

# Why optional stopping can be a problem for Bayesians<sup>1</sup>

—  
*Bayesian Machine Learning*

Boëzennec Robin, Hugo Simon

24/03/2022

école —  
normale —  
supérieure —  
paris — saclay —

université  
PARIS-SACLAY

---

<sup>1</sup>Rianne de Heide et al. (Nov. 2020). "Why optional stopping can be a problem for Bayesians". In: *Psychonomic Bulletin and Review* 28.3, pp. 795–812.

## 1 Introduction

p-hacking

Continuous monitoring

## 2 Definitions

Definitions

Post-odds calibration

## 3 Results

Experiment 0

Experiment 1: *Type 0 prior*

Experiment 2: *Type I prior*

*Type II prior*

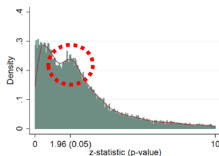
## 4 Conclusions?

- Misusing data analysis to find patterns in data that can be presented as statistically significant.
- Because of incentives to publish positive results, conflict of interest or simply statistical misunderstanding.

## Economics

Brodeur et al (*AEJ:A*, in press)

"Star Wars: The empirics strike back"

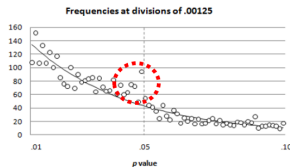


(b) De-rounded distribution of z-statistics.

## Psychology

Masicampo Lalande (*QJEP*, 2012)

"A peculiar prevalence of p values just below .05"



## Biology

Head et al (*PLOS Biology* 2015)

"Extent and Consequences of P-Hacking in Science"

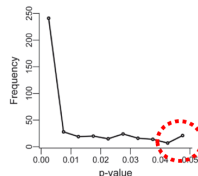


Figure 1: Dotted circles highlight excess of .05, though most p-values are way smaller.

- Experimenter can choose to arbitrarily stop aggregating data depending on already observed ones.
- One practice to be careful with using p-value in NHST.
- E.g. stopping as soon as statistical significance is reached leads to rejecting null hypothesis at a higher alpha risk than the p-value estimates.

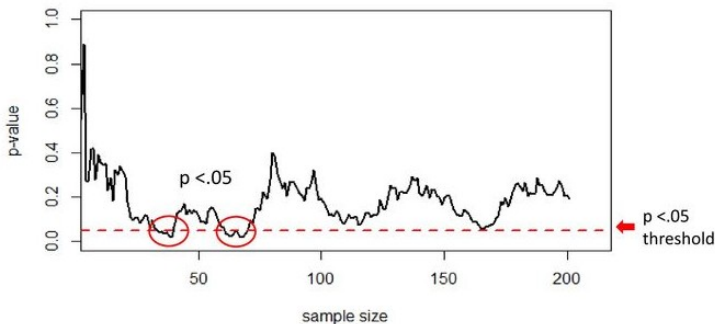


Figure 2: p-hacking with continuous monitoring on non-effect data.

## 1 Introduction

p-hacking

Continuous monitoring

## 2 Definitions

Definitions

Post-odds calibration

## 3 Results

Experiment 0

Experiment 1: *Type 0 prior*

Experiment 2: *Type I prior*

*Type II prior*

## 4 Conclusions?

## Definition (Post-odds and FDR)

Let be tested pair of hypotheses  $(H_0, H_1)$  and measurable event  $A \in \mathcal{X}$ . Then we define the odds and False Discovery rate (FDR) of  $H$  given  $A$  as

$$\odot(H | A) = \frac{\mathbb{P}(H_1 | A)}{\mathbb{P}(H_0 | A)} = \frac{1}{\mathbb{P}(H_0 | A)} - 1 = \frac{1}{\text{FDR}(H | A)} - 1$$

$\odot = 9 \iff \text{FDR} = 0.1$  means "discovering" (rejecting  $H_0$ , **accepting**  $H_1$ ) would be false with proba. 10%.

## Definition (Bayes Factor)

Applying Bayes theorem to the previously define post-odds leads to

$$\odot(H | A) = \frac{\mathbb{P}(H_1 | A)}{\mathbb{P}(H_0 | A)} = \frac{\mathbb{P}(A | H_1)}{\mathbb{P}(A | H_0)} \cdot \frac{\mathbb{P}(H_1)}{\mathbb{P}(H_0)} = \beta_H(A) \cdot \odot(H)$$

where  $\beta_H$  is called the Bayes Factor (BF) or likelihood ratio.

- Quantifies how prior belief is updated from observations.
- No need to compute data evidence  $\mathbb{P}(A)$ .

## Definition (Post-odds calibration)

A functional  $\gamma : \mathcal{X} \mapsto \mathbb{R}^+$  is said to be post-odds calibrated for the model  $H$  if  $\forall c$  with non zero probability,  $\mathbb{O}(H \mid \gamma(X) = c) = c$

We know that for fixed  $n$

- $\mathbb{O}(H \mid \mathbb{O}(H \mid X^n) = c) = c$
- $\mathbb{P}(\text{p-value}(X^n) \leq c \mid H_0) = c$

But does calibration hold with optional stopping ?

$$\underbrace{\mathbb{O}(H \mid \mathbb{O}(H \mid X^\tau))}_{\substack{\text{observed post-odds} \\ \text{(under "nominal post-odds claim")}}} \stackrel{?}{=} \underbrace{\mathbb{O}(H \mid X^\tau)}_{\text{nominal post-odds}}$$

**Yes**<sup>23</sup> ! (As long as  $H \perp\!\!\!\perp \tau$ )

<sup>2</sup>Alex Deng et al. (Oct. 2016). *Continuous Monitoring of A/B Tests without Pain: Optional Stopping in Bayesian Testing*.

<sup>3</sup>Allard Hendriksen et al. (2021). "Optional Stopping with Bayes Factors: A Categorization and Extension of Folklore Results, with an Application to Invariant Situations". In: *Bayesian Analysis* 16.3.

## 1 Introduction

- p-hacking

- Continuous monitoring

## 2 Definitions

- Definitions

- Post-odds calibration

## 3 Results

- Experiment 0

- Experiment 1: *Type 0 prior*

- Experiment 2: *Type I prior*

- Type II prior*

## 4 Conclusions?



$$\mathbb{Q}(H | X) = \frac{\exp\left(\frac{n^2 \bar{x}^2}{2(n+1)}\right)}{\sqrt{n+1}} \cdot \frac{1}{1}$$

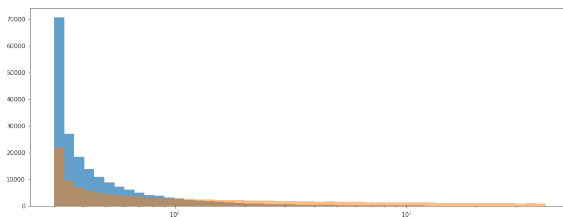
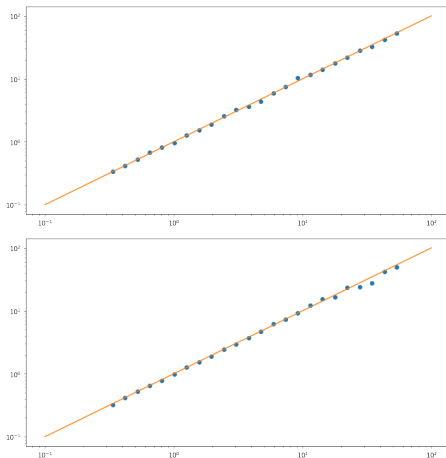


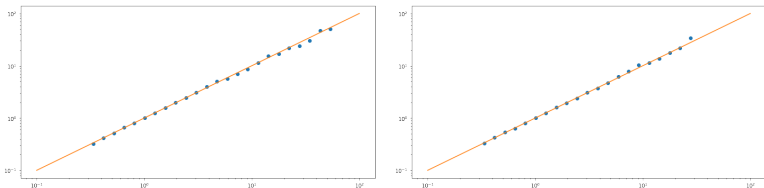
Figure 3: Histogram for the two settings



**Figure 4:** Experiment 0 calibration with  $n = 10$  on the top and with optional stopping on the bottom

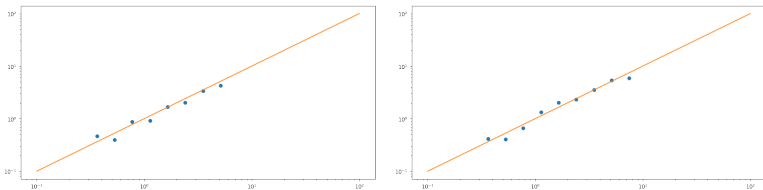
But what happens if the prior is not fully believed ?

$$\odot (H | X) = \frac{1}{\sqrt{n+1}} \left( 1 - \frac{\left( \frac{1}{n+1} \sum_{i=1}^n x_i \right)^2}{\frac{1}{n+1} \sum_{i=1}^n x_i^2} \right)$$

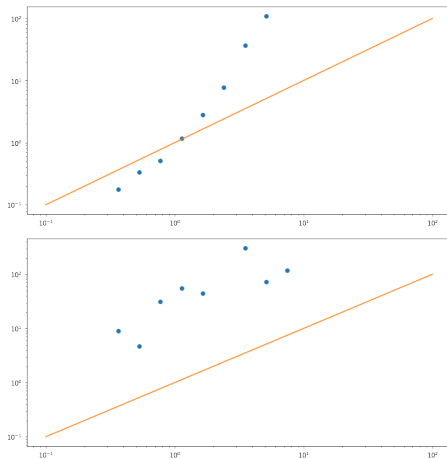


**Figure 5:** Experiment 1 calibration with  $n = 10$  on the left and with optional stopping on the right

$$\odot (H | X) = \frac{\exp\left(-\frac{\sum_{i=1}^n x_i^2}{2\sigma^2}\right)}{\int_{-\infty}^{\infty} \frac{1}{\pi(1+\mu^2)} \exp\left(-\frac{\sum_{i=1}^n (x_i - \mu)^2}{2\sigma^2}\right) d\mu}$$

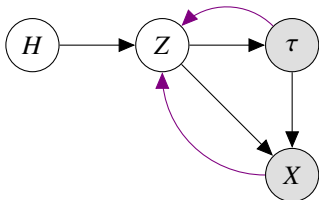


**Figure 6:** Experiment 2 with Cauchy prior and calibration with  $n = 10$  on the left and with optional stopping on the right (our results)



**Figure 7:** Experiment 2 with constant prior and calibration with  $n = 10$  on the left and with optional stopping on the right

- Priors that depend on the sample size and sometimes data itself.
- Common in some Bayesian literature though they do not define a generative model.
- Furthermore, the introduction of a stopping time makes it unclear which prior we should use
  - prior with fixed  $n$  or with  $\tau$ ?
  - Jeffreys prior considering optional stopping?



## 1 Introduction

p-hacking

Continuous monitoring

## 2 Definitions

Definitions

Post-odds calibration

## 3 Results

Experiment 0

Experiment 1: *Type 0 prior*

Experiment 2: *Type I prior*

*Type II prior*

## 4 Conclusions?

## Objective Bayesian

- *Type 0 prior* ☒
  - works fine
- *Type I prior* ☐ X
  - robustness to priors is crucial
  - $\tau$  just amplifies calibration loss
- *Type II prior* ☐ X
  - no satisfying calibration definition

## Subjective Bayesian

- *Type 0 prior* ☒
  - works fine
- *Type I prior* ☐  $\approx$ 
  - well if this is what you believe...
  - ...but might look for consensus
- *Type II prior* ☐ ?
  - how can this formalize any belief ?



- Deng, Alex et al. (Oct. 2016). *Continuous Monitoring of A/B Tests without Pain: Optional Stopping in Bayesian Testing*.
- Heide, Rianne de et al. (Nov. 2020). “Why optional stopping can be a problem for Bayesians”. In: *Psychonomic Bulletin and Review* 28.3, pp. 795–812.
- Hendriksen, Allard et al. (2021). “Optional Stopping with Bayes Factors: A Categorization and Extension of Folklore Results, with an Application to Invariant Situations”. In: *Bayesian Analysis* 16.3.

Thank you!

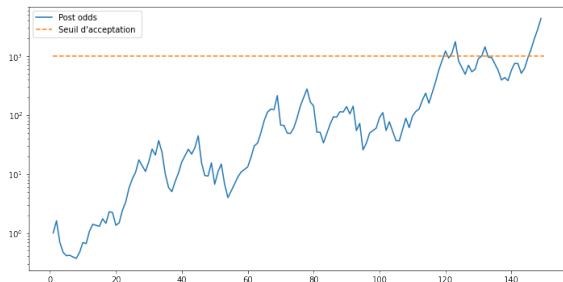


Figure 8: Post odds for our experiment on real data