

ROBUSTNESS OF NULL HYPOTHESIS BAYESIAN TESTING UNDER OPTIONAL STOPPING

Jorge N. Tendeiro
IMPS 2020 – July 14

University of Groningen

Joint work with Henk Kiers and Don van Ravenzwaaij.

AKA:

Sequential testing.

Definition:

Continuously testing a null hypothesis (\mathcal{H}_0) as data are collected until \mathcal{H}_0 is rejected.

Procedure:

- 1 Collect some data.
- 2 Perform the test (α and n_{\max} chosen in advance):
Compute p and...
 - ...if $p < \alpha$: STOP and retain \mathcal{H}_1 .
 - ...if $p > \alpha$: Back to 1.
- 3 Continue until either conclusive evidence is found or n_{\max} is reached.

Known for a long time to be *very* problematic:

- Based on null hypothesis significance testing (NHST).
- NHST has a **lot** of problems.¹
- In particular^{2,3} : **Too high** proportions of false positives ($\gg \alpha$).

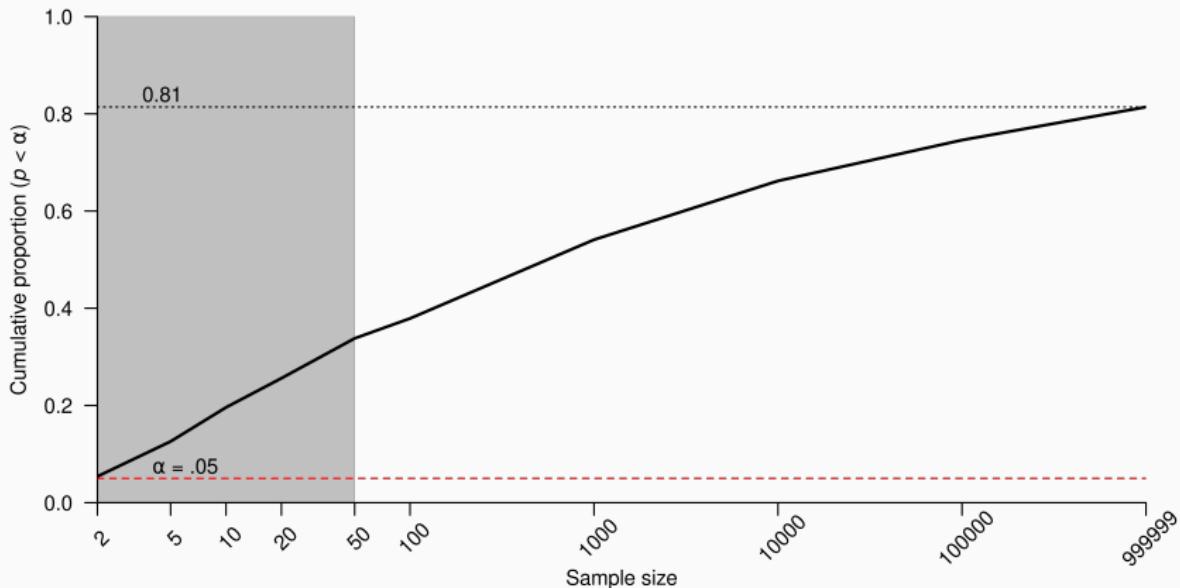
¹Wasserstein, Schirm, and Lazar (2019).

²Armitage, McPherson, and Rowe (1969).

³Jennison and Turnbull (1990).

Example:

- One-sample t -test: $\mathcal{H}_0 : \mu = 0$ vs $\mathcal{H}_1 : \mu \neq 0$.
- Repeat 1,000 times:
 - Sampling plan: $n = 2(1)1,000,000$ from $\mathcal{N}(0, 1)$.
 - Stop if $p < \alpha = .05$.



Some ways to avoid this problem:

- Using corrections^{1,2,3,4,5,6,7,8}.
Not commonly used in psychology.
- Not using optional stopping
(i.e., fixed sample size, sample until completion).
- Turning to the Bayesian paradigm.

¹Armitage (1960).

²Botella et al. (2006).

³Fitts (2010).

⁴Frick (1998).

⁵Jennison and Turnbull (1999).

⁶Lakens (2014).

⁷Pocock (1983).

⁸Wald (1945).

NHBT^{1,2,3,4} is the Bayesian counterpart to NHST.

It uses the **Bayes factor** in place of the *p*-value.

Definition 1:

The Bayes factor quantifies the *update* in our relative belief about the likelihood of two hypotheses (\mathcal{H}_0 , \mathcal{H}_1) in light of the observed data (D):

$$\underbrace{\frac{p(\mathcal{H}_1)}{p(\mathcal{H}_0)}}_{\text{prior odds}} \times \underbrace{\frac{p(D|\mathcal{H}_1)}{p(D|\mathcal{H}_0)}}_{BF_{10}} = \underbrace{\frac{p(\mathcal{H}_1|D)}{p(\mathcal{H}_0|D)}}_{\text{posterior odds}}$$

Definition 2:

The Bayes factor indicates the relative predictive value of each model.

- E.g., if the observed data are better predicted under \mathcal{H}_1 than under \mathcal{H}_0 then $p(D|\mathcal{H}_1) > p(D|\mathcal{H}_0)$ and so $BF_{10} > 1$.

¹Jeffreys (1961).

²Kass and Raftery (1995).

³Tendeiro and Kiers (2019).

⁴van de Schoot et al. (2017).

Procedure:^{1,2,3}

- 1 Collect some data.
- 2 Perform the test (BF_L , BF_U , and n_{\max} chosen in advance):
Compute BF_{10} and...
 - ...if $BF_{10} < BF_L$: Stop and retain \mathcal{H}_0 .
 - ...if $BF_{10} > BF_U$: Stop and retain \mathcal{H}_1 .
 - ...if $BF_L < BF_{10} < BF_U$: Back to 1.
- 3 Continue until either conclusive evidence is found or n_{\max} is reached.

One *major* improvement of Bayesian over frequentist optional stopping:

The Bayesian procedure can stop due to sufficiently strong evidence in favor of \mathcal{H}_0 .

¹Lindley (1957).

²Edwards, Lindman, and Savage (1963).

³Kass and Raftery (1995).

- It has been argued through the years that optional stopping under the Bayesian paradigm is allowed.^{1,2,3,4,5}
- It has even been further developed and used in practice.^{6,7,8,9,10}
- However, two recent papers disputed this state of affairs^{11,12} (also¹³).

Rouder offered a rebuttal to these ideas in 2014.

Title: 'Optional stopping: No problem for Bayesians'.¹⁴

¹Edwards, Lindman, and Savage (1963).

²Kass and Raftery (1995).

³Wagenmakers (2007).

⁴Wagenmakers et al. (2010).

⁵Francis (2012).

⁶Matzke et al. (2015).

⁷Schönbrodt et al. (2017).

⁸Schönbrodt and Wagenmakers (2018).

⁹Wagenmakers et al. (2012).

¹⁰Wagenmakers et al. (2015).

¹¹Yu et al. (2014).

¹²Sanborn and Hills (2014).

¹³de Heide and Grünwald (2017).

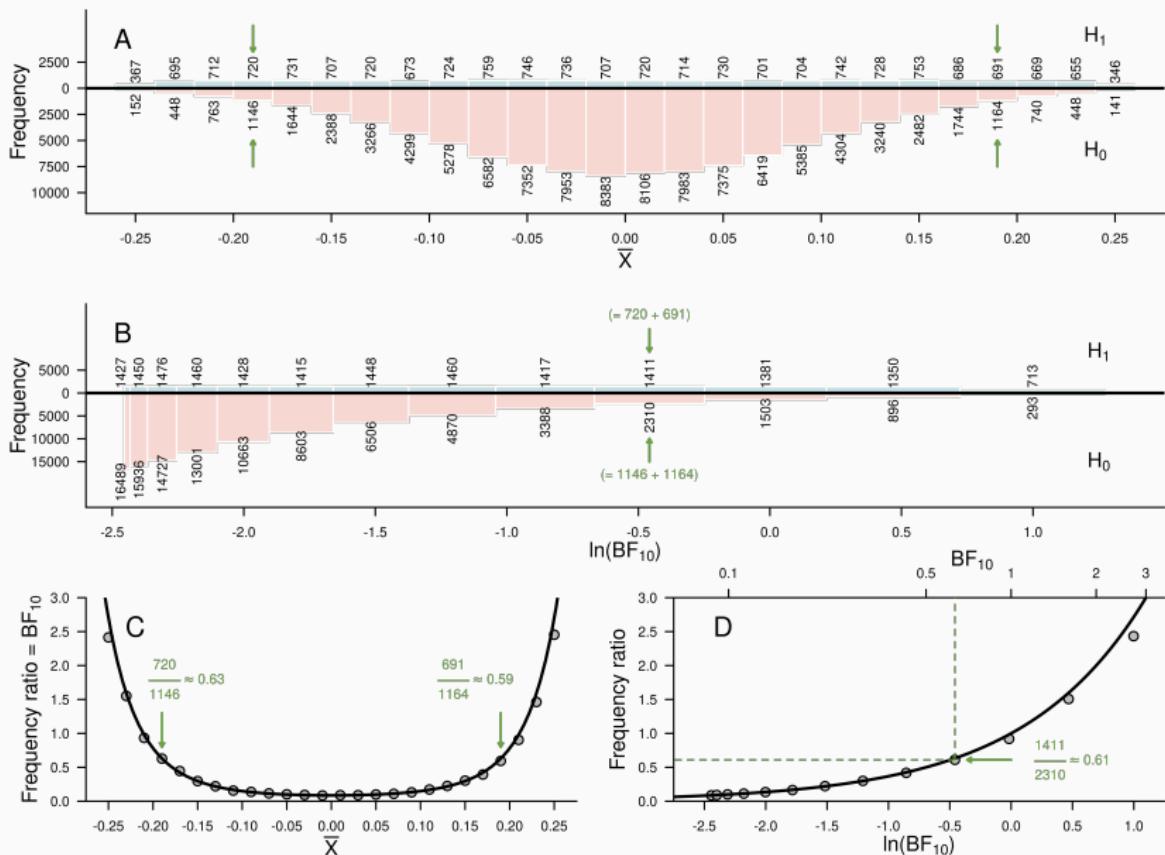
¹⁴Rouder (2014).

Yu et al. (2014) and Sanborn & Hills (2014) questioned the *long run properties* of the Bayesian optional stopping procedure.

Rouder (2014) argued that there was no problem *in a particular sense*.

Let's **visualize** the argument:

- Data: $X_i \sim \mathcal{N}(\mu, \sigma^2)$, for $i = 1, \dots, n$ and σ known.
- $\mathcal{H}_0 : \mu = 0$.
- $\mathcal{H}_1 : \mu \sim \mathcal{N}(0, \sigma_1^2)$, for σ_1 known.



Rouder claimed that Bayes factors are well calibrated under optional stopping.

The argument goes as follows:

- Assume prior odds equal to 1, so:

$$\underbrace{\frac{p(D|\mathcal{H}_1)}{p(D|\mathcal{H}_0)}}_{BF_{10}} = \underbrace{\frac{p(\mathcal{H}_1|D)}{p(\mathcal{H}_0|D)}}_{\text{posterior odds}}.$$

- By definition of posterior odds:

\mathcal{H}_1 is BF_{10} times more likely than \mathcal{H}_0 after considering the data.

Rouder made two assertions:

- 1 Assertion 1: For any given value BF_{10} ,

\mathcal{H}_1 is BF_{10} times more likely than \mathcal{H}_0 to have generated BF_{10} .

- 2 Assertion 2: The above statement also holds under optional stopping.

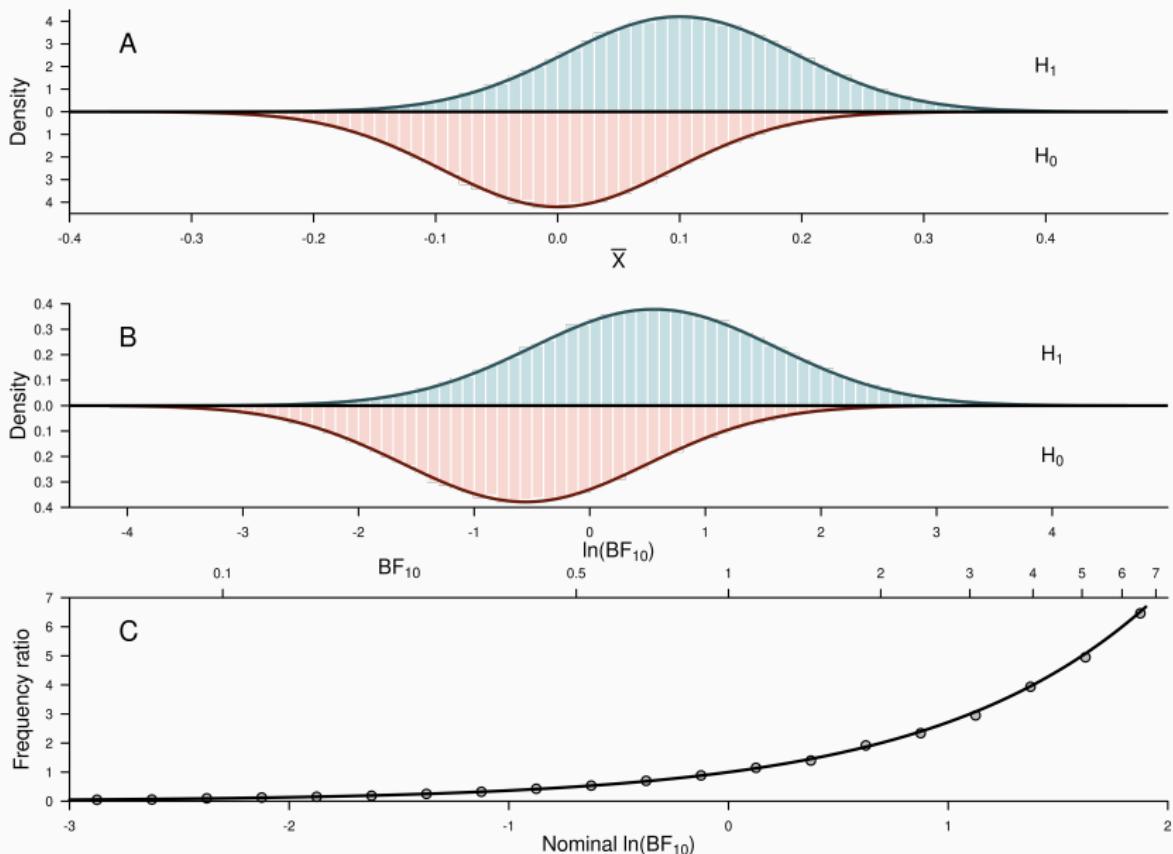
In our paper, we:

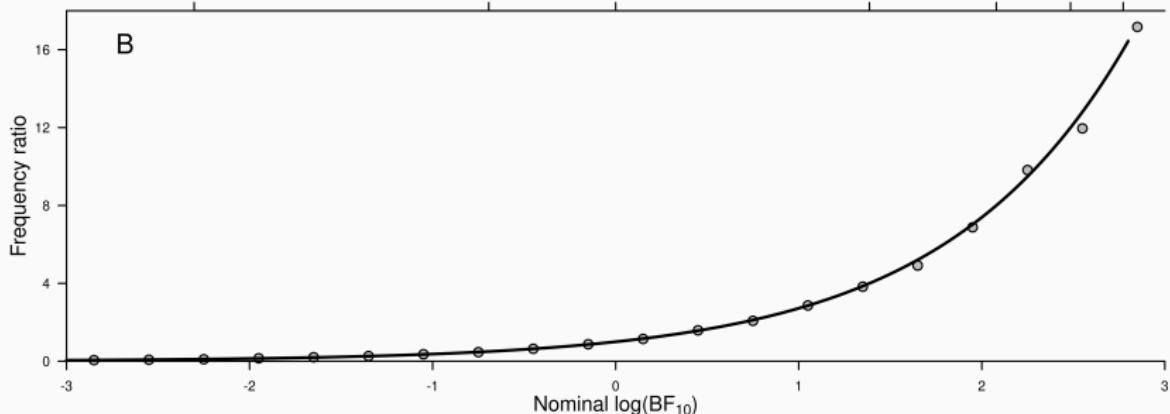
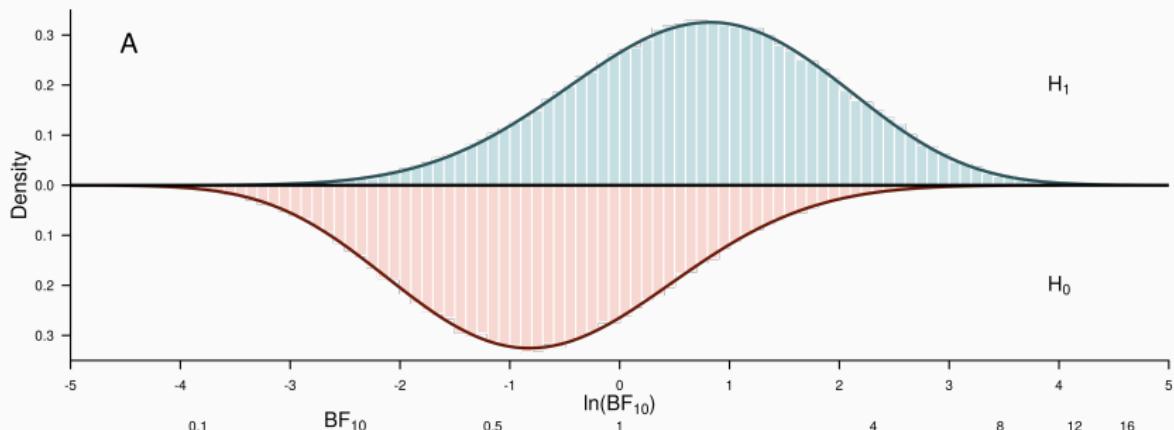
- Considered the same two tests as Rouder (2014):
 - Both tests about the mean of a normal distribution $\mathcal{N}(\mu, \sigma^2)$, σ known.
 - First test:
 $\mathcal{H}_0 : \mu = 0$ versus $\mathcal{H}_1 : \mu = \mu_1$.
 - Second test:
 $\mathcal{H}_0 : \mu = 0$ versus $\mathcal{H}_1 : \mu \sim \mathcal{N}(0, \sigma_1^2)$, σ_1 known.
- Derived **exact** probability distributions for BF_{10} .
- Proved Assertion 1 for n fixed.
- Proved Assertion 2 after one step of the optional stopping procedure.

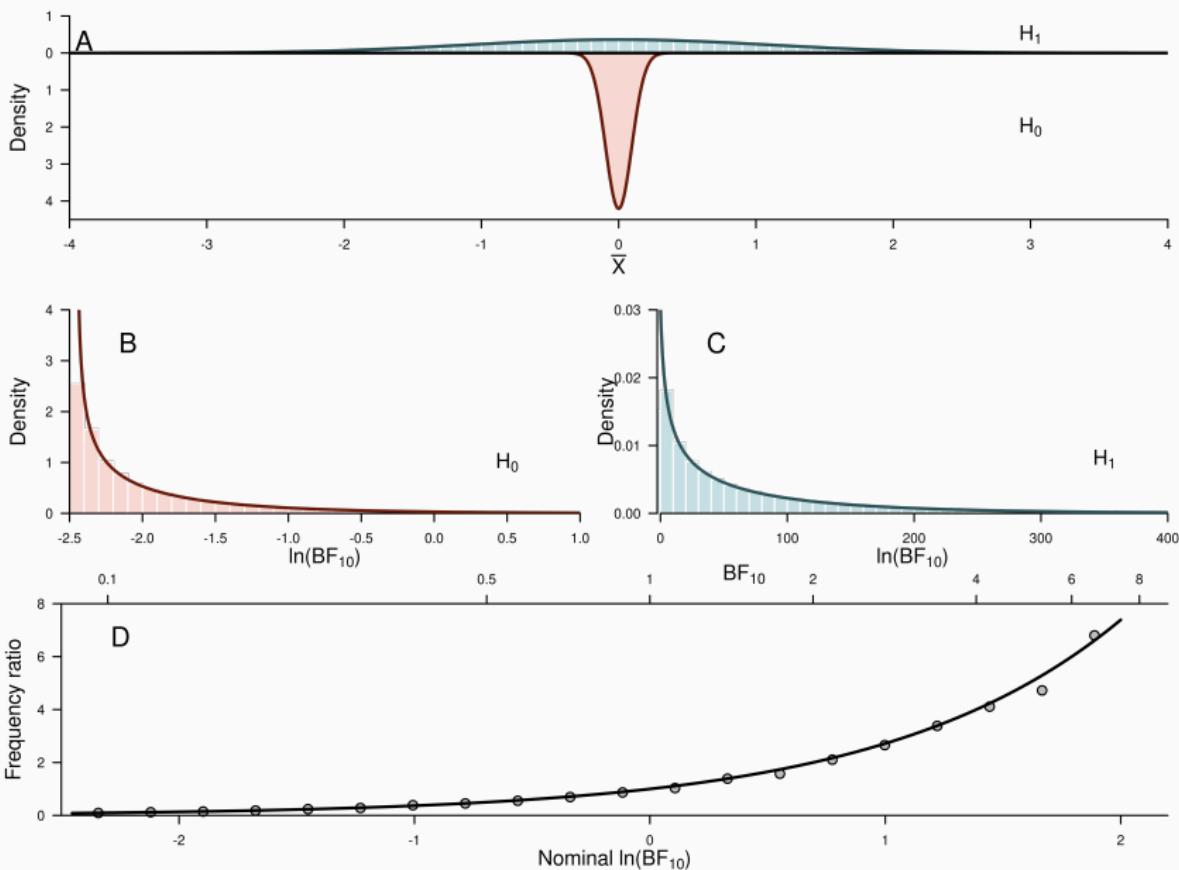
RESULTS

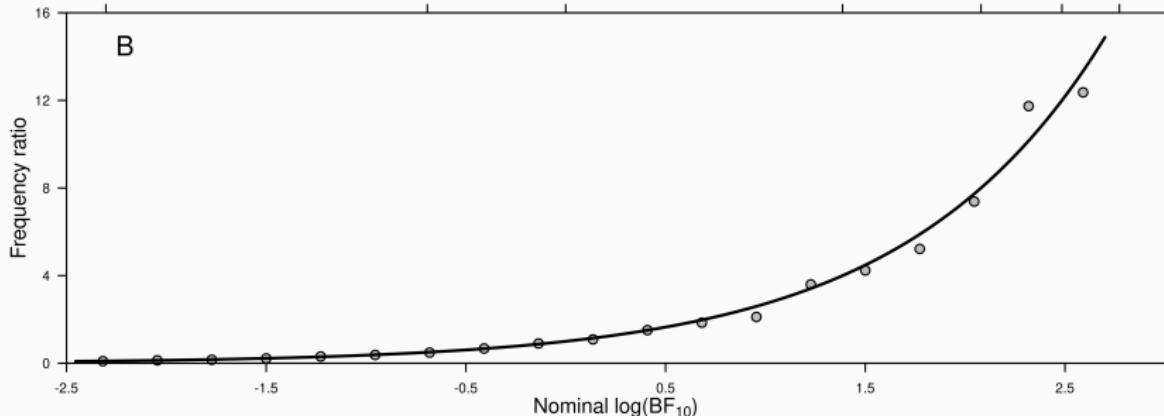
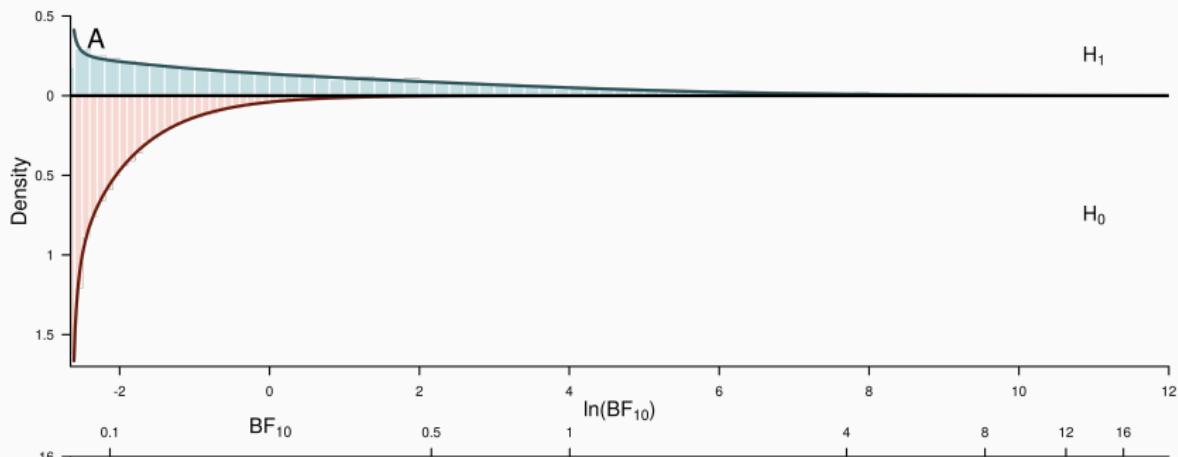
$\mathcal{H}_0 : \mu = 0$ vs $\mathcal{H}_1 : \mu = \mu_1, n$ OBSERVATIONS (ASSERTION 1)

13 / 17









DISCUSSION

We offer a mathematical proof to a Bayes factor property suggested by Rouder (2014).

Is this conclusive evidence that Bayesian optional stopping is allowed?
Well, not just yet.¹

However, in a very recent reply, Rouder again disagrees...

<https://psyarxiv.com/m6dhw/>

To be continued...

¹de Heide and Grünwald (2017).

THANK YOU!