, Asma Khalil, Basky Thilaganathan, Khalid S Khan, Javier Zamora, and Shakila Thangaratinam. Risks of stillbirth and neonatal death with advancing gestation at term: A systematic review and meta-analysis of cohort studies of 15 million pregnancies. *PLoS medicine*, 16(7):e1002838, 2019.
- Luigi Pace and Alessandra Salvan. Likelihood, replicability and Robbins’ confidence sequences. *International Statistical Review*, 2019.
- Edward Posner. Random coding strategies for minimum entropy. *IEEE Transactions on Information Theory*, 21(4):388–391, 1975.
- Aaditya Ramdas, Johannes Ruf, Martin Larsson, and Wouter Koolen. Admissible anytime-valid sequential inference must rely on nonnegative martingales. *arXiv preprint arXiv:2009.03167*, 2020.
- Herbert Robbins. Statistical methods related to the law of the iterated logarithm. *The Annals of Mathematical Statistics*, 41(5):1397–1409, 1970.
- Glenn Shafer. The language of betting as a strategy for statistical and scientific communication. *Journal of the Royal Statistical Society, Series A*, 2021.
- Glenn Shafer, Alexander Shen, Nikolai Vereshchagin, and Vladimir Vovk. Test martingales, Bayes factors and p-values. *Statistical Science*, pages 84–101, 2011.
- David Siegmund. *Sequential analysis: tests and confidence intervals*. Springer Science & Business Media, 2013.
- J. ter Schure, M. F. Perez-Ortiz, A. Ly, and P. Grünwald. The safe log rank test: Error control under continuous monitoring with unlimited horizon. *arXiv preprint arXiv:1906.07801*, 2021.
- Rosanne J Turner. Safe tests for 2 x 2 contingency tables and the Cochran-Mantel-Haenszel test. Master’s thesis, Leiden University, 2019.

- S. van der Pas and P. Grünwald. Almost the best of three worlds: Risk, consistency and optional stopping for the switch criterion in nested model selection. *Statistica Sinica*, 28(1):229–255, 2018.
- Vladimir Vovk and Ruodu Wang. E-values: Calibration, combination, and applications. *Annals of Statistics*, 2021.
- Eric-Jan Wagenmakers and Alexander Ly. Bayesian scepticism about swepis: Quantifying the evidence that early induction of labour prevents perinatal deaths, 2020. URL psyarxiv.com/5ydpb.
- Ulla-Britt Wennerholm, Sissel Saltvedt, Anna Wessberg, Mårten Alkmark, Christina Bergh, Sophia Brismar Wendel, Helena Fadl, Maria Jonsson, Lars Ladfors, Verena Sengpiel, et al. Induction of labour at 41 weeks versus expectant management and induction of labour at 42 weeks (SWEdish Post-term Induction Study, swepis): multicentre, open label, randomised, superiority trial. *British Medical Journal*, 367, 2019.