# An empirically-driven guide on using Bayes Factors for M/EEG decoding

Lina Teichmann[1*], Denise Moerel[2], Chris Baker[1], Tijl Grootswagers[3]

[1] Laboratory of Brain and Cognition, National Institute of Mental Health, Bethesda, MD, USA

[2] Department of Cognitive Science, Macquarie University, Sydney, Australia

[3] The MARCS Institute for Brain, Behaviour & Development, Western Sydney University, Sydney, Australia

*Corresponding author: lina.teichmann@nih.gov

# Abstract

Bayes Factors can be used to provide quantifiable evidence for contrasting hypotheses and have thus become increasingly popular in cognitive science. However, Bayes Factors are rarely used to statistically assess the results of neuroimaging experiments. Here, we provide an empirically-driven guide on implementing Bayes Factors for time-series neural decoding results. Using real and simulated Magnetoencephalography (MEG) data, we examine how parameters such as the shape of the prior and data size affect Bayes Factors. Additionally, we discuss benefits Bayes Factors bring to analysing multivariate pattern analysis data and show how using Bayes Factors can be used instead or in addition to traditional frequentist approaches.

# 1.   Introduction

The goal of multivariate decoding in cognitive neuroscience is to infer whether information is represented in the brain (Hebart & Baker, 2018). To draw meaningful conclusions in this information-based framework, we need to statistically assess whether the conditions of interest evoke different data patterns. In the context of time-resolved neuroimaging data, activation patterns are extracted across MEG or EEG sensors and classification accuracies are used to estimate information at every timepoint (see Figure 1 for an example). Currently, null hypothesis statistical testing (NHST) and p-values are the de-facto method of choice for statistically assessing classification accuracies, but recent studies have started using Bayes Factors (Grootswagers et al., 2021; e.g., Grootswagers, Robinson, & Carlson, 2019b; Grootswagers, Robinson, Shatek, et al., 2019; Kaiser et al., 2018; Karimi-Rouzbahani et al., 2021; Mai et al., 2019; Proklova et al., 2019; Robinson et al., 2019, 2021). Under the null hypothesis, the mean equals chance decoding and under the alternative hypothesis the mean is larger than chance decoding. The direct comparison of the predictions of two hypotheses is one of the strengths of the Bayesian framework of hypothesis testing (Jeffreys, 1939, 1935). Bayes Factors describe the probability of one hypothesis over the other given the observed data. In the multivariate pattern analysis (MVPA) context, this means we use Bayes Factors to test the probability of above-chance classification versus at-chance classification given the decoding results across participants at each timepoint. The goal of the current paper is to present and discuss Bayes Factors from a practical standpoint in the context of time-series decoding, while referring the reader to published work focusing on the theoretical and technical background of Bayes Factors.
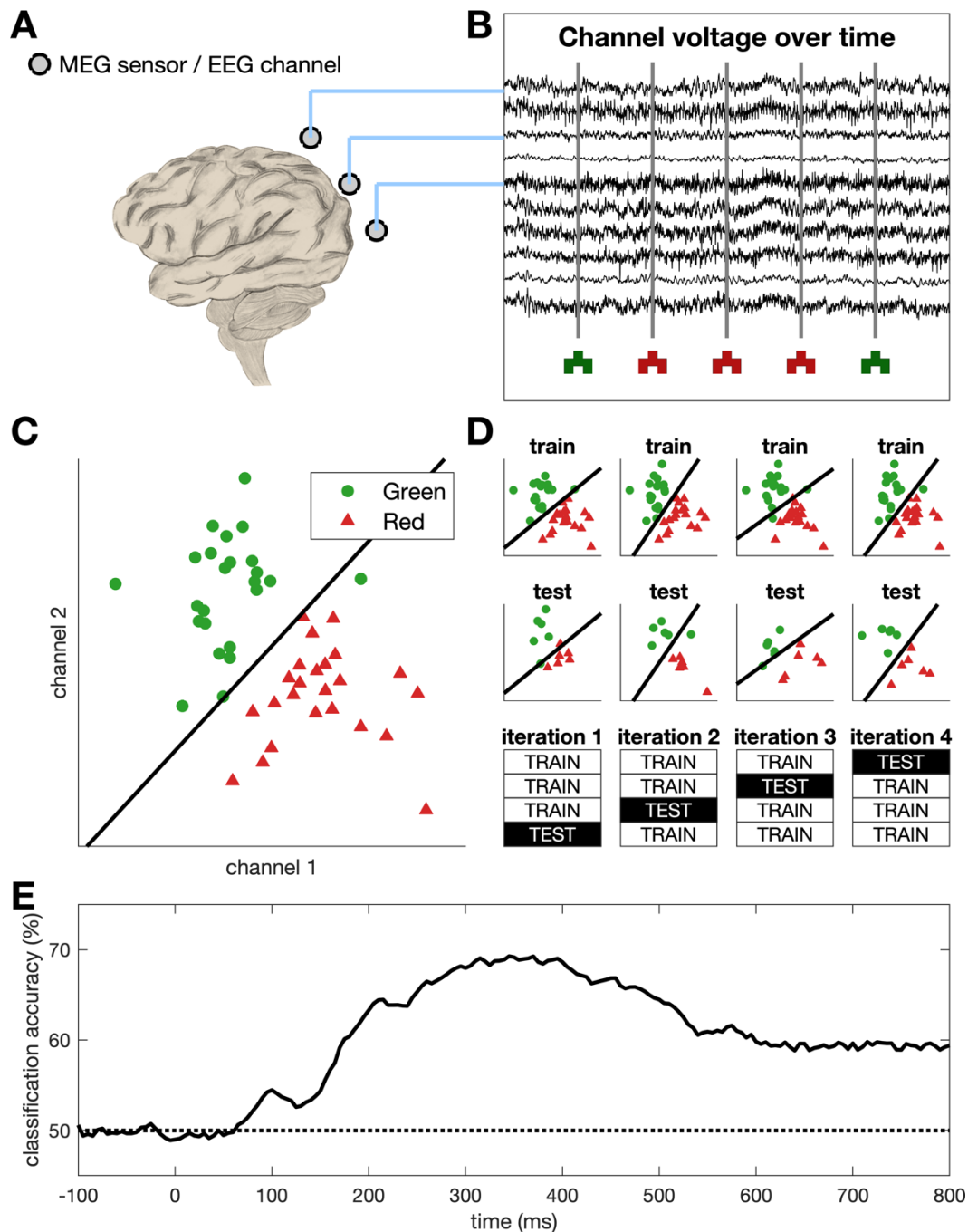
The Bayesian approach brings several advantages over the traditional NHST framework (Dienes, 2011, 2014, 2016b; Keysers et al., 2020; Morey et al., 2016; Wagenmakers et al., 2018). In addition to allowing us to contrast evidence for above-chance versus at-chance decoding directly, Bayes Factors are a measure of strength of evidence for one hypothesis versus another. That means, we can *directly* assess how much evidence we have for different analyses. For example, if we were interested in testing whether viewing different colours evokes different neural responses, we could examine differences in the neural signal evoked by seeing red, green, and yellow objects. Using Bayes Factors, we could then directly compare whether red versus green can be decoded as well as red versus yellow. Larger Bayes Factors reflect more evidence which makes the interpretation of statistical results across analyses more intuitive. Another advantage is that Bayes Factors can be calculated iteratively while more data are being collected and that testing can be stopped when there is a sufficient

70    amount of evidence (Keysers et al., 2020; Wagenmakers et al., 2018). Such stopping-rules

71    could be accompanied by a pre-specified acquisition plan and potentially an (informal)

72    preregistration via portals such as the Open Science Framework (Foster & Deardorff, 2017).

73    Using the data to determine when enough evidence has been collected is particularly relevant

74    for neuroimaging experiments, as it might significantly reduce research costs and reduce the

75    risk of having underpowered studies. Thus, using a Bayesian approach to statistically assess

76    time-series classification results can be beneficial both from a theoretical as well as an

77    economic standpoint and might ease the ability to interpret and communicate scientific

78    findings.

79

80    While Bayes Factors provide an alternative to the more traditional NHST framework,

81    incorporating Bayes Factors into existing time-series decoding pipelines may seem daunting.

82    Introductory papers often focus on mathematical aspects, and on relatively straightforward

83    behavioural experiments (e.g., Ly et al., 2016; Morey et al., 2016; Rouder et al., 2009). We

84    present an example based on a previously published time-series decoding study (Teichmann

85    et al., 2019) and will present results from simulations to show the influence of certain

86    parameters on Bayes Factors. We make use of the established Bayes Factor R package

87    (Morey et al., 2015) to calculate the Bayes Factors but provide sample codes along with this

88    paper showing how to access the Bayes Factor R package via Matlab and Python

89    (https://github.com/LinaTeichmann1/BFF_repo). We also show how the Bayes Factors in our

90    example compare to p-values. Based on empirical evidence, we will give recommendations

91    for Bayesian analysis applied to M/EEG classification results. The aim of this paper is to

92    provide a broad introduction to Bayes Factors from a viewpoint of time-series neuroimaging

93    decoding. We aim to do so without going into the technical or mathematical detail, and instead

94    provide pointers to relevant literature on the specifics.

95

96

**Figure 1**. **Overview of MVPA for time-series neural data.** (A) Example MEG sensors / EEG channels. (B) Simulated time-series neuroimaging data for a few sensors/channels. Vertical lines show stimulus onsets with example stimuli plotted below. Data is first epoched from -100 to 800 ms relative to stimulus onset, resulting in multiple time-series chunks associated with seeing a red or a green shape. (C) Using the epoched data, we can extract the sensor/channel activation pattern across the different sensors/channels (only 2 displayed for simplicity) for every trial at every timepoint. Then a classifier (black line) is trained to differentiate between the activation patterns evoked by red and green trials. The shape of the stimuli is not relevant in this context. (D) Example of a 4-fold cross validation where the classifier is trained on three quarters of the data and tested on the left-out quarter. This process is repeated at every timepoint. (E) We can calculate how often the classifier accurately predicts the colour of the stimulus at each timepoint by averaging across all testing folds. Theoretical chance level is 50% as there are two conditions in the simulated data (red and green). During the period

110  before stimulus onset, we expect decoding to be at chance, and thus the baseline period can
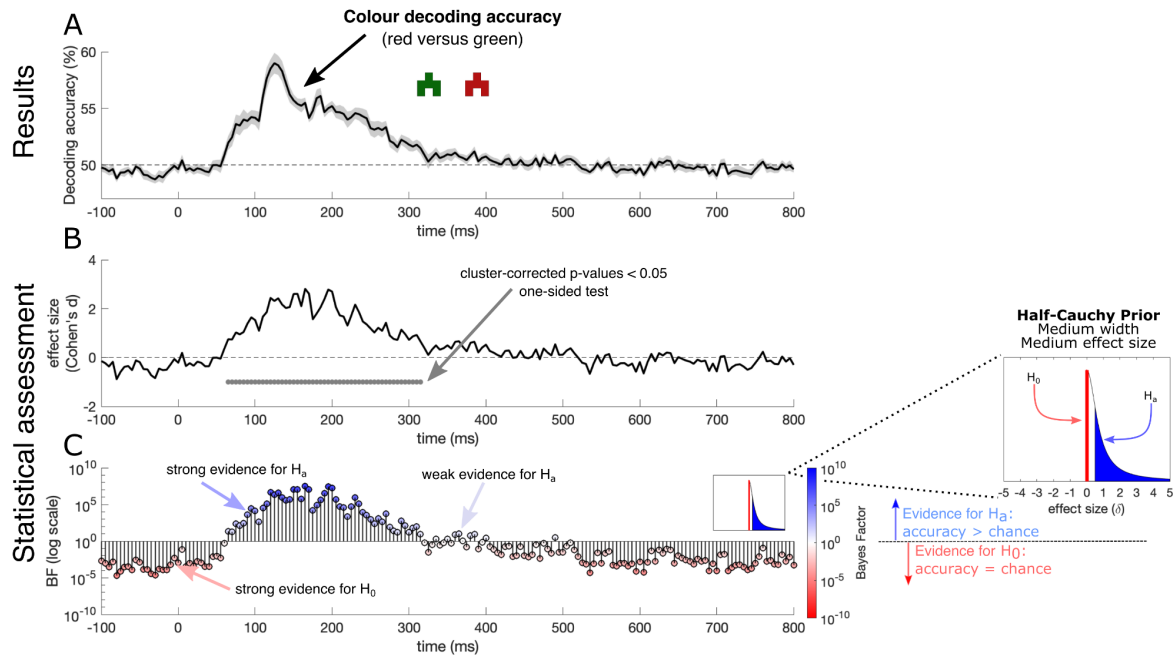111  serve as a sanity check.

112

# 2.   Methods & Results

114  2.1 Example dataset & inferences based of Bayes Factors

115  The aim of the current paper is to show how to use Bayes Factors when assessing time-series
116  neuroimaging classification results and test what effect different analysis parameters have on
117  the results. We have used a practical example of previously published MEG data (Teichmann
118  et al., 2019), which we re-analysed using Bayes Factors. In the original experiment, eighteen
119  participants viewed coloured shapes and grayscale objects in separate blocks while the neural
120  signal was recorded using MEG. Here, we only considered the coloured shape trials ("real
121  colour blocks", 1600 trials in total). Identical shapes were coloured in red or green and were
122  shown for 100 ms followed by an inter-stimulus-interval of 800-1100 ms. The data was
123  epoched from -100 ms to 800 ms (200 Hz resolution) relative to stimulus onset and a linear
124  classifier was used to differentiate between the neural responses evoked by red and green
125  shapes. A 5-fold cross-validation was used with the classifier being trained on 80% of the data
126  and tested on the remaining 20%. This classification analysis resulted in decoding accuracies
127  over time for each participant. In the original study, permutation tests and cluster-corrected p-
128  values were used to assess decoding accuracies as implemented in CoSMoMVPA (Oosterhof
129  et al., 2016). Here, we calculated Bayes Factors instead and examined how parameter
130  changes affected the results.

131

132  When running statistical tests on classification results, we are interested in whether decoding
133  accuracy is above-chance at each timepoint. To test this using a frequentist approach, we can
134  use permutation tests to establish whether there is enough evidence to reject $H_0$ which states
135  that decoding is equal to chance. If there is enough evidence, we can reject $H_0$ and conclude
136  that decoding is different from chance. Given that below-chance decoding accuracies are not
137  meaningful, we usually are interested only in above-chance decoding (directional hypothesis).
138  In contrast to the frequentist approach, Bayes Factors quantify how much the plausibility of
139  two hypotheses changes, given the data (see e.g., Ly et al., 2016). Here, we ran a Bayesian
140  t-test of Bayes Factor R package (Morey et al., 2015) at each timepoint, testing whether the
141  data is more consistent with $H_a$ (decoding is larger than chance) over $H_0$ (decoding is equal to
142  chance). The resulting Bayes Factors center around 1 with numbers smaller than 1
143  representing evidence for $H_0$ and numbers larger than 1 representing evidence for $H_a$. In
144  contrast to p-values, Bayes Factors are directly interpretable and comparable (cf. Keysers et

145    al., 2020; Morey et al., 2016; Wagenmakers et al., 2016). That is, a Bayes Factor of 10 means

146    the data is 10 times more likely to be observed under $H_a$ as opposed to $H_0$. Similarly, a Bayes

147    Factor of 1/10 means the data is 10 times more likely to be observed under $H_0$ as opposed to

148    $H_a$ . Thus, in the context of time-series decoding, Bayes Factors allow us to directly assess

149    whether and how much evidence there is at a given timepoint for the alternative over the null

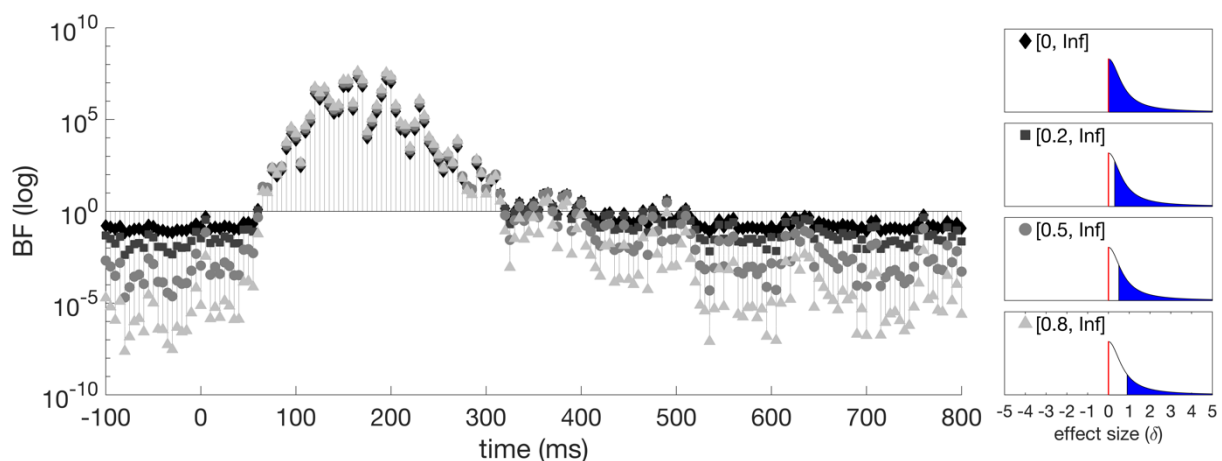150    hypothesis and *vice versa* (Figure 2C).

151



152

153    **Figure 2**. **Decoding results of our practical example dataset with statistical**
154    **assessments.** (A) Colour decoding over time (black line). The dashed line shows theoretical
155    chance decoding (50%). The grey shaded area represents the standard error across
156    participants. (B) Effect size over time with the cluster-corrected p-values at each timepoint
157    printed below in grey. (C) Bayes Factors over time for this dataset on a logarithmic scale. Blue,
158    upwards pointing stems indicate evidence for above-chance decoding and red, downwards
159    pointing stems show evidence for at-chance decoding at every timepoint. We used a hybrid
160    one-sided model comparing evidence for above-chance decoding versus a point-nil at $\delta = 0$
161    (no effect). For the alternative hypothesis, we used a half-Cauchy prior with medium width (r
162    = 0.707) covering an interval from $\delta = 0.5$ to $\delta = \infty$. The half-Cauchy prior assumes that small
163    effect sizes are more likely than large ones, but the addition of the interval deems very small
164    effects $\delta < 0.5$ as irrelevant. During the baseline period (i.e., before stimulus onset), the Bayes
165    Factors strongly support the null hypothesis, confirming the sanity check expectation.
166

167    2.2 Adjusting the prior range to account observed chance decoding

168    Bayes Factors represent the plausibility that the data emerged from one hypothesis compared

169    to another. In the example dataset, the two hypotheses are that decoding is at chance (i.e.,

170    $H_0$, no colour information present) or that decoding is above chance (i.e., $H_a$, colour

171    information present). To deal with the fact that observed chance decoding can be different

172    than the theoretical chance level, we can adjust the prior range of the alternative hypothesis
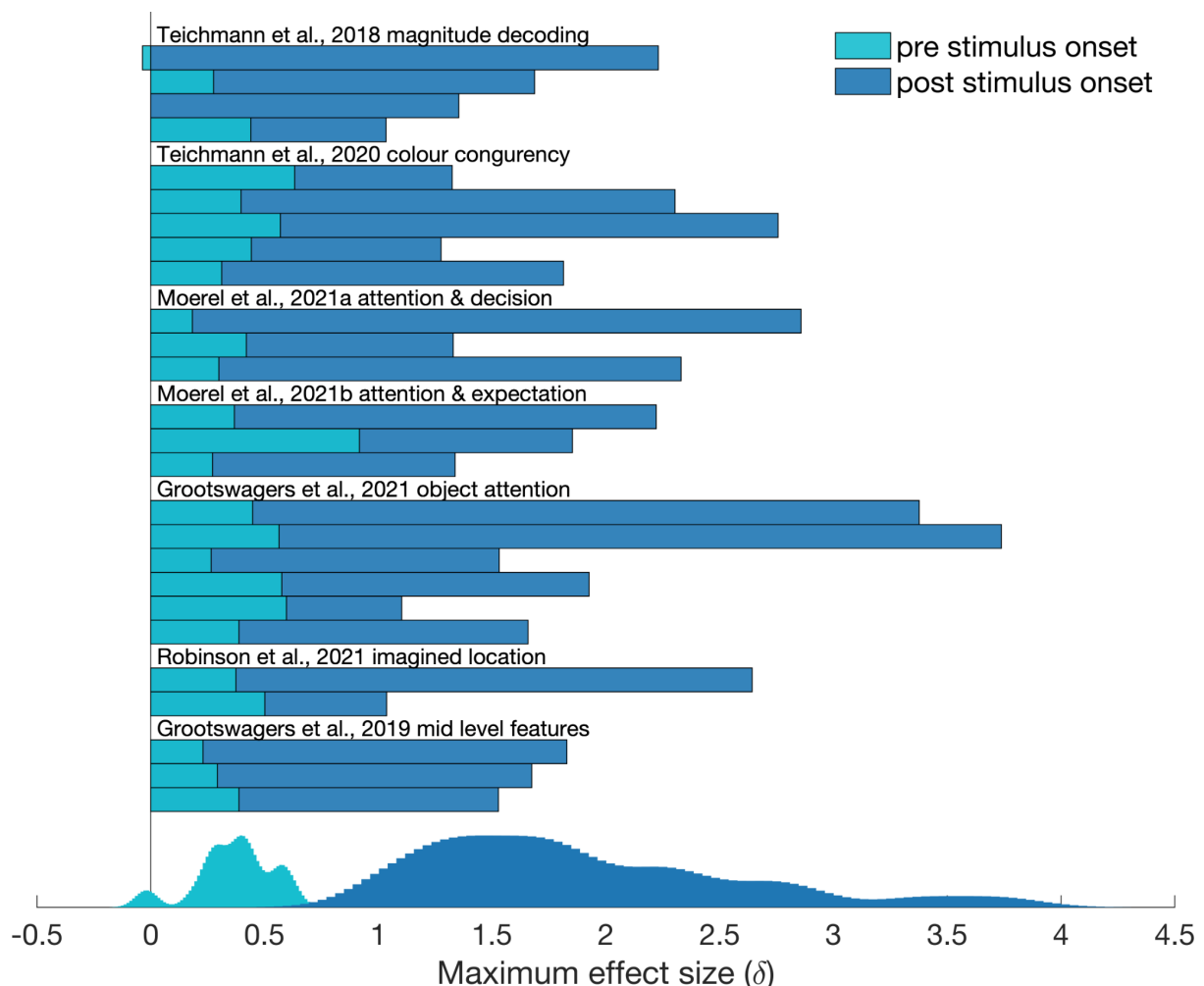
173    to allow for small effects under the null hypothesis (Rouder et al., 2009). The prior range (called

174    "null interval" in the R package) is defined in standardized effect sizes and consists of a lower

175    and upper bound. To incorporate the differences between observed and theoretical chance

176    level, we can define a range of relevant effect sizes for the alternative hypothesis, for example,

177    from $\delta = 0.5$ to $\delta = \infty$. To determine which values are reasonable as the lower bound of this

178    interval, we changed the prior range systematically and examined the effect on the resulting

179    Bayes Factors (Figure 3). We found that smaller lower bounds at $\delta = 0$ and $\delta = 0.2$ resulted in

180    weaker evidence supporting the null hypothesis than ranges starting at $\delta = 0.5$ and $\delta = 0.8$.

181    The range did not have a large effect on timepoints with strong evidence for $H_a$. The effect of

182    changing the prior range is larger for the null hypothesis than the alternative as chance

183    decoding is not exactly 50% but distributed around chance. Changing the lower bound of the

184    prior range means that the effects that are just larger than $\delta = 0$ can support the null

185    hypothesis. Thus, the results here demonstrate that we can compensate for the differences

186    between theoretical and observed chance by adjusting the prior range and effectively

187    considering small effect sizes as evidence for the null hypothesis rather than the alternative.

188



189

190    **Figure 3. The effect of changing the prior range (null interval) on Bayes Factors in our**
191    **example data.** Intervals starting at larger effect sizes led to more timepoints showing
192    conclusive evidence for $H_0$. This is due to the fact that theoretical and observed chance levels
193    are not the same. The panels on the right show the prior distributions with the different null
194    intervals.

195

196    To further examine what a reasonable lower bound of the prior range is, we looked at effect

197    sizes observed during the baseline window (before stimulus onset) in a selection of our

198    previous studies (Grootswagers et al., 2021; Grootswagers, Robinson, & Carlson, 2019a;

199    Moerel, Grootswagers, et al., 2021; Moerel, Rich, et al., 2021; Teichmann et al., 2018, 2020).

200    Using the baseline window allows us to quantify the difference between theoretical and

201    observed chance, as we do not expect any meaningful effects before stimulus onset (e.g.,

202    stimulus colour is not decodable before the stimulus is presented). Thus, the baseline period

203     can effectively tell us which effect sizes can be expected by chance. Using this method, we

204     estimated maximum effect sizes for different analyses in each paper (see different bars in

205     Figure 4). Across our selection of prior studies, we found an average maximum effect size of

206     $\delta$ = 0.39 before stimulus onset and an average maximum effect size of $\delta$ = 1.91 after stimulus

207     onset (Figure 4). This survey shows that effect sizes as large as $\delta$ = 0.5 can be observed when

208     no meaningful information is in the signal. Thus, this supports the conclusions from the

209     example dataset showing that prior ranges with a lower bound of $\delta$ = 0.5 may be a sensible

210     choice when using Bayes Factors to examine time-series M/EEG decoding results.

211



212

213 **Figure 4. Estimated maximum effect sizes during baseline and after stimulus onset for**
214 **prior decoding studies that used visual stimuli.** Using already published data, we
215 calculated the maximum effect sizes during the baseline (light blue) and post-stimulus (dark
216 blue) to estimate typical peak effect sizes in visual decoding studies. Each bar represents a
217 unique analysis within the paper. The estimations show that a reasonable range for $H_a$ would
218 start at $\delta$ = 0.5 or above, as during baseline decoding accuracies corresponding to
219 standardized effect sizes as high as $\delta$ = 0.5 were observed.
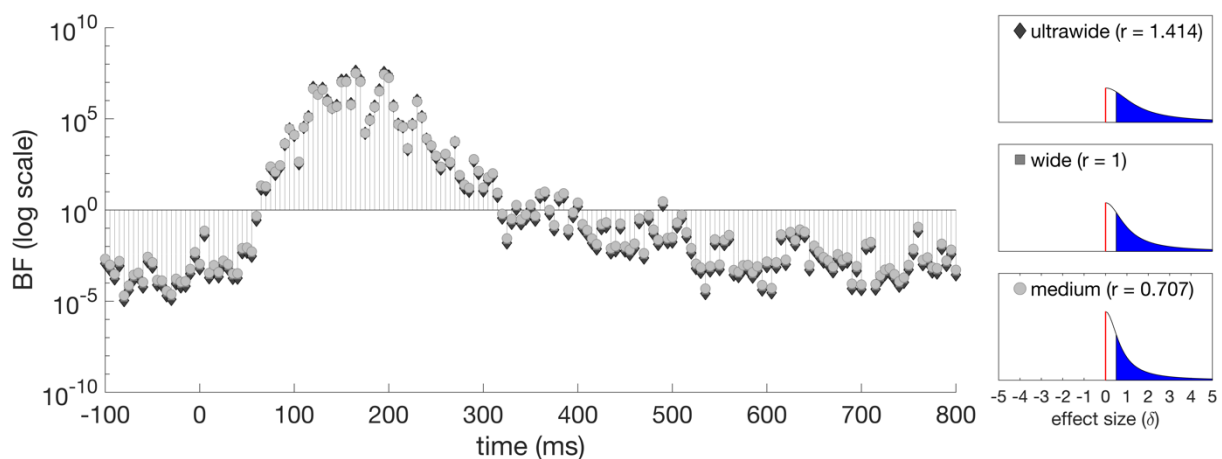220

221    2.3 Changing the prior width to capture different effect sizes

222    Another feature that can be changed in the Bayesian t-test is the width of the half-Cauchy

223    distribution (referred to as r-value in the Bayes Factor Package). Small r-values create a

224    narrower, sharply peaking distribution, whereas larger values make the distribution wider with

225    a prolonged peak. Standard prior widths incorporated in the Bayes Factor R package are

226    medium (r = 0.707), wide (r = 1), and ultrawide (r = 1.414). Keeping the prior range consistent

227    ([0.5, Inf]) while using the three prior widths implemented into the R Bayes Factor Package

228    (medium = 0.707; wide = 1; ultrawide = 1.414). We found that changing the width of the

229    Cauchy prior did not have a pronounced effect on the Bayes Factors (Figure 5). In our specific

230    example, this is probably the case because the effect sizes quickly rose to $\delta > 2$ (Figure 2b)

231    which means that the subtle differences between the different prior widths do not have a

232    substantial effect on the likelihood of the data arising from $H_a$ over $H_0$. Thus, using the default

233    prior width (r = 0.707) for the decoding context seems like a reasonable choice.

234



235

**Figure 5**. **Bayes Factors over time for the example data set when the prior width is changed.** The width of the prior had no pronounced effect on the Bayes Factors we calculated. The panels on the right show the prior distributions with the different widths.

240    2.4 The effect of data size on statistical inferences

241    In a lot of cases, there are financial and time limits on how many participants can be tested

242    and for how long. To obtain an estimate of how much data is needed to draw conclusions and

243    avoid ending up with underpowered studies, we used the example dataset and reduced the

244    data size for analysis. As classification analyses are usually run at the subject level but

245    statistical assessment is run at the group level, we tested how changing data size both by trial

246    numbers and participant numbers influences Bayes Factors in the time-series decoding

247    context (Figure 6). In the original example dataset, the classifier was trained on 1408 trials

248    and tested on 352 trials (5-fold cross-validation). There were five different shapes in the red

249    and the green condition (160 repetitions for each coloured shape) and the cross-validation

250    schema was based on leaving all trials of one shape out for testing. Statistical inferences were

251    drawn on the group level which contained data from 18 participants. To examine the effect of

252    data size (and effectively noise level) on the Bayes Factor calculations, we re-ran the analysis

253    reducing the data size first by retaining the first 1200 (75%), 800 (50%), 400 (25%), or 160

254    (10%) trials participants completed. We cross-validated in the same way as in the original

255    paper, with the only difference being how many trials of each shape were included. In addition,

256    we subsampled from the whole group, retaining data from the first 6, 12, or all 18 participants

257    and re-ran the statistical analysis. We then compared the results from the reduced-size colour

258    datasets using Bayes Factors and cluster-corrected p-values[1].
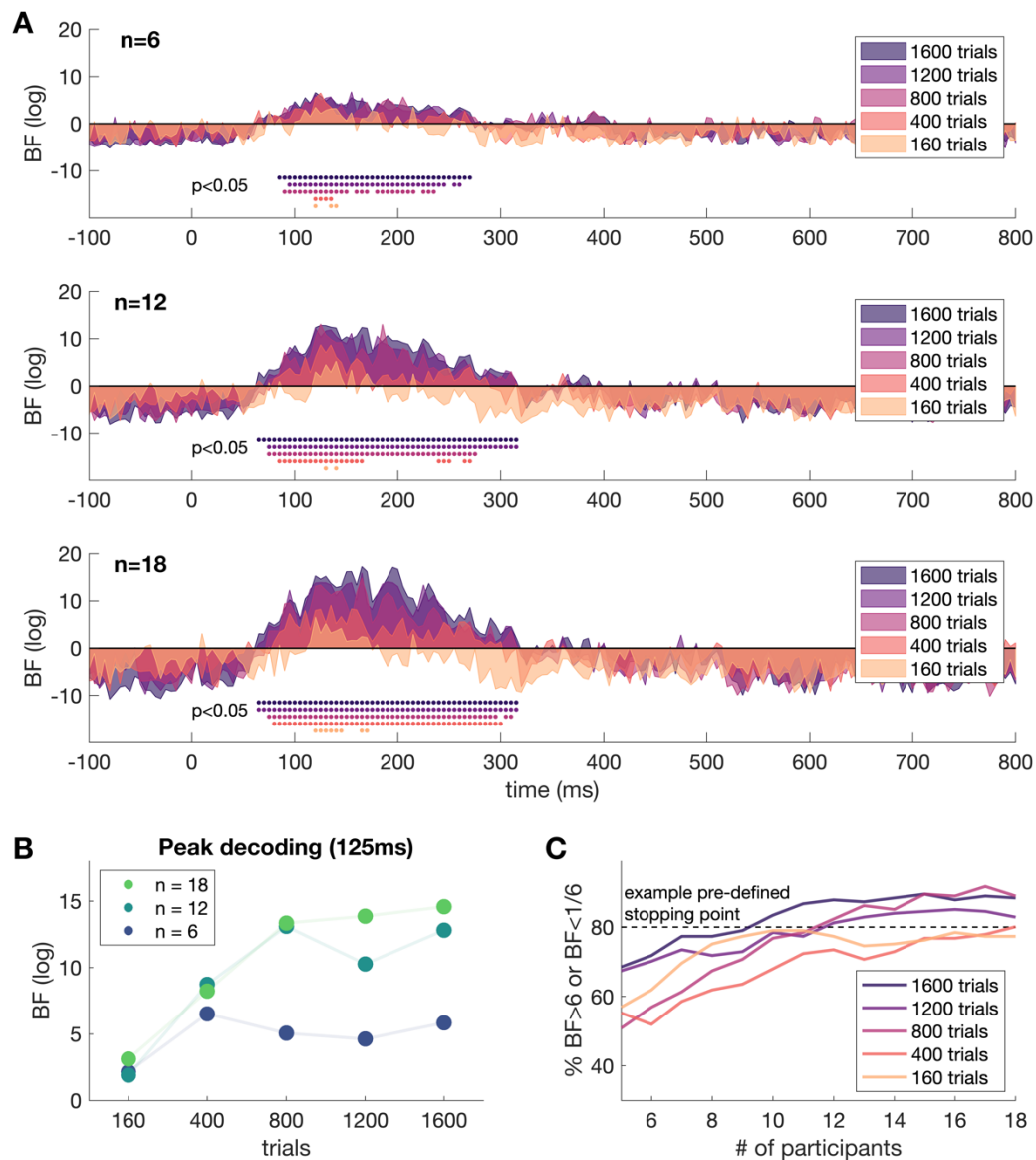
259

260    Overall, our analyses highlight that we need to have a large enough number of trials and a

261    large enough number of participants to draw firm conclusions about our time-resolved

262    decoding results. Testing more participants resulted in stronger evidence for $H_a$ and $H_0$, with

263    fewer timepoints in the inconclusive range (Bayes Factors) and more significant above-chance

264    decoding timepoints (p-values). Similarly, running the classification with more trials, led to

265    more timepoints with large Bayes Factors supporting $H_a$ and more above-chance decoding

266    timepoints. However, one of the key advantages of using Bayes Factors instead of p-values

267    is that we can potentially obtain a good idea of how many trials are needed even if we run a

268    pilot experiment with a limited number of participants. A reasonable strategy would be to

269    overpower the subject-level data (i.e., number of trials) for the pilot sample and then sub-

270    sample to explore how many trials are needed. In our example, we can see that the amount

271    of evidence for $H_a$ at peak decoding is not sufficient when we only use 160 trials (10% of the

272    original sample), regardless of the number of subjects. Increasing the trials to 400 or 800 (25%

273    or 50% of the original sample) leads to similar conclusions as using all 1600 trials. As Bayesian

274    statistics allow for sequential sampling, we could collect data from more participants until a

275    criterion is reached. For example, if we had pre-defined a stopping criterion as 80% of the

276    timepoints being in the conclusive range (Bayes Factors larger than 6 or smaller than 1/6), we

277    would have been able to stop collecting data after 9 participants completed 1600 trials or after

278    18 participants completed 400 (Figure 6c). Overall, the data suggest that insufficient data at

279    the subject-level ultimately leads to inconclusive evidence, highlighting that a large number of

280    trials is just as, if not more important, than large numbers of participants.

281

282

---

[1] In comparison to the original paper, we did not use trial label permutations. Instead, we performed sign-flip permutations (which reduces the computational time) as implemented in CoSMoMVPA to generate the null distribution.

**Figure 6. Results of the colour MEG decoding study, using a limited number of trials and participant data to simulate a piloting scenario.** (A) The first three plots show Bayes Factors over time along with cluster-corrected p-values. The colour in all plots reflects the number of trials used to train and test the classifier. (B) Compares Bayes Factors at peak decoding (125ms) for the different data sizes. (C) Compares how many participants would have needed to be tested given the different number of trials with an example pre-defined stopping point. For example, with 1600 trials and >9 participants, 80% of the Bayes factors (at different time points) exceeded 6 or 1/6. With fewer trials, more participants are needed to reach this example stopping point.

294    The example dataset provides insight into the effect of parameters such as data size and prior

295    shape on Bayes Factors. However, it is possible that different studies find different effect sizes.

296    We simulated larger datasets with fixed effect sizes between $\delta = 0$ and $\delta = 1$ to examine the

297    interaction of sample size with different prior ranges for different effect sizes (Figure 7). We

298    simulated 1000 datasets with specific effect sizes for each sample size and calculated the

299    Bayes Factors. We then calculated the median Bayes Factor for each sample- and effect size

300    combination to show how prior range choices interact with the possibility of finding evidence

301    for effects of different sizes. Specifically, we compared a prior range of 0.5 to infinity (Figure

302    7A) to a prior range of zero to infinity (Figure 7B).

303

304    When specifying the prior range to 0.5 to infinity (Figure 7A), our results show that small

305    sample sizes are sufficient to draw solid conclusions when the effect sizes are near the

306    extremes. For example, the simulations showed that there is substantial evidence for $H_0$ from

307    a small sample size if the true effect is very small. In contrast, if the effect size fell in between

308    the specified ranges for the prior of $H_a$ and $H_0$ (i.e., between 0 and 0.5), we found that small

309    sample sizes tended to result in inconclusive Bayes Factors neither supporting $H_a$ or $H_0$.

310    However, if the sample size increased, the confidence that these effects were "real" also

311    increased and therefore resulted in stronger confidence supporting one of the hypotheses.

312    Importantly, however, large sample sizes did not automatically lead to an interpretable Bayes

313    Factor if the effect was truly in between the specified prior ranges of $H_a$ and $H_0$, indicating that

314    sample size had no effect on Bayes Factors in this case.
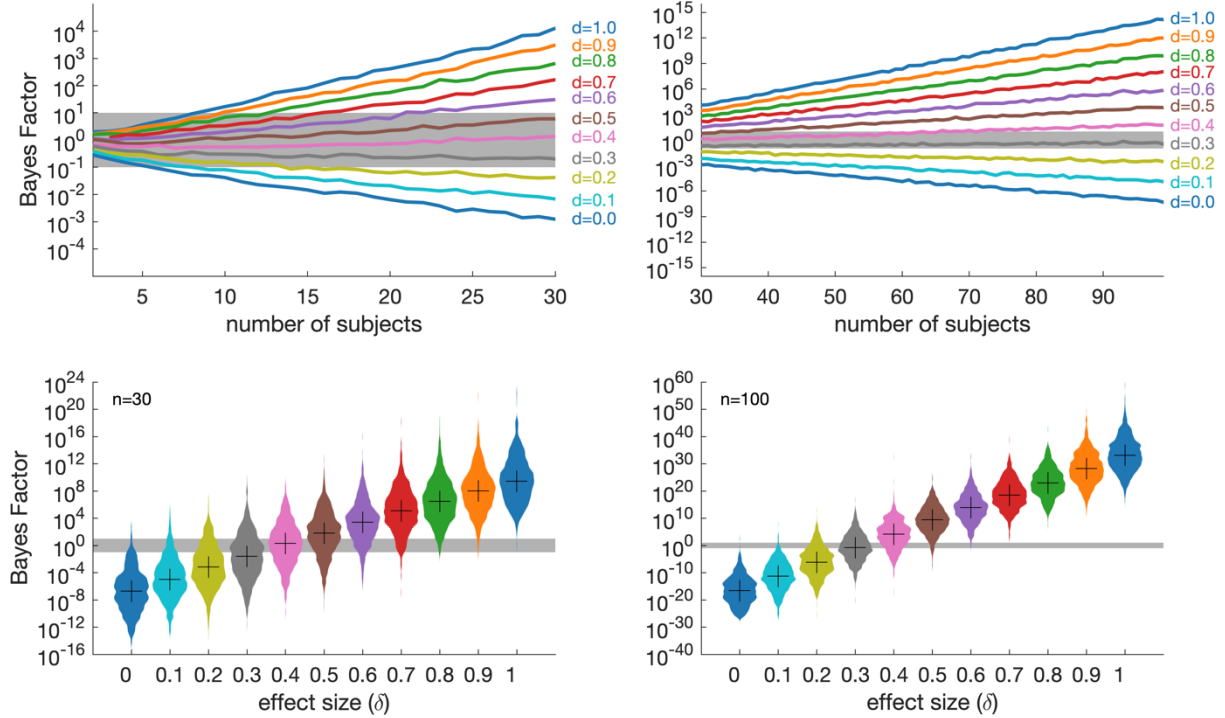
315

316    Consistent with our results for the example data, the simulations also showed that changing

317    the range of the prior has a strong effect on finding substantial evidence for $H_0$. If the prior

318    range for the alternative is specified to start at zero (Figure 7B), it was almost impossible to

319    find any evidence for $H_0$, even if the effect size was truly zero. Thus, the simulations show that

320    defining the prior range with a gap between effects expected under $H_0$ and $H_a$ is critical and

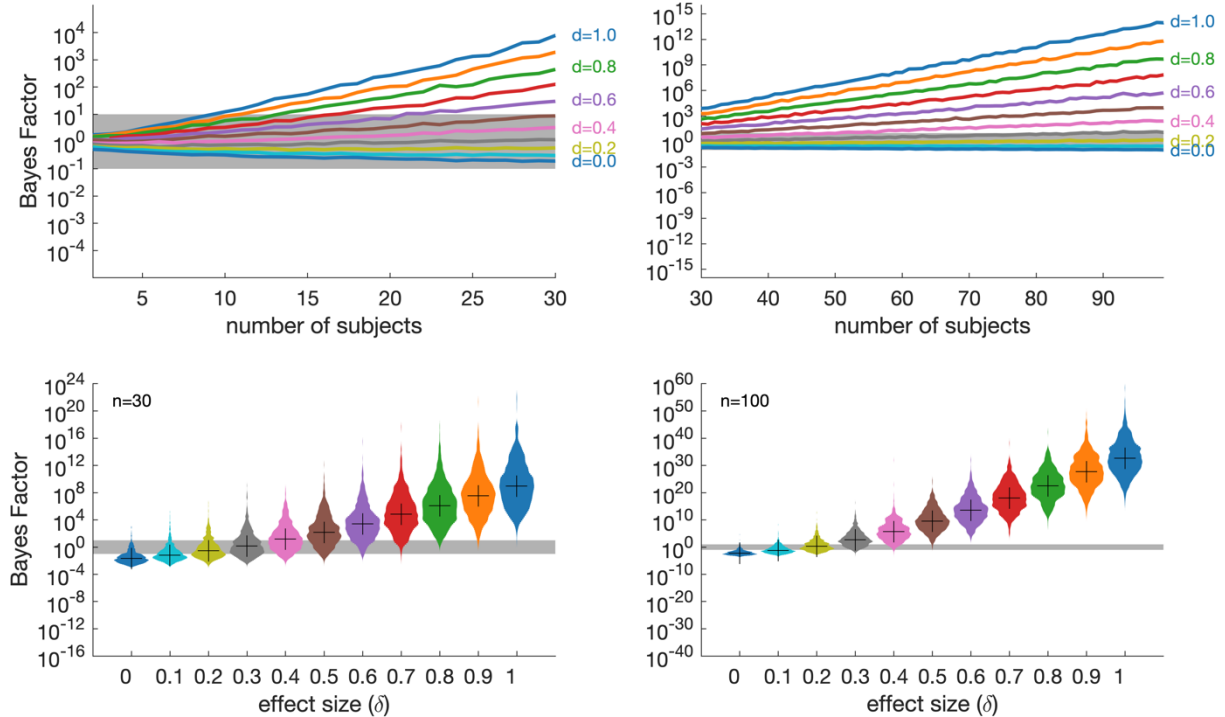321    that more data leads to larger Bayes Factors, but only if there is a true underlying effect.

322

323

324

**Figure 7. Simulated data varying effect sizes and numbers of participants highlight the rationale for using an interval.** We performed 1000 simulations to demonstrate how the Bayes Factors behave with different sample sizes given different effect sizes. A shows Bayes Factors obtained by using a half-Cauchy prior with an interval [0.5 Inf]. B shows Bayes Factors obtained by using a half-Cauchy prior without an interval. The first and third rows show the median Bayes Factors of 1000 simulations as a function of the number of participants. The

333    second and fourth rows show the distribution of the Bayes Factors from 1000 simulations

334    using 30 participants (left panels) and 100 participants (right panels). The distributions of the

335    Bayes Factors highlight the rationale for using an interval, as without an interval it is nearly

336    impossible to find substantial evidence for the null hypothesis even when the effect size equals

337    zero.

338

339

340


# 341    Discussion

342    Bayes Factors have seen a recent increase in popularity in cognitive science, as they can be

343    used to provide quantifiable evidence for contrasting hypotheses. However, their uptake has

344    to date been slow for neuroimaging experiments. To facilitate their adoption, we have provided

345    an empirically-driven guide on implementing Bayes Factors for time-series neuroimaging

346    decoding, using both real and simulated data. We showed that using Bayes Factors and

347    cluster-corrected p-values lead to similar results when statistically assessing time-series

348    neuroimaging decoding results. However, the key advantages of using Bayes Factors are the

349    ability to compare evidence for $H_a$ with evidence for $H_0$ and having results that are quantifiable

350    (e.g., Dienes, 2014; Wagenmakers et al., 2016). Our results show that for time-series

351    decoding data, half-Cauchy priors with default width and an interval ranging from effect sizes

352    of 0.5 to infinity provide sensible results. We also show that even a small number of

353    participants can yield informative Bayes Factors, which can be useful for making decisions on

354    experimental design parameters (e.g., number of trials) during piloting stages of a study.

355

356    Our results showed that the overall conclusions derived from Bayes Factors and p-values

357    were quite similar, highlighting that theoretical considerations should be the deciding factor

358    when choosing a statistical approach to analyse neural time-series data. In the decoding

359    context, p-values afford a dichotomous decision of whether there is enough evidence to reject

360    the hypothesis that decoding is at chance at a given timepoint. Rejecting the null hypothesis

361    is decoupled from any prior beliefs or theories (Dienes, 2011) and is linked to an accepted

362    overall error rate such as $\alpha$ = 0.05. P-values allow us to test for the presence of an effect at a

363    given timepoint using widely accepted thresholds for evidence. While Bayes Factors can in

364    principle also be thresholded to draw dichotomous conclusions, one of the added benefits of

365    Bayes Factors over p-values is the ability to quantify the evidence. Another useful benefit of

366    using Bayes Factors to analyse time-series decoding data is that Bayes Factors allow us to

367    accrue evidence for above-chance as well as at-chance decoding. For time-series analyses
368    in particular, this is a useful feature as the time period prior to stimulus onset can be considered
369    as a control period where we would expect evidence for the null hypothesis. Testing both
370    hypotheses simultaneously can also be a beneficial feature when the research question
371    involves hypotheses predicting certain time-periods without any information in the neural
372    signal (e.g., "X happens before Y" versus "Y happens before X"). Thus, depending on the
373    research question it may be clear which statistical approach suits the time-series decoding
374    analysis best. Otherwise, as overall conclusions do not differ, Bayes Factors and p-values can
375    be used in a complementary way to provide quantifiable evidence for and against the tested
376    hypotheses as well as definitive decisions (see also Lakens et al., 2020; van Dongen et al.,
377    2019; Wagenmakers et al., 2018).

378

379    Through our results, we provide an empirical, straightforward guide to help implement Bayes
380    Factors and demonstrate the extent of practical benefits when using Bayes Factors for time-
381    series neural decoding. Using a data-driven approach, we showed which analysis parameters
382    are most suitable for statistical assessment of time-series decoding data with Bayes Factors.
383    While the Bayes Factors in our example MEG decoding dataset were robust against changes
384    in the predefined width of the prior, defining the prior range so that there is a gap between $H_a$
385    and $H_0$ was critical for finding evidence for the $H_0$. This strong effect of the prior range on the
386    resulting Bayes Factors is particularly relevant in the decoding context, as classification
387    accuracies under the null are not symmetrically distributed around chance (cf. Allefeld et al.,
388    2016). Thus, a gap between $H_0$ and the lower bound of $H_a$ ensures that small above-chance
389    classification accuracies are not treated as evidence for $H_a$. Furthermore, we systematically
390    varied dataset size and showed that using Bayes Factors for time-series decoding data is
391    particularly beneficial when there is limited, noisy data such as in a piloting scenario, as
392    quantifiable evidence for one hypothesis over another gives a stronger sense of whether it is
393    worth pursuing the research question with the piloted design, or make changes (e.g., modify
394    trial numbers or add/remove conditions). Finally, Bayes Factors can be calculated sequentially
395    while evidence accumulation is monitored to stop once a criterion is reached (Dienes, 2011;
396    Rouder, 2014), which can save resources and avoid underpowered studies (Wagenmakers et
397    al., 2018). One possibility is to define a stopping criterion in terms of a percentage of timepoints
398    where evidence is in the conclusive range of Bayes Factors (e.g., 80% of Bayes Factors are
399    above 6 or below 1/6). As longer baselines can artificially increase the percentage of
400    conclusive timepoints, only timepoints after stimulus onset should be considered or the
401    duration of the baseline period should be pre-defined. As researchers generally do not have
402    unlimited resources, it is possible to also pre-register an upper limit for the sample size (e.g.,
403    maximum 50 participants).

404

405     An open question is to what extent our parameter choices generalize to different paradigms,

406     analysis approaches, and modalities. The Bayes Factor parameters used here were optimized

407     for time-series decoding. It is in principle possible to use Bayes Factors in a similar way to

408     analyse other time-series data such as event related potentials, oscillations or regressions,

409     however, the Bayes Factor parameters might have to be adjusted. Similarly, the analysis

410     pipeline discussed here could be extended to other neural decoding modalities such as fMRI

411     (see e.g., Moerel, Rich, et al., 2021). Pilot data or analyses of previous data can be used to

412     examine how parameters have to be modified in order to get sensible results.

413

414     A final consideration is the multiple comparisons problem arising from statistically testing many

415     time points. When using Bayes Factors, as long as the evidence for each hypothesis is

416     interpreted at face value (and not thresholded for 'significance'), we do not need to control for

417     multiple comparisons (Dienes, 2011, 2016a; Świątkowski & Carrier, 2020). That is because

418     once we have established a prior and collected the data, we examine how much we have to

419     adjust our prior beliefs given the data and compare the adjustment required for both

420     hypotheses. This idea is not related to overall error rates and thus does not change if we

421     sample data sequentially or run multiple tests (Dienes, 2016a). If a research question strongly

422     depends on a dichotomous decision on multiple tests, then we advise to report corrected p-

423     values (for which correction methods are well established) alongside the Bayes Factors.

424

425     In conclusion, we have provided an empirically-driven guide on how to use and interpret Bayes

426     Factors for time-series neuroimaging decoding data. We show that Bayes Factors bring

427     several advantages to interpreting time-series decoding results such as quantifiable evidence

428     and an ability to compare evidence for above-chance with evidence for at-chance decoding.

429     We      hope      this      guide,      and      the      accompanying      example      code

430     (https://github.com/LinaTeichmann1/BFF_repo) can serve as a starting point to incorporate

431     Bayesian statistics to existing analysis pipelines.

432

# References

Allefeld, C., Görgen, K., & Haynes, J.-D. (2016). Valid population inference for information-based imaging: From the second-level t-test to prevalence inference. *Neuroimage*, *141*, 378–392.

Dienes, Z. (2011). Bayesian versus orthodox statistics: Which side are you on? *Perspectives on Psychological Science*, *6*(3), 274–290.

Dienes, Z. (2014). Using Bayes to get the most out of non-significant results. *Frontiers in Psychology*, *5*, 781.

Dienes, Z. (2016a). How Bayes factors change scientific practice. *Journal of Mathematical Psychology*, *72*, 78–89.

Dienes, Z. (2016b). How Bayes factors change scientific practice. *Journal of Mathematical Psychology*, *72*, 78–89. https://doi.org/10.1016/j.jmp.2015.10.003

Foster, E. D., & Deardorff, A. (2017). Open Science Framework (OSF). *Journal of the Medical Library Association : JMLA*, *105*(2), 203–206. https://doi.org/10.5195/jmla.2017.88

Grootswagers, T., Robinson, A. K., & Carlson, T. A. (2019a). The representational dynamics of visual objects in rapid serial visual processing streams. *NeuroImage*, *188*, 668–679.

Grootswagers, T., Robinson, A. K., & Carlson, T. A. (2019b). The representational dynamics of visual objects in rapid serial visual processing streams. *NeuroImage*, *188*, 668–679. https://doi.org/10.1016/j.neuroimage.2018.12.046

Grootswagers, T., Robinson, A. K., Shatek, S. M., & Carlson, T. A. (2019). Untangling featural and conceptual object representations. *NeuroImage*, *202*, 116083. https://doi.org/10.1016/j.neuroimage.2019.116083

Grootswagers, T., Robinson, A. K., Shatek, S. M., & Carlson, T. A. (2021). The neural dynamics underlying prioritisation of task-relevant information. *Neurons, Behavior, Data Analysis, and Theory*, *5*(1), 1–17. https://doi.org/10.51628/001c.21174

460   Hebart, M. N., & Baker, C. I. (2018). Deconstructing multivariate decoding for the study of

461         brain function. *Neuroimage*, *180*, 4–18.

462   Jeffreys, H. (1939). The Theory of Probability. *The Theory of Probability*.

463   Jeffreys, H. (1935). Some tests of significance, treated by the theory of probability.

464         *Mathematical Proceedings of the Cambridge Philosophical Society*, *31*(2), 203–222.

465   Kaiser, D., Moeskops, M. M., & Cichy, R. M. (2018). Typical retinotopic locations impact the

466         time course of object coding. *NeuroImage*.

467         https://doi.org/10.1016/j.neuroimage.2018.05.006

468   Karimi-Rouzbahani, H., Woolgar, A., & Rich, A. N. (2021). Neural signatures of vigilance

469         decrements predict behavioural errors before they occur. *ELife*, *10*, e60563.

470   Keysers, C., Gazzola, V., & Wagenmakers, E.-J. (2020). Using Bayes factor hypothesis

471         testing in neuroscience to establish evidence of absence. *Nature Neuroscience*,

472         *23*(7), 788–799.

473   Lakens, D., McLatchie, N., Isager, P. M., Scheel, A. M., & Dienes, Z. (2020). Improving

474         inferences about null effects with Bayes factors and equivalence tests. *The Journals*

475         *of Gerontology: Series B*, *75*(1), 45–57.

476   Ly, A., Verhagen, J., & Wagenmakers, E.-J. (2016). Harold Jeffreys's default Bayes factor

477         hypothesis tests: Explanation, extension, and application in psychology. *Journal of*

478         *Mathematical Psychology*, *72*, 19–32.

479   Mai, A.-T., Grootswagers, T., & Carlson, T. A. (2019). In search of consciousness:

480         Examining the temporal dynamics of conscious visual perception using MEG time-

481         series data. *Neuropsychologia*, *129*, 310–317.

482         https://doi.org/10.1016/j.neuropsychologia.2019.04.015

483   Moerel, D., Grootswagers, T., Robinson, A. K., Shatek, S. M., Woolgar, A., Carlson, T. A., &

484         Rich, A. N. (2021). Undivided attention: The temporal effects of attention dissociated

485         from decision, memory, and expectation. *BioRxiv*, 2021.05.24.445376.

486         https://doi.org/10.1101/2021.05.24.445376

Moerel, D., Rich, A. N., & Woolgar, A. (2021). Selective attention and decision-making have separable neural bases in space and time. *BioRxiv*, 2021.02.28.433294. https://doi.org/10.1101/2021.02.28.433294

Morey, R. D., Romeijn, J.-W., & Rouder, J. N. (2016). The philosophy of Bayes factors and the quantification of statistical evidence. *Journal of Mathematical Psychology*, *72*, 6–18.

Morey, R. D., Rouder, J. N., Jamil, T., & Morey, M. R. D. (2015). Package 'bayesfactor.' *URLh Http://Cran/r-Projectorg/Web/Packages/BayesFactor/BayesFactor Pdf i (Accessed 1006 15)*.

Oosterhof, N. N., Connolly, A. C., & Haxby, J. V. (2016). CoSMoMVPA: Multi-modal multivariate pattern analysis of neuroimaging data in Matlab/GNU Octave. *Frontiers in Neuroinformatics*, *10*. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4956688/

Proklova, D., Kaiser, D., & Peelen, M. V. (2019). MEG sensor patterns reflect perceptual but not categorical similarity of animate and inanimate objects. *NeuroImage*, *193*, 167–177. https://doi.org/10.1016/j.neuroimage.2019.03.028

Robinson, A. K., Grootswagers, T., & Carlson, T. A. (2019). The influence of image masking on object representations during rapid serial visual presentation. *NeuroImage*, *197*, 224–231. https://doi.org/10.1016/j.neuroimage.2019.04.050

Robinson, A. K., Grootswagers, T., Shatek, S. M., Gerboni, J., Holcombe, A., & Carlson, T. A. (2021). Overlapping neural representations for the position of visible and imagined objects. *Neurons, Behavior, Data Analysis, and Theory*, *4*(1), 1–28. https://doi.org/10.51628/001c.19129

Rouder, J. N. (2014). Optional stopping: No problem for Bayesians. *Psychonomic Bulletin & Review*, *21*(2), 301–308.

Rouder, J. N., Speckman, P. L., Sun, D., Morey, R. D., & Iverson, G. (2009). Bayesian t tests for accepting and rejecting the null hypothesis. *Psychonomic Bulletin & Review*, *16*(2), 225–237.

514    Świątkowski, W., & Carrier, A. (2020). There is Nothing Magical about Bayesian Statistics:

515        An Introduction to Epistemic Probabilities in Data Analysis for Psychology Starters.

516        *Basic and Applied Social Psychology*, *42*(6), 387–412.

517    Teichmann, L., Grootswagers, T., Carlson, T., & Rich, A. N. (2018). Decoding digits and dice

518        with magnetoencephalography: Evidence for a shared representation of magnitude.

519        *Journal of Cognitive Neuroscience*, *30*(7), 999–1010.

520    Teichmann, L., Grootswagers, T., Carlson, T., & Rich, A. N. (2019). Seeing versus knowing:

521        The temporal dynamics of real and implied colour processing in the human brain.

522        *NeuroImage*, *200*, 373.

523    Teichmann, L., Quek, G. L., Robinson, A. K., Grootswagers, T., Carlson, T. A., & Rich, A. N.

524        (2020). The influence of object-colour knowledge on emerging object representations

525        in the brain. *Journal of Neuroscience*.

526    van Dongen, N. N., van Doorn, J. B., Gronau, Q. F., van Ravenzwaaij, D., Hoekstra, R.,

527        Haucke, M. N., Lakens, D., Hennig, C., Morey, R. D., & Homer, S. (2019). Multiple

528        perspectives on inference for two simple statistical scenarios. *The American*

529        *Statistician*, *73*(sup1), 328–339.

530    Wagenmakers, E.-J., Marsman, M., Jamil, T., Ly, A., Verhagen, J., Love, J., Selker, R.,

531        Gronau, Q. F., Šmíra, M., & Epskamp, S. (2018). Bayesian inference for psychology.

532        Part I: Theoretical advantages and practical ramifications. *Psychonomic Bulletin &*

533        *Review*, *25*(1), 35–57.

534    Wagenmakers, E.-J., Morey, R. D., & Lee, M. D. (2016). Bayesian benefits for the pragmatic

535        researcher. *Current Directions in Psychological Science*, *25*(3), 169–176.

536