

Hypothesis Testing Using Bayes Factor in Behavioral Sciences

Timo B. Roettger¹ and Michael Franke²

¹Department of Linguistics & Scandinavian Studies, University of Oslo

²Department of Linguistics, University of Tübingen

Abstract

Recent times have seen a surge of Bayesian inference across the behavioral sciences. However, the process of testing hypothesis is often conceptually challenging or computationally costly. This tutorial provides an accessible, non-technical introduction that covers the most common scenarios in experimental sciences: Testing the evidence for an alternative hypothesis using Bayes Factor through the Savage Dickey approximation. This method is conceptually easy to understand and computationally cheap.

Keywords: statistics, Bayes, Bayes Factor, Savage Dickey, hypothesis testing, ROPE

1 Introduction

To date, the most common quantitative approach across the experimental sciences is to run an experiment with one or more predictors and statistically test whether the predictors affect the measured variables. Traditionally, these statistical tests have been done within the null hypothesis significance testing framework. Over the last decade or so, however, we have seen more and more statistical approaches within an alternative inferential framework: Bayesian inference. Testing hypothesis within the Bayesian framework is often considered either conceptually challenging, computationally too costly, or both. This tutorial provides an accessible, non-technical introduction to Bayesian hypothesis testing that is easy to understand and computationally cheap.

2 Motivation and intended audience

This tutorial provides a very basic introduction to the topic using R (R Core Team, 2025). We wrote this tutorial with a particular reader in mind. If you have used R before and if you have a basic understanding of linear regression, and Bayesian inference, this tutorial is for you. We will remain mostly conceptual to provide you with a conceptual tool to approach hypothesis testing within Bayesian inference. The form of hypothesis testing that we would like to introduce to you is, however, different from the traditional null hypothesis significance testing in that it requires more thinking about the quantitative nature of your data. This is not a bug but, at least for us, a feature that will allow you to understand both your data and what you can learn from them better.

If you don't have any experience with regression modeling, you will probably still be able to follow, but you might also want to consider doing a crash course. To bring you up to speed, we recommend the excellent tutorial by Bodo Winter (2013) on mixed effects regression in a non-Bayesian —a.k.a. frequentist— paradigm. To then make the transition to Bayesian versions of these regression models, we shamelessly suggest our own tutorial on “Bayesian Regression for Factorial Designs” as a natural follow-up using the same data and Winter (Franke & Roettger, 2019). In a sense, the present tutorial on hypothesis testing could be considered the long-awaited sequel of the series started by Winter. For continuity, we will continue to use the original data set.

To actively follow this tutorial, you should have R installed on your computer (<https://www.r-project.org>). Unless you already have a favorite editor for tinkering with R scripts, we recommend to try out RStudio (<https://www.rstudio.com>). You will also need some packages, which you can import with the following code:

```
# package for convenience functions (e.g. plotting)
library(tidyverse)
library(ggdist)

# package for Bayesian regression modeling
library(brms)

# package for posterior wrangling and plotting
library(tidybayes)

# package for BF calculation and plotting
library(bayestestR)
```

3 Data, research questions & hypotheses

This tutorial looks at a data set relevant for investigating whether voice pitch differs across social contexts in Korean. Korean is a language in which the social distance between speakers plays a central role. The way Korean speakers speak depends for example on whether they are in a formal context (e.g. during a consultation with a professor) or an informal context (e.g. chatting with a friend about the holidays) (Winter & Grawunder, 2012).

To load the data into your R environment, run the following code

```
# TO DO: STORE ONLINE
# TO DO: SIMPLIFY STORED DATA
polite = read_csv("../data/polite.csv") |>
  # remove men
  filter(gender == "F") |>
  # transform context to factor
  mutate(context = as.factor(context))

polite

# A tibble: 126 x 4
  subject gender context pitch
  <chr>    <chr>  <fct>   <dbl>
```

```

1 F1      F      formal  215.
2 F1      F      informal 211.
3 F1      F      formal  285.
4 F1      F      informal 266.
5 F1      F      formal  211.
6 F1      F      informal 286.
7 F1      F      formal  252.
8 F1      F      informal 282.
9 F1      F      formal  230.
10 F1     F      informal 250.
# i 116 more rows

```

This data set contains anonymous identifiers for individual speakers stored in the variable `subject`. In this tutorial we will only be looking at female speakers btw. Subjects produced different sentences, and the experiment manipulated whether the sentences were produced in a formal or an informal social context, indicated by the variable `context`. Crucially, each row contains a measurement of pitch in Hz stored in the variable `pitch`.

For most analyses of behavioral experiments, researchers are interested in whether an outcome variable is meaningfully affected by at least one manipulated variable and if so how the outcome variable is affected by it. In this case, Winter and Grawunder (2012) wanted to test whether pitch is meaningfully affected by the social context of the utterance.

As a first step, we can explore this question visually:

Figure 1

Empirical distribution of pitch values across contexts

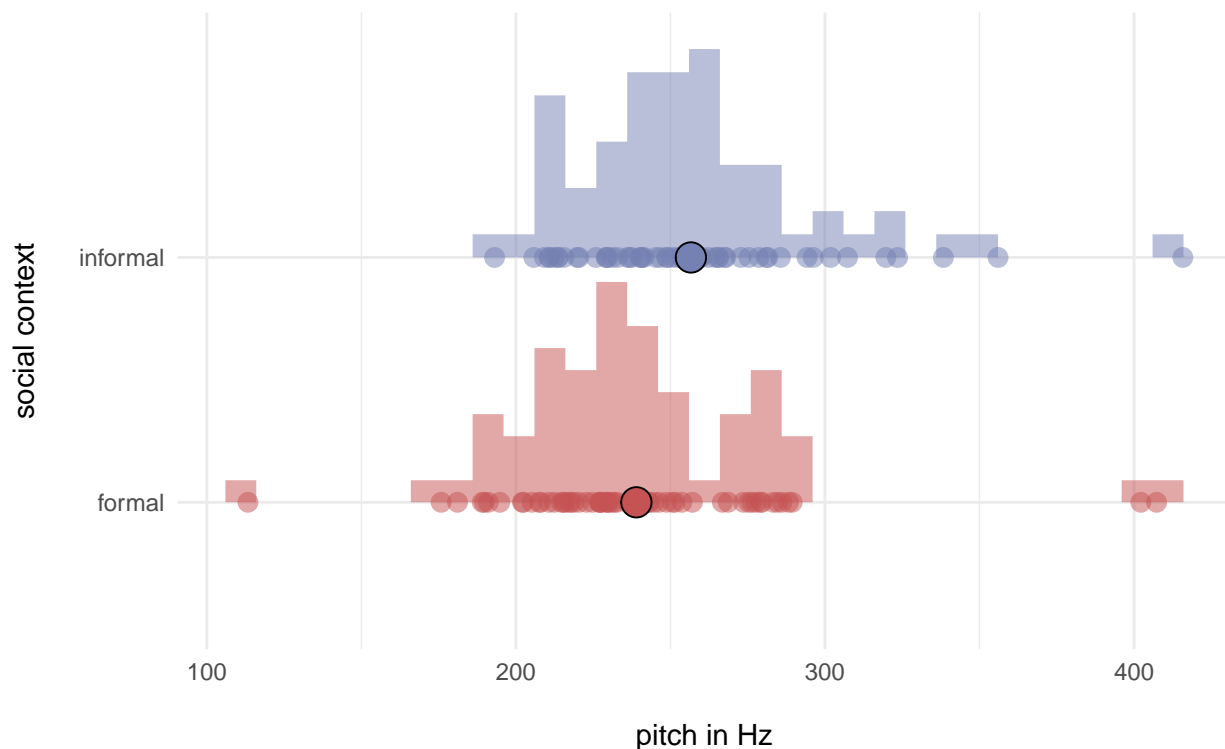


Figure 1 displays the pitch values for all utterances in the dataset across contexts (semi-transparent points). The solid points indicate the average pitch values across all sentences and speakers. Looking at the plot, we can see that voice pitch from utterances in formal contexts are on average slightly lower than those in informal contexts. The red distribution is slightly shifted to the left of the blue distribution by around 1.3 semitones. In other words, speakers tend to slightly lower their voice pitch when speaking in a formal context. But there is also a lot of overlap between the two contexts. Now as Bayesians, we would like to translate the data into an expression of evidence: does the data provide evidence for our research hypotheses?

Let us build a Bayesian linear model to approach an answer to this question. Our first step is to specify the model formula and check which priors need to be specified:

```
# contrast code predictor

contrasts(polite$context) <- c(-0.5,0.5)

# define linear model formula
# predict pitch by context and allow for that relationship
# to vary between subjects
formula <- bf(pitch ~ context + (1 + context | subject))

# get priors for this model
get_prior(formula, polite)
```

The default priors that brms picks for the Intercept and the variance parameters are mostly reasonable as they are derived from the data, weakly informative and symmetrical. However the prior for our critical parameter `context1` should also be weakly informative (Gelman et al., 2017), i.e. the prior assumption about the difference between informal and formal contexts should be that we don't know, but our best guess is that it is close to zero and equally likely to be more or less than zero. So we specify a normal distribution centered on zero for this parameter.

Note: Only for demonstration purposes, we will use default priors for the other parameters, but you always should critically reflect on all of your priors.

```
# pick a weakly informative prior for the critical parameter
priors <- prior(normal(0, 20),
               class = b,
               coef = "context1")
```

Now we do a so-called prior predictive check, in other words we want to know what the posterior distribution looks like before having seen the data, based on the priors only. This is a useful exercise to make sure that the priors results in reasonable quantitative assumptions. We usually do it for all parameters, but here we will focus only on the critical parameter `context1`, i.e. the difference between formal and informal contexts. Let us also have a look at the predictions for the prior-only model.

```
# NOTE: CAN WE STORE THE SAMPLING PARAMETERS (seed, iter, chains, cores, backend, data)?
#       TO MAKE THE CODE CHUNKS SMALLER?

# run the model
fit_prior <- brm(formula,
                 prior = priors,
```

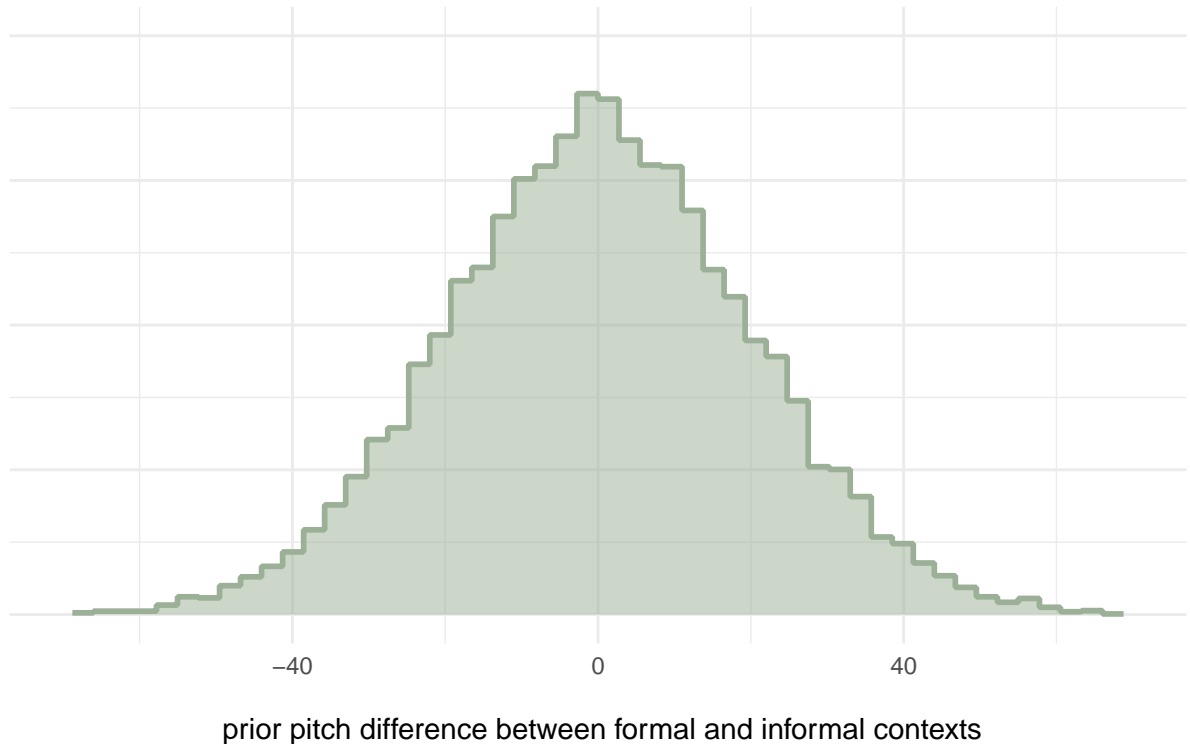
```
family = gaussian(),
# sample prior only
sample_prior = "only",
# store / load model output
file = "../models/fit_prior",
# common sampling specifications
seed = 1234,
iter = 8000,
chains = 4,
cores = 4,
backend = "cmdstanr",
data = polite)
```

```
# extract prior samples
prior_samples <-
  fit_prior |>
  spread_draws(b_context1)

# plot
ggplot(prior_samples,
       aes(x = b_context1)) +
  stat_histinterval(slab_color = project_colors[11],
                   slab_fill = alpha(project_colors[11], 0.5),
                   fill = NA,
                   color = NA,
                   outline_bars = FALSE) +
  labs(x = "\n prior pitch difference between formal and informal contexts",
       y = "") +
  scale_x_continuous(limits = c(-70, 70)) +
  theme_minimal() +
  theme(axis.text.y = element_blank())
```

Figure 2

Prior probability of the effect of context on pitch, i.e. before seeing the data



Looking at the distribution, the priors for the effect of context on pitch seems sensible. The most plausible value is zero. Values that are smaller or larger than zero become less plausible the further they are away from zero and values being smaller or larger than zero are equally likely. Good. Before we have seen the data, our model is somewhat pessimistic about the effect of context on on pitch. Now we can run the full model that integrates the likelihood (our data) with the priors and visualize the posteriors for the critical parameter.

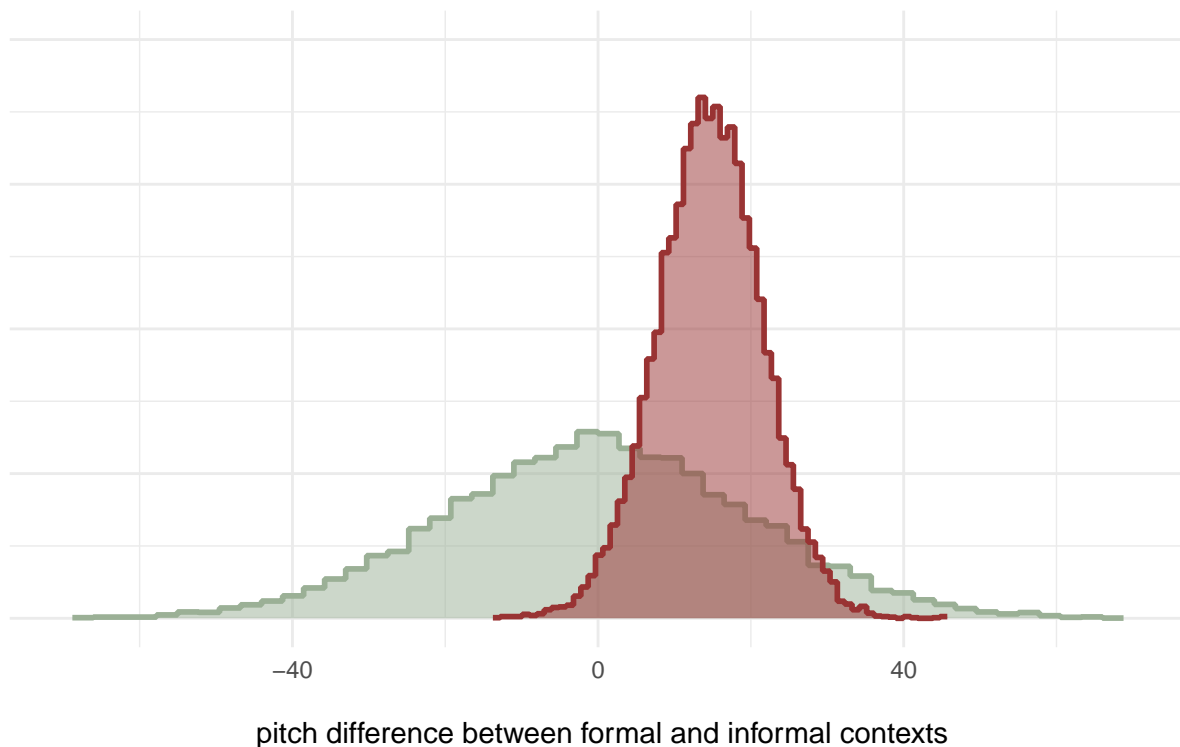
```
# run the model
fit <- brm(formula,
  prior = priors,
  family = gaussian(),
  # store / load model output
  file = "../models/fit",
  # common sampling specifications
  seed = 1234,
  iter = 8000,
  chains = 4,
  cores = 4,
  backend = "cmdstanr",
  data = polite)
```

```
posterior_plot <- fit |>
  spread_draws(b_context1) |>
  ggplot(aes(x = b_context1)) +
    stat_histinterval(data = prior_samples,
                      slab_color = project_colors[11],
                      slab_fill = alpha(project_colors[11], 0.5),
                      fill = NA,
                      color = NA,
                      outline_bars = FALSE) +
    stat_histinterval(slab_color = project_colors[14],
                      slab_fill = alpha(project_colors[14], 0.5),
                      color = NA,
                      outline_bars = FALSE) +
    scale_thickness_shared() +
    labs(x = "\n pitch difference between formal and informal contexts",
         y = "") +
    scale_x_continuous(limits = c(-70, 70)) +
    theme_minimal() +
    theme(axis.text.y = element_blank())

posterior_plot
```

Figure 3

Posterior probability of the effect of context on pitch, i.e. after seeing the data



The posterior samples (red distribution) suggests that the majority of plausible values after seeing the data are positive, or in other words, informal contexts elicit larger pitch values. Negative values are not very plausible posterior values, but also not completely implausible. Compared to our prior probability (green distribution) for which roughly 50% of posteriors are negative, this decrease in plausibility of negative values is quite noteworthy already.

What we have done here should be quite familiar. We compare our model predictions to a reference point. It is a single point value: zero. But do we really care that much for such point hypotheses? Is zero really that special? We might think so because years of using null hypothesis significance testing has conditioned us to think that way. But this tutorial would like to break this cycle and move forward. Bear with us and let's approach hypothesis testing a bit differently today.

3.1 Grounding hypotheses in regions of practical equivalences

Above we claimed that we wanted to test “whether pitch is **meaningfully affected** by the social context of the utterance”. We snuck the word meaningfully in there for a reason. But what does “meaningful” mean? This is really a good question and (un)fortunately requires quite a bit of thinking. This tutorial deals with speech data. Speech is, in spoken languages at least, THE vehicle to transmit linguistic information in order to communicate with each other. Speech is also very complex and very noisy: Not everything that can be measured in the acoustic signal matters for the listener. For example, if something cannot be perceived reliably, it is at least conceivable that it might play little to no role in communication. While speech sciences has a rich research tradition to estimate what can and what cannot be reliably heard, exact estimates depends on a lot of moving parts. Such thresholds are referred to as Just Noticeable Differences (JNDs) and can be used to define what constitutes meaningful differences when we look at speech data.

For example, Liu (2013) report on JNDs ranging from 3 to 14 Hz. Jongman et al. (2017) report on JNDs between 6 and 9 Hz. Turner et al. (2019) reported on JNDs between 17 and 25 Hz for non-speech stimuli and between 35 and 40 Hz for speech stimuli. While these studies are hard to compare, they give us at least a the range of JND values to work with.

So we could interpret the original hypothesis the following way: If a pitch difference is below the JND, it is not meaningful. So instead of testing against a point-zero hypothesis, we can test against a range of parameter values that are equivalent to the null value for practical purposes. In our case, let us begin with the lowest reported JND of the above studies on pitch perception in speech (3 Hz), but be extra conservative and double the reported JND to 6 Hz. We then assume that pitch values between -6 and 6 are meaningless. Such ranges are sometimes called regions of practical equivalence (ROPEs), range of equivalence, equivalence margin, smallest effect size of interest, or good-enough belt (see Kruschke, 2018).

```
rope <- c(-6, 6)
```

With a ROPE being defined, we can now test our hypothesis “whether pitch is **meaningfully affected** by the social context of the utterance” using Bayes Factor:

4 Testing hypothesis using Bayes Factor

4.1 What is Bayes Factor

Bayes Factors (henceforth: BFs) allow us to quantify relative evidence of one model compared to another.

4.2 Approximating Bayes Factor with Savage Dickey

4.3 Calculating Bayes Factor for a specified Region of Practical Equivalence (ROPE)

Instead of doing it by hand, we can calculate the Savage Dickey ratio with the `bayesfactor_parameters()` function from the `bayesfactorR` package. What happens behind the scenes is that the function will sample posteriors from your specified model based on priors only (so before seeing any data) and calculates the posterior probability of the specified null hypothesis (here the range specified by our ROPE).

```
#|warning: FALSE
#|message: FALSE

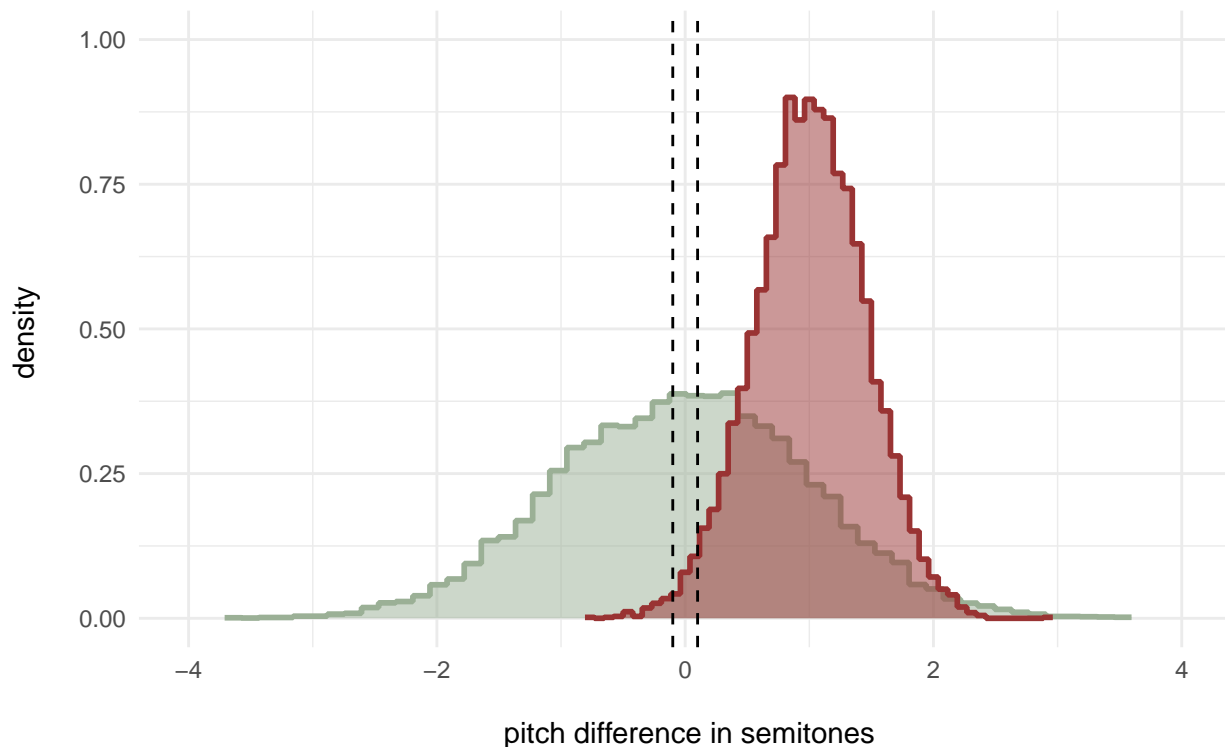
BF_1 <- bayesfactor_parameters(posterior = fit,
                              null = rope,
                              parameter = "b_context1")
```

Before interpreting the number we get, let us visually explore what our BF corresponds to.

```
posterior_plot +
  geom_vline(xintercept = c(rope[1], rope[2]),
            lty = "dashed")
```

Figure 4

Prior and posterior probability of the effect of context on pitch relative to the ROPE (-0.1, 0.1)



What the BF does is relating two numbers: (a) The prior probability of parameter values outside the rope, i.e. the proportion of the green distribution that falls outside the dashed lines, and (b) the posterior probability of parameter values outside the rope, i.e. the proportion of the red distribution that falls outside the dashed lines. Eye-balling the plot, we can maybe already see that more of the red distribution is outside the ROPE than of the green distribution.

BF_1

Bayes Factor (Null-Interval)

Parameter	BF
contextinformal	5.76

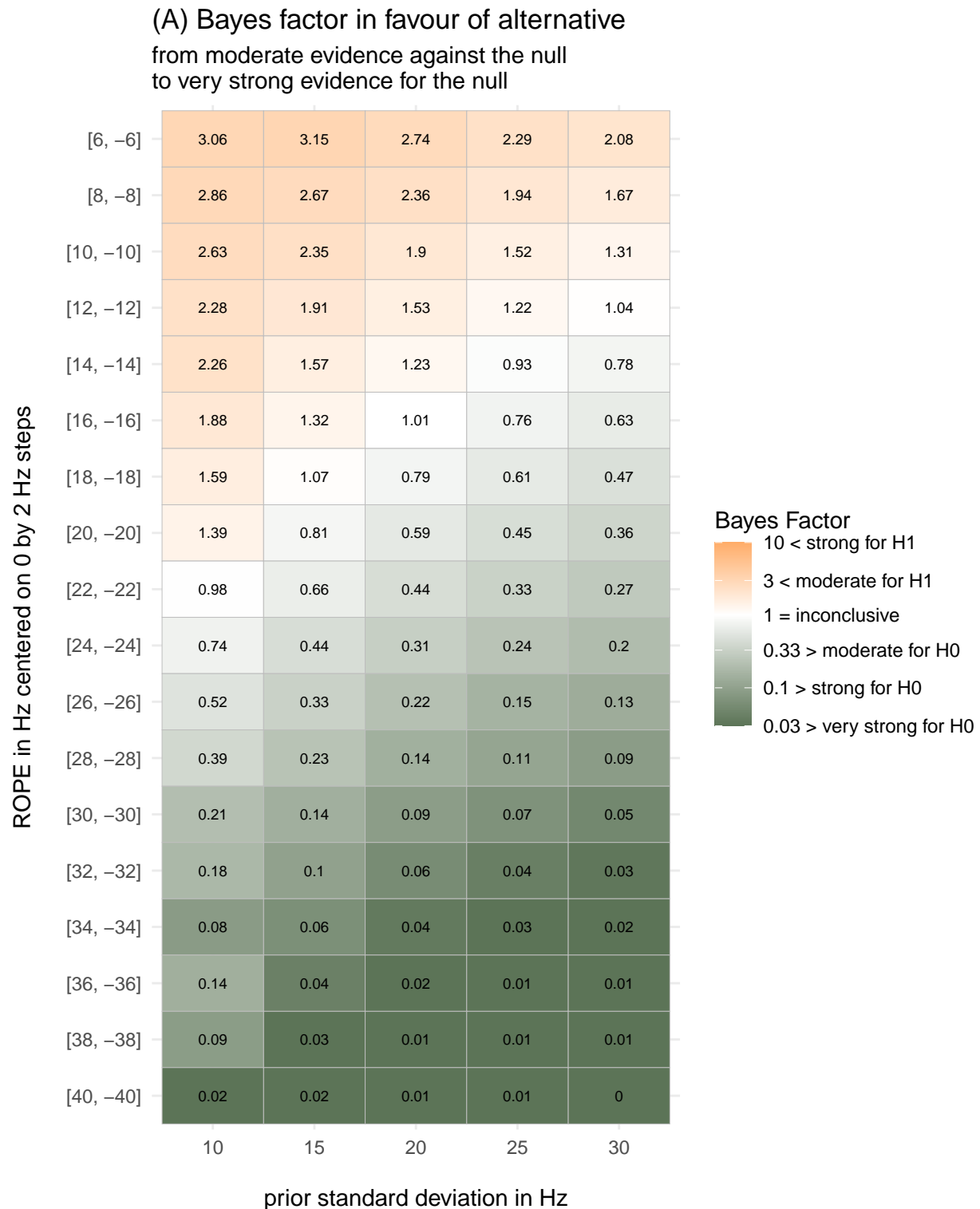
* Evidence Against The Null: [-0.100, 0.100]

To be exact, 2.7 times for of the red distribution is outside of the ROPE than of the green distribution.

That means the model that has seen the data provide 2.7 times more evidence for pitch being outside of the ROPE, or in other words, it is 2.7 times more likely (after having seen the data), that context affects pitch meaningfully. According to Lee and Wagenmakers (2014) criteria for interpreting BFs, this value corresponds to only anecdotal evidence for the alternative hypothesis.

4.4 Sensitivity analysis for different priors and ROPEs

Now as you probably have guessed already, all these probabilities are very much dependent on the priors of the model, so it is important to evaluate the robustness of our Bayes Factor-based interpretation across a range of sensible priors. And as long as we are not a 100% sure about what a meaningful difference is, we might as well explore the robustness of the Bayes Factor across different ROPEs. We won't bore you with the code for that process, but you can follow it along in our scripts. Let us explore the following ROPE intervals as informed by the three studies cited above on pitch perception: we test a range of ROPE intervals from 6 Hz to 40 Hz. We also assume the following five prior values for the width of the standard deviation of the critical parameter (centered on zero): 10, 15, 20, 25, 30. These are all sensible prior widths assuming that medium to strong effects in either direction are plausible.

Figure 5*Bayes Factors for a range of priors and a range of ROPEs*

The combination of Bayes Factors is visualized in Figure X. Orange cells indicate evidence for the

alternative. Green cells indicate evidence for the null. It becomes clear that the conclusions we can draw from our data are rather dependent on the choices we made along the way.

By comparing the Bayes Factors along the y-axis, we can see that they are heavily dependent on the chosen ROPE. We here chose (theoretically speaking) a quite large range of ROPEs, all of which are informed by psychoacoustic studies of what pitch differences can be reliably heard and thus likely are meaningful for communication. In light of this range of possible definitions what constitutes meaningful differences, our data seem not very robust, as illustrated by the shift from orange to green. Even the smallest ROPE intervals provide only anecdotal to moderate evidence for the alternative. And the most conservative ROPEs, following Turner et al. (2019), leads to moderate to very strong evidence against the alternative hypothesis.

Additionally, when comparing the Bayes Factors along the x-axis, we can see that they are comparatively consistent for different standard deviations of the critical prior. However, we can also see that the Bayes Factors decrease with the width of the priors (from left to right). This is not surprising and a known phenomenon, often discussed under the Jeffreys-Lindley paradox (Lindley, 1957): The more diffuse the priors are (i.e. wider priors), the larger is the probability that a specific parameter values is not compatible with the data.

Combined, we can see that the larger the ROPE and the wider the priors, the more likely becomes the null hypothesis. In an ideal world, the evidence provided by the data should be robust across these choices. However, this exploration of our inference is a fantastic opportunity to assess the boundaries of our conclusions. In this case, the original conclusions by Winter and Grawunder (2012) was based on the null hypothesis significance testing and traditionally tested the compatibility of the data with a point-null hypothesis. They concluded “that in formal speech, Korean [...] female speakers lowered their average fundamental frequency [...]”. This statement is still true according to their inferential criteria, but thinking more deeply about the theoretical consequences of differences in pitch, it might be less clear that these differences are truly meaningful.

4.5 BF for point hypothesis

don't lol

5 How to write things up

6 Some words of encouragement

Bayesian inference in general and this form of hypothesis testing in particular require much more thinking than we might be used to. We think this is a good thing. Many voices have criticized the lack of engagement that we behavioral scientists invest into thinking how our theoretical ideas connect to concrete predictions in the quantitative systems under investigation (Coretta et al., 2023; e.g. Scheel, 2022; Woensdregt et al., 2024). The presented form of hypothesis testing is easy to understand, but does require to think deeply about prior quantitative assumptions as well as what it means for observations to be meaningfully different. That is neither trivial nor easy. But we would like to encourage you to engage in exactly this thinking to better understand our data and how they might link with our understanding of cognition and behavior.

7 Other Resources

There are many fantastic resources out there to help you learn about the wonderful world of statistics. Here are a few recommendations. - A very accessible introduction to linear models in R is Winter (2019). - ...

8 References

R version 4.4.3 (2025-02-28)

Platform: x86_64-apple-darwin20

Running under: macOS Sequoia 15.5

Matrix products: default

BLAS: /Library/Frameworks/R.framework/Versions/4.4-x86_64/Resources/lib/libRblas.0.dylib

LAPACK: /Library/Frameworks/R.framework/Versions/4.4-x86_64/Resources/lib/libRlapack.dylib; L

locale:

[1] en_US.UTF-8/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8

time zone: Europe/Oslo

tzcode source: internal

attached base packages:

[1] stats graphics grDevices utils datasets methods base

other attached packages:

[1] ggdist_3.3.2	rstan_2.32.6	StanHeaders_2.32.10
[4] see_0.11.0	bayestestR_0.16.1.2	tidybayes_3.0.6
[7] brms_2.22.0	Rcpp_1.1.0	ggribes_0.5.6
[10] lubridate_1.9.4	forcats_1.0.0	stringr_1.5.1
[13] dplyr_1.1.4	purrr_1.1.0	readr_2.1.5
[16] tidyr_1.3.1	tibble_3.3.0	ggplot2_3.5.1
[19] tidyverse_2.0.0		

loaded via a namespace (and not attached):

[1] svUnit_1.0.6	tidyselect_1.2.1	farver_2.1.2
[4] loo_2.8.0	fastmap_1.2.0	tensorA_0.36.2.1
[7] digest_0.6.36	timechange_0.3.0	lifecycle_1.0.4
[10] processx_3.8.4	magrittr_2.0.3	posterior_1.6.0
[13] compiler_4.4.3	rlang_1.1.6	tools_4.4.3
[16] utf8_1.2.6	yaml_2.3.10	data.table_1.17.8
[19] knitr_1.48	labeling_0.4.3	bridgesampling_1.1-2
[22] bit_4.6.0	pkgbuild_1.4.4	cmdstanr_0.8.1
[25] abind_1.4-5	withr_3.0.2	datawizard_1.2.0
[28] grid_4.4.3	stats4_4.4.3	colorspace_2.1-1
[31] inline_0.3.19	scales_1.3.0	insight_1.3.1.14
[34] cli_3.6.5	mvtnorm_1.3-1	rmarkdown_2.27
[37] crayon_1.5.3	generics_0.1.4	RcppParallel_5.1.8
[40] rstudioapi_0.16.0	tzdb_0.5.0	bayesplot_1.11.1
[43] parallel_4.4.3	matrixStats_1.3.0	vctr_0.6.5
[46] Matrix_1.7-2	jsonlite_2.0.0	hms_1.1.3
[49] arrayhelpers_1.1-0	bit64_4.6.0-1	logspline_2.1.22
[52] glue_1.8.0	codetools_0.2-20	ps_1.7.7
[55] distributional_0.4.0	stringi_1.8.7	gtable_0.3.6

[58] QuickJSR_1.3.0	munsell_0.5.1	pillar_1.11.0
[61] htmltools_0.5.8.1	Broddingnag_1.2-9	R6_2.6.1
[64] vroom_1.6.5	evaluate_0.24.0	lattice_0.22-6
[67] backports_1.5.0	rstantools_2.4.0	coda_0.19-4.1
[70] gridExtra_2.3	nlme_3.1-167	checkmate_2.3.1
[73] xfun_0.45	pkgconfig_2.0.3	

[[1]]

Bürkner P (2017). "brms: An R Package for Bayesian Multilevel Models Using Stan." *_Journal of Statistical Software_*, *80*(1), 1-28. doi:10.18637/jss.v080.i01 <<https://doi.org/10.18637/jss.v080.i01>>.

Bürkner P (2018). "Advanced Bayesian Multilevel Modeling with the R Package brms." *_The R Journal_*, *10*(1), 395-411. doi:10.32614/RJ-2018-017 <<https://doi.org/10.32614/RJ-2018-017>>.

Bürkner P (2021). "Bayesian Item Response Modeling in R with brms and Stan." *_Journal of Statistical Software_*, *100*(5), 1-54. doi:10.18637/jss.v100.i05 <<https://doi.org/10.18637/jss.v100.i05>>.

[[2]]

Makowski D, Ben-Shachar M, Lüdtke D (2019). "bayestestR: Describing Effects and their Uncertainty, Existence and Significance within the Bayesian Framework." *_Journal of Open Source Software_*, *4*(40), 1541. doi:10.21105/joss.01541 <<https://doi.org/10.21105/joss.01541>>, <<https://joss.theoj.org/papers/10.21105/joss.01541>>.

[[3]]

Kay M (2023). *_tidybayes: Tidy Data and Geoms for Bayesian Models_*. doi:10.5281/zenodo.1308151 <<https://doi.org/10.5281/zenodo.1308151>>, R package version 3.0.6, <<http://mjskay.github.io/tidybayes/>>.

[[4]]

Wickham H, Averick M, Bryan J, Chang W, McGowan LD, François R, Golemund G, Hayes A, Henry L, Hester J, Kuhn M, Pedersen TL, Miller E, Bache SM, Müller K, Ooms J, Robinson D, Seidel DP, Spinu V, Takahashi K, Vaughan D, Wilke C, Woo K, Yutani H (2019). "Welcome to the tidyverse." *_Journal of Open Source Software_*, *4*(43), 1686. doi:10.21105/joss.01686 <<https://doi.org/10.21105/joss.01686>>.

[[5]]

Kay M (2024). "ggdist: Visualizations of Distributions and Uncertainty in the Grammar of Graphics." *_IEEE Transactions on Visualization and Computer Graphics_*, *30*(1), 414-424. doi:10.1109/TVCG.2023.3327195 <<https://doi.org/10.1109/TVCG.2023.3327195>>.

Kay M (2024). *_ggdist: Visualizations of Distributions and Uncertainty_*. doi:10.5281/zenodo.3879620

<<https://doi.org/10.5281/zenodo.3879620>>, R package version 3.3.2,
 <<https://mjskay.github.io/ggdist/>>.

- Coretta, S., Casillas, J. V., Roessig, S., Franke, M., Ahn, B., Al-Hoorie, A. H., Al-Tamimi, J., Alotaibi, N. E., AlShakhori, M. K., Altmiller, R. M., et al. (2023). Multidimensional signals and analytic flexibility: Estimating degrees of freedom in human-speech analyses. *Advances in Methods and Practices in Psychological Science*, 6(3), 25152459231162567.
- Franke, M., & Roettger, T. (2019). *Bayesian regression modeling (for factorial designs): A tutorial*. OSF. <https://doi.org/10.31234/osf.io/cdxv3>
- Gelman, A., Simpson, D., & Betancourt, M. (2017). The Prior Can Often Only Be Understood in the Context of the Likelihood. *Entropy*, 19(10), 555. <https://doi.org/10.3390/e19100555>
- Jongman, A., Qin, Z., Zhang, J., & Sereno, J. A. (2017). Just noticeable differences for pitch direction, height, and slope for mandarin and english listeners. *The Journal of the Acoustical Society of America*, 142(2), EL163–EL169.
- Kruschke, J. K. (2018). Rejecting or Accepting Parameter Values in Bayesian Estimation. *Advances in Methods and Practices in Psychological Science*, 1(2), 270–280. <https://doi.org/10.1177/2515245918771304>
- Lee, M. D., & Wagenmakers, E.-J. (2014). *Bayesian Cognitive Modeling: A Practical Course*. Cambridge University Press.
- Lindley, D. V. (1957). A Statistical Paradox. *Biometrika*, 44(1/2), 187–192. <https://doi.org/10.2307/2333251>
- Liu, C. (2013). Just noticeable difference of tone pitch contour change for english-and chinese-native listeners. *The Journal of the Acoustical Society of America*, 134(4), 3011–3020.
- Scheel, A. M. (2022). Why most psychological research findings are not even wrong. *Infant and Child Development*, 31(1), e2295.
- Turner, D. R., Bradlow, A. R., & Cole, J. S. (2019). Perception of pitch contours in speech and nonspeech. *INTERSPEECH*, 2275–2279.
- Winter, B. (2013). *Linear models and linear mixed effects models in R with linguistic applications* (arXiv:1308.5499). arXiv. <https://doi.org/10.48550/arXiv.1308.5499>
- Winter, B. (2019). *Statistics for linguists: An introduction using r*. Routledge.
- Winter, B., & Grawunder, S. (2012). The phonetic profile of Korean formal and informal speech registers. *Journal of Phonetics*, 40(6), 808–815. <https://doi.org/10.1016/j.wocn.2012.08.006>
- Woensdregt, M., Fusaroli, R., Rich, P., Modrák, M., Kolokolova, A., Wright, C., & Warlaumont, A. S. (2024). Lessons for theory from scientific domains where evidence is sparse or indirect. *Computational Brain & Behavior*, 7(4), 588–607.