

1 An Introduction to Bayesian Multilevel Models Using brms: A Case Study of Gender Effects
2 on Vowel Variability in Standard Indonesian

3 Ladislav Nalborczyk^{1,2}, Cédric Batailler³, Hélène Løevenbruck¹, Anne Vilain^{4,5}, &
4 Paul-Christian Bürkner⁶

5 ¹ Univ. Grenoble Alpes, CNRS, LPNC, 38000, Grenoble, France

6 ² Department of Experimental Clinical and Health Psychology, Ghent University, Ghent,
7 Belgium

8 ³ Univ. Grenoble Alpes, LIP/PC2S, 38000, Grenoble, France

9 ⁴ Univ. Grenoble Alpes, CNRS, Grenoble INP, GIPSA-lab, 38000, Grenoble, France

10 ⁵ Institut Universitaire de France, France

11 ⁶ Department of Computer Science, Aalto University, Finland

12 **Disclosure:** The authors have declared that no competing interests existed at the time of
13 publication. **Funding:** The first author of the manuscript is funded by a fellowship from
14 Univ. Grenoble Alpes.

15 Author Note

16 Correspondence concerning this article should be addressed to Ladislav Nalborczyk,
17 Laboratoire de Psychologie et Neurocognition, Univ. Grenoble Alpes, 1251 avenue centrale,
18 38058 Grenoble Cedex 9, France. E-mail: ladislav.nalborczyk@univ-grenoble-alpes.fr

Abstract

19

20 **Purpose:** Bayesian multilevel models are increasingly used to overcome the limitations of
21 frequentist approaches in the analysis of complex structured data. This paper introduces
22 Bayesian multilevel modelling for the specific analysis of speech data, using the brms
23 package developed in R. **Method:** In this tutorial, we provide a practical introduction to
24 Bayesian multilevel modelling, by reanalysing a phonetic dataset containing formant (F1 and
25 F2) values for five vowels of Standard Indonesian (ISO 639-3:ind), as spoken by eight
26 speakers (four females), with several repetitions of each vowel. **Results:** We first give an
27 introductory overview of the Bayesian framework and multilevel modelling. We then show
28 how Bayesian multilevel models can be fitted using the probabilistic programming language
29 Stan and the R package brms, which provides an intuitive formula syntax. **Conclusions:**
30 Through this tutorial, we demonstrate some of the advantages of the Bayesian framework for
31 statistical modelling and provide a detailed case study, with complete source code for full
32 reproducibility of the analyses (<https://osf.io/dpzcb/>).

33

Keywords: Bayesian data analysis, multilevel models, mixed models, brms, Stan

- 34 An Introduction to Bayesian Multilevel Models Using brms: A Case Study of Gender Effects
35 on Vowel Variability in Standard Indonesian
36 wordcount (excluding abstract, references, tables and figures): 8637

1 Introduction

37

38 The last decade has witnessed noticeable changes in the way experimental data are
39 analysed in phonetics, psycholinguistics, and speech sciences in general. In particular, there
40 has been a shift from analysis of variance (ANOVA) to *linear mixed models*, also known as
41 *hierarchical models* or *multilevel models* (MLMs), spurred by the spreading use of
42 data-oriented programming languages such as R (R Core Team, 2017), and by the
43 enthusiasm of its active and ever growing community. This shift has been further sustained
44 by the current transition in data analysis in social sciences, with researchers evolving from a
45 widely criticised point-hypothesis mechanical testing (e.g., Bakan, 1966; Gigerenzer, Krauss,
46 & Vitouch, 2004; Kline, 2004; Lambdin, 2012; Trafimow et al., 2018) to an approach that
47 emphasises parameter estimation, model comparison, and continuous model expansion (e.g.,
48 Cumming, 2012, 2014; Gelman et al., 2013; Gelman & Hill, 2007; Kruschke, 2015; Kruschke
49 & Liddell, 2017a, 2017b; McElreath, 2016).

50 MLMs offer great flexibility in the sense that they can model statistical phenomena
51 that occur on different levels. This is done by fitting models that include both constant and
52 varying effects (sometimes referred to as *fixed* and *random* effects). Among other advantages,
53 this makes it possible to generalise the results to unobserved levels of the *groups* existing in
54 the data (e.g., stimulus or participant, Janssen, 2012). The multilevel strategy can be
55 especially useful when dealing with repeated measurements (e.g., when measurements are
56 nested into participants) or with unequal sample sizes, and more generally, when handling
57 complex dependency structures in the data. Such complexities are frequently found in the
58 kind of experimental designs used in speech science studies, for which MLMs are therefore
59 particularly well suited.

60 The standard MLM is usually fitted in a frequentist framework, with the `lme4` package
61 (Bates et al., 2015b) in R (R Core Team, 2017). However, when one tries to include the
62 maximal varying effect structure, this kind of model tends either not to converge, or to give

63 aberrant estimations of the correlation between varying effects (e.g., Bates et al., 2015a)¹.
64 Yet, fitting the maximal varying effect structure has been explicitly recommended (e.g., Barr,
65 Levy, Scheepers, & Tily, 2013). In contrast, the maximal varying effect structure can
66 generally be fitted in a Bayesian framework (Bates et al., 2015a; Eager & Roy, 2017;
67 Nicenboim & Vasishth, 2016; Sorensen, Hohenstein, & Vasishth, 2016).

68 Another advantage of Bayesian statistical modelling is that it fits the way researchers
69 intuitively understand statistical results. Widespread misinterpretations of frequentist
70 statistics (like p-values and confidence intervals) are often attributable to the wrong
71 interpretation of these statistics as resulting from a Bayesian analysis (e.g., Dienes, 2011;
72 Gigerenzer et al., 2004; Hoekstra, Morey, Rouder, & Wagenmakers, 2014; Kruschke &
73 Liddell, 2017a; Morey, Hoekstra, Rouder, Lee, & Wagenmakers, 2015). However, the
74 intuitive nature of the Bayesian approach might arguably be hidden by the predominance of
75 frequentist teaching in undergraduate statistical courses.

76 Moreover, the Bayesian approach offers a natural solution to the problem of multiple
77 comparisons, when the situation is adequately modelled in a multilevel framework (Gelman,
78 Hill, & Yajima, 2012; Scott & Berger, 2010), and allows *a priori* knowledge to be
79 incorporated in data analysis via the prior distribution. The latter feature is particularly
80 relevant when dealing with constraint parameters or for the purpose of incorporating expert
81 knowledge.

82 The aim of the current paper is to introduce Bayesian multilevel models, and to
83 provide an accessible and illustrated hands-on tutorial for analysing typical phonetic data.
84 This paper will be structured in two main parts. First, we will briefly introduce the Bayesian
85 approach to data analysis and the multilevel modelling strategy. Second, we will illustrate
86 how Bayesian MLMs can be implemented in R by using the `brms` package (Bürkner, 2017b)
87 to reanalyse a dataset from McCloy (2014) available in the `phonR` package (McCloy, 2016).

¹ In this context, the *maximal varying effect structure* means that any potential source of systematic influence should be explicitly modelled, by adding appropriate varying effects.

88 We will fit Bayesian MLMs of increasing complexity, going step by step, providing
 89 explanatory figures and making use of the tools available in the `brms` package for model
 90 checking and model comparison. We will then compare the results obtained in a Bayesian
 91 framework using `brms` with the results obtained using frequentist MLMs fitted with `lme4`.
 92 Throughout the paper, we will also provide comments and recommendations about the
 93 feasibility and the relevance of such analysis for the researcher in speech sciences.

94 1.1 Bayesian data analysis

95 The Bayesian approach to data analysis differs from the frequentist one in that each
 96 parameter of the model is considered as a random variable (contrary to the frequentist
 97 approach which considers parameter values as unknown and fixed quantities), and by the
 98 explicit use of probability to model the uncertainty (Gelman et al., 2013). The two
 99 approaches also differ in their conception of what *probability* is. In the Bayesian framework,
 100 probability refers to the experience of uncertainty, while in the frequentist framework it
 101 refers to the limit of a relative frequency (i.e., the relative frequency of an event when the
 102 number of trials approaches infinity). A direct consequence of these two differences is that
 103 Bayesian data analysis allows researchers to discuss the probability of a parameter (or a
 104 vector of parameters) θ , given a set of data y :

$$p(\theta|y) = \frac{p(y|\theta)p(\theta)}{p(y)}$$

105 Using this equation (known as Bayes' theorem), a probability distribution $p(\theta|y)$ can
 106 be derived (called the *posterior distribution*), that reflects knowledge about the parameter,
 107 given the data and the prior information. This distribution is the goal of any Bayesian
 108 analysis and contains all the information needed for inference.

109 The term $p(\theta)$ corresponds to the *prior distribution*, which specifies the prior
 110 information about the parameters (i.e., what is known about θ before observing the data) as
 111 a probability distribution. The left hand of the numerator $p(y|\theta)$ represents the *likelihood*,

112 also called the *sampling distribution* or *generative model*, and is the function through which
113 the data affect the posterior distribution. The likelihood function indicates how likely the
114 data are to appear, for each possible value of θ .

115 Finally, $p(y)$ is called the *marginal likelihood*. It is meant to normalise the posterior
116 distribution, that is, to scale it in the “probability world”. It gives the “probability of the
117 data”, summing over all values of θ and is described by $p(y) = \sum_{\theta} p(\theta)p(y|\theta)$ for discrete
118 parameters, and by $p(y) = \int p(\theta)p(y|\theta)d\theta$ in the case of continuous parameters.

119 All this pieced together shows that the result of a Bayesian analysis, namely the
120 posterior distribution $p(\theta|y)$, is given by the product of the information contained in the data
121 (i.e., the likelihood) and the information available before observing the data (i.e., the prior).
122 This constitutes the crucial principle of Bayesian inference, which can be seen as an updating
123 mechanism (as detailed for instance in Kruschke & Liddell, 2017a). To sum up, Bayes’
124 theorem allows a prior state of knowledge to be updated to a posterior state of knowledge,
125 which represents a compromise between the prior knowledge and the empirical evidence.

126 The process of Bayesian analysis usually involves three steps that begin with setting up
127 a probability model for all the entities at hand, then computing the posterior distribution,
128 and finally evaluating the fit and the relevance of the model (Gelman et al., 2013). In the
129 context of linear regression, for instance, the first step would require to specify a likelihood
130 function for the data and a prior distribution for each parameter of interest (e.g., the
131 intercept or the slope). We will go through these three steps in more details in the
132 application section, but we will first give a brief overview of the multilevel modelling strategy.

133 1.2 Multilevel modelling

134 MLMs can be considered as “multilevel” for at least two reasons. First, an MLM can
135 generally be conceived as a regression model in which the parameters are themselves
136 modelled as outcomes of another regression model. The parameters of this second-level
137 regression are known as *hyperparameters*, and are also estimated from the data (Gelman &

138 Hill, 2007). Second, the multilevel structure can arise from the data itself, for instance when
139 one tries to model the second-language speech intelligibility of a child, who is considered
140 within a particular class, itself considered within a particular school. In such cases, the
141 hierarchical structure of the data itself calls for hierarchical modelling. In both conceptions,
142 the number of levels that can be handled by MLMs is virtually unlimited (McElreath, 2016).
143 When we use the term *multilevel* in the following, we will refer to the structure of the model,
144 rather than to the structure of the data, as non-nested data can also be modelled in a
145 multilevel framework.

146 As briefly mentioned earlier, MLMs offer several advantages compared to single-level
147 regression models, as they can handle the dependency between units of analysis from the
148 same group (e.g., several observations from the same participant). In other words, they can
149 account for the fact that, for instance, several observations are not independent, as they
150 relate to the same participant. This is achieved by partitioning the total variance into
151 variation due to the groups (level-2) and to the individual (level-1). As a result, such models
152 provide an estimation of the variance component for the second level (i.e., the variability of
153 the participant-specific estimates) or higher levels, which can inform us about the
154 generalisability of the findings (Janssen, 2012; McElreath, 2016).

155 Multilevel modelling allows both *fixed* and *random* effects to be incorporated. However,
156 as pointed out by Gelman (2005), we can find at least five different (and sometimes
157 contradictory) ways of defining the meaning of the terms *fixed* and *random* effects. Moreover,
158 Gelman and Hill (2007) remarked that what is usually called a *fixed* effect can generally be
159 conceived as a *random* effect with a null variance. In order to use a consistent vocabulary,
160 we follow the recommendations of Gelman and Hill (2007) and avoid these terms. We
161 instead use the more explicit terms *constant* and *varying* to designate effects that are
162 constant, or that vary by groups².

² Note that MLMs are sometimes called *mixed models*, as models that comprise both *fixed* and *random* effects.

163 A question one is frequently faced with in multilevel modelling is to know which
 164 parameters should be considered as varying, and which parameters should be considered as
 165 constant. A practical answer is provided by McElreath (2016), who states that “any batch of
 166 parameters with *exchangeable* index values can be and probably should be pooled”. For
 167 instance, if we are interested in the categorisation of native versus non-native phonemes and
 168 if for each phoneme in each category there are multiple audio stimuli (e.g., multiple
 169 repetitions of the same phoneme), and if we do not have any reason to think that, for each
 170 phoneme, audio stimuli may differ in intelligibility in any systematic way, then repetitions of
 171 the same phoneme should be pooled together. The essential feature of this strategy is that
 172 *exchangeability* of the lower units (i.e., the multiple repetitions of the same phoneme) is
 173 achieved by conditioning on indicator variables (i.e., the phonemes) that represent groupings
 174 in the population (Gelman et al., 2013).

175 To sum up, multilevel models are useful as soon as there are predictors at different
 176 levels of variation (Gelman et al., 2013). One important aspect is that this
 177 varying-coefficients approach allows each subgroup to have a different mean outcome level,
 178 while still estimating the global mean outcome level. In an MLM, these two estimations
 179 inform each other in a way that leads to the phenomenon of *shrinkage*, that will be discussed
 180 in more detail below (see section 2.3).

181 As an illustration, we will build an MLM starting from the ordinary linear regression
 182 model, and trying to predict an outcome y_i (e.g., second-language (L2) speech-intelligibility)
 183 by a linear combination of an intercept α and a slope β that quantifies the influence of a
 184 predictor x_i (e.g., the number of lessons received in this second language):

$$y_i \sim \text{Normal}(\mu_i, \sigma_e)$$

$$\mu_i = \alpha + \beta x_i$$

185 This notation is strictly equivalent to the (maybe more usual) following notation:

$$y_i = \alpha + \beta x_i + \epsilon_i$$

$$\epsilon_i \sim \text{Normal}(0, \sigma_e)$$

186 We prefer to use the first notation as it generalises better to more complex models, as
 187 we will see later. In Bayesian terms, these two lines describe the *likelihood* of the model,
 188 which is the assumption made about the generative process from which the data is issued.
 189 We make the assumption that the outcomes y_i are normally distributed around a mean μ_i
 190 with some error σ_e . This is equivalent to saying that the errors are normally distributed
 191 around 0, as illustrated by the above equivalence. Then, we can extend this model to the
 192 following multilevel model, adding a varying intercept:

$$y_i \sim \text{Normal}(\mu_i, \sigma_e)$$

$$\mu_i = \alpha_{j[i]} + \beta x_i$$

$$\alpha_j \sim \text{Normal}(\alpha, \sigma_\alpha)$$

193 where we use the notation $\alpha_{j[i]}$ to indicate that each group j (e.g., class) is given a
 194 unique intercept, issued from a Gaussian distribution centered on α , the grand intercept³,
 195 meaning that there might be different mean scores for each class. From this notation we can
 196 see that in addition to the residual standard deviation σ_e , we are now estimating one more
 197 variance component σ_α , which is the standard deviation of the distribution of varying
 198 intercepts. We can interpret the variation of the parameter α between groups j by
 199 considering the *intra-class correlation* (ICC) $\sigma_\alpha^2 / (\sigma_\alpha^2 + \sigma_e^2)$, which goes to 0, if the grouping
 200 conveys no information, and to 1, if all observations in a group are identical (Gelman & Hill,
 201 2007, p. 258).

202 The third line is called a *prior* distribution in the Bayesian framework. This prior

³ Acknowledging that these individual intercepts can also be seen as adjustments to the grand intercept α , that are specific to group j .

203 distribution describes the population of intercepts, thus modelling the dependency between
204 these parameters.

205 Following the same strategy, we can add a varying slope, allowed to vary according to
206 the group j :

$$y_i \sim \text{Normal}(\mu_i, \sigma_e)$$

$$\mu_i = \alpha_{j[i]} + \beta_{j[i]}x_i$$

$$\alpha_j \sim \text{Normal}(\alpha, \sigma_\alpha)$$

$$\beta_j \sim \text{Normal}(\beta, \sigma_\beta)$$

207 Indicating that the effect of the number of lessons on L2 speech intelligibility is allowed
208 to differ from one class to another (i.e., the effect of the number of lessons might be more
209 beneficial to some classes than others). These varying slopes are assigned a prior distribution
210 centered on the grand slope β , and with standard deviation σ_β .

211 In this introductory section, we have presented the foundations of Bayesian analysis
212 and multilevel modelling. Bayes' theorem allows prior knowledge about parameters to be
213 updated according to the information conveyed by the data, while MLMs allow complex
214 dependency structures to be modelled. We now move to a detailed case study in order to
215 illustrate these concepts.

Box 1. Where are my random effects ?

In the Bayesian framework, every unknown quantity is considered as a random variable that we can describe using probability distributions. As a consequence, there is no such thing as a "fixed effect" or a "random effects distribution" in a Bayesian framework. However, these semantic quarrels disappear when we write down the model.

Suppose we have a dependent continuous variable y and a dichotomic categorical predictor x (assumed to be contrast-coded). Let y_{ij} denote the score of the i^{th} participant in the j^{th} condition. We can write a "mixed effects" model (as containing both fixed and random effects) as follows:

$$y_{ij} = \alpha + \alpha_i + \beta x_j + e_{ij}, \quad e_{ij} \sim \text{Normal}(0, \sigma_e^2), \quad \alpha_i \sim \text{Normal}(0, \sigma_a^2)$$

Where the terms α and β represent the "fixed effects" and denote the overall mean response and the condition difference in response, respectively. In addition, e_{ij} are random errors assumed to be normally distributed with unknown variance σ_e^2 , and α_i 's are individual specific random effects normally distributed in the population with unknown variance σ_a^2 .

We can rewrite this model to make apparent that the so-called "random effects distribution" can actually be considered a prior distribution (from a Bayesian standpoint), since by definition, distributions on unknown quantities are considered as priors:

$$y_{ij} \sim \text{Normal}(\mu_{ij}, \sigma_e^2)$$

$$\mu_{ij} = \alpha + \beta x_j$$

$$\alpha_i \sim \text{Normal}(\alpha, \sigma_a^2)$$

where the parameters of this prior are learned from the data. As we have seen, the same mathematical entity can be conceived either as a "random effects distribution" or as a prior distribution, depending on the framework.

217 1.3 Software programs

218 Sorensen et al. (2016) provided a detailed and accessible introduction to Bayesian
 219 MLMs (BMLMs) applied to linguistics, using the probabilistic language **Stan** (Stan
 220 Development Team, 2016). However, discovering BMLMs and the **Stan** language all at once
 221 might seem a little overwhelming, as **Stan** can be difficult to learn for users that are not
 222 experienced with programming languages. As an alternative, we introduce the **brms** package
 223 (Bürkner, 2017b), that implements BMLMs in **R**, using **Stan** under the hood, with an
 224 **lme4**-like syntax. Hence, the syntax required by **brms** will not surprise the researcher
 225 familiar with **lme4**, as models of the following form:

$$y_i \sim \text{Normal}(\mu_i, \sigma_e)$$

$$\mu_i = \alpha + \alpha_{\text{subject}[i]} + \beta x_i$$

226 are specified in **brms** (as in **lme4**) with: `y ~ 1 + x + (1|subject)`. In addition to
 227 linear regression models, **brms** allows generalised linear and non-linear multilevel models to
 228 be fitted, and comes with a great variety of distribution and link functions. For instance,
 229 **brms** allows fitting robust linear regression models, or modelling dichotomous and categorical
 230 outcomes using logistic and ordinal regression models. The flexibility of **brms** also allows for
 231 distributional models (i.e., models that include simultaneous predictions of all response
 232 parameters), Gaussian processes or non-linear models to be fitted, among others. More
 233 information about the diversity of models that can be fitted with **brms** and their
 234 implementation is provided in Bürkner (2017b) and Bürkner (2017a).

235 2 Application example

236 To illustrate the use of BMLMs, we reanalysed a dataset from McCloy (2014), available
 237 in the **phonR** package (McCloy, 2016). This dataset contains formant (F1 and F2) values for
 238 five vowels of Standard Indonesian (ISO 639-3:ind), as spoken by eight speakers (four

239 females), with approximately 45 repetitions of each vowel. The research question we
240 investigated here is the effect of gender on vowel production variability.

241 2.1 Data pre-processing

242 Our research question was about the different amount of variability in the respective
243 vowel productions of male and female speakers, due to cognitive or social differences. To
244 answer this question, we first needed to get rid of the differences in vowel production that
245 are due to physiological differences between males and females (e.g., shorter vocal tract
246 length for females). More generally, we needed to eliminate the inter-individual differences
247 due to physiological characteristics in our groups of participants. For that purpose, we first
248 applied the Watt & Fabricius formant normalisation technique (Watt & Fabricius, 2002).
249 The principle of this method is to calculate for each speaker a “centre of gravity” S in the
250 F1/F2 plane, from the formant values of point vowels [i, a, u], and to express the formant
251 values of each observation as ratios of the value of S for that formant.

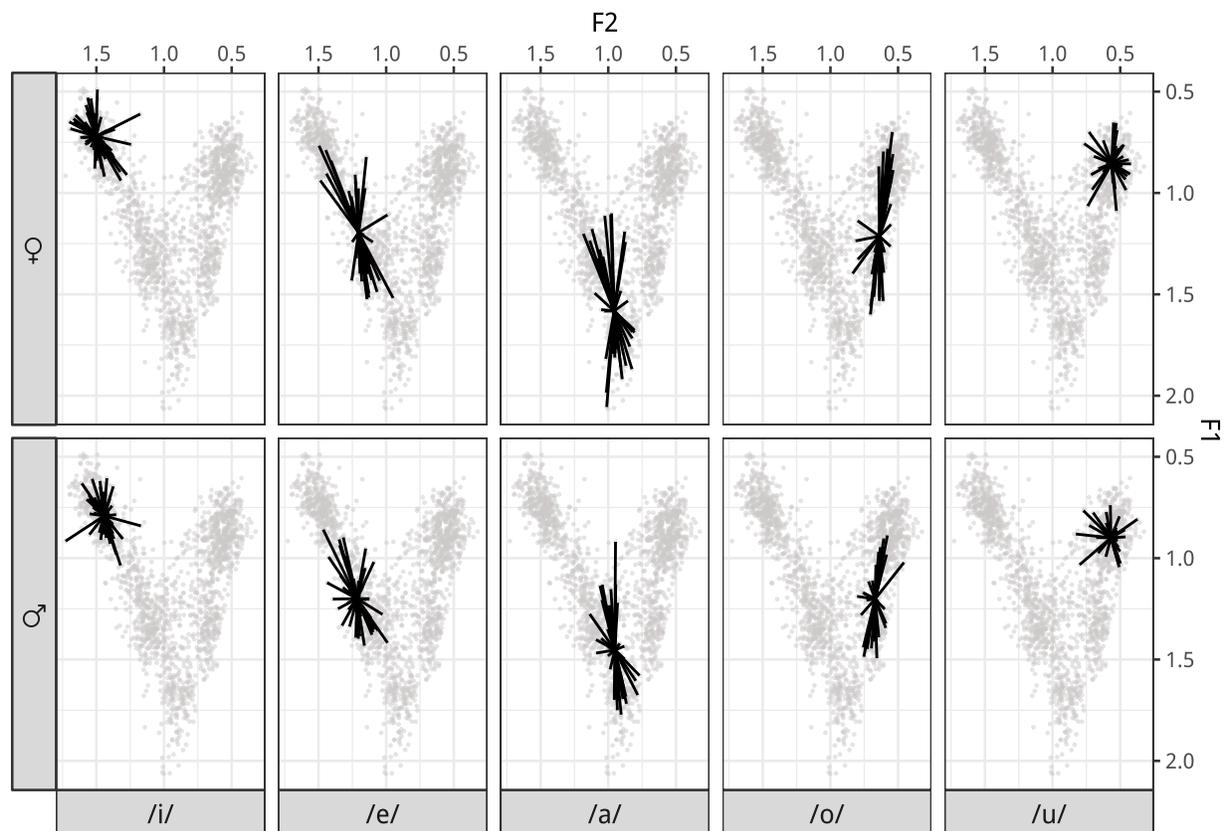


Figure 1. Euclidean distances between each observation and the centres of gravity corresponding to each vowel across all participants, by gender (top row: female, bottom row: male) and by vowel (in column), in the normalised F1-F2 plane. The grey background plots represent the individual data collapsed for all individuals (male and female) and all vowels. Note that, for the sake of clarity, this figure represents a unique center of gravity for each vowel for all participants, whereas in the analysis, one center of gravity was used for each vowel and each participant.

252 Then, for each vowel and participant, we computed the Euclidean distance between
 253 each observation and the centre of gravity of the whole set of observations in the F1-F2
 254 plane for that participant and that vowel. The data obtained by this process are illustrated
 255 in Figure 1, and a sample of the final dataset can be found in Table 1.

Table 1

Ten randomly picked rows from the data.

subj	gender	vowel	f1	f2	f1norm	f2norm	distance	repetition
M02	m	/e/	534	1724	1.143	1.113	0.118	7
F09	f	/i/	468	2401	0.943	1.447	0.223	16
F04	f	/a/	885	1413	1.636	0.804	0.223	12
M01	m	/a/	671	1262	1.615	0.823	0.176	25
F04	f	/a/	700	1951	1.294	1.109	0.237	36
F04	f	/e/	614	2100	1.135	1.194	0.070	42
M04	m	/i/	338	2163	0.803	1.432	0.040	16
F04	f	/o/	649	1357	1.200	0.772	0.154	12
M04	m	/a/	524	1573	1.245	1.041	0.146	20
M02	m	/u/	411	762	0.879	0.492	0.134	25

2.2 Constant effect of gender on vowel production variability

We then built a first model with constant effects only and vague priors on α and β , the intercept and the slope. We contrast-coded **gender** (f = -0.5, m = 0.5). Our dependent variable was therefore the distance from each individual vowel centre of gravity, which we will refer to as *formant distance* in the following. The formal model can be expressed as:

$$\text{distance}_i \sim \text{Normal}(\mu_i, \sigma_e)$$

$$\mu_i = \alpha + \beta \times \text{gender}_i$$

$$\alpha \sim \text{Normal}(0, 10)$$

$$\beta \sim \text{Normal}(0, 10)$$

$$\sigma_e \sim \text{HalfCauchy}(10)$$

261 where the first two lines of the model describe the likelihood and the linear model⁴.
 262 The next three lines define the prior distribution for each parameter of the model, where α
 263 and β are given a vague (weakly informative) Gaussian prior centered on 0, and the residual
 264 variation is given a Half-Cauchy prior (Gelman, 2006; Polson & Scott, 2012), thus restricting
 265 the range of possible values to positive ones. As depicted in Figure 2, the Normal(0, 10)
 266 prior is weakly informative in the sense that it grants a relative high weight to α and β
 267 values, between -25 and 25. This corresponds to very large (given the scale of our data)
 268 values for, respectively, the mean distance value α , and the mean difference between males
 269 and females β . The HalfCauchy(10) prior placed on σ_e also allows very large values of σ_e , as
 270 represented in the right panel of Figure 2.

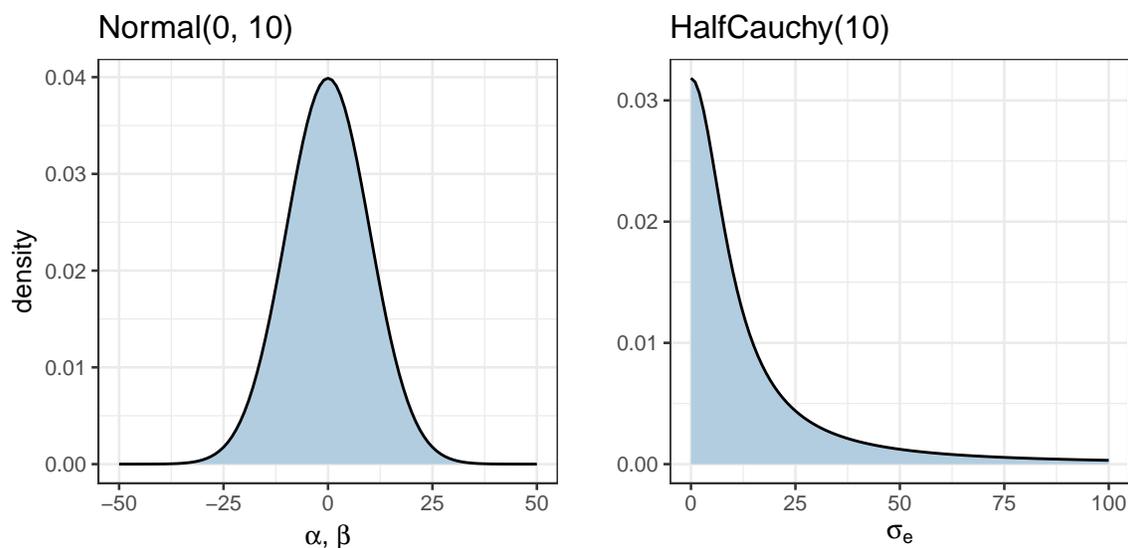


Figure 2. Prior distributions used in the first model, for α and β (left panel) and for the residual variation σ_e (right panel).

271 These priors can be specified in numerous ways (see `?set_prior` for more details),
 272 among which the following:

⁴ Note that –for the sake of simplicity– throughout this tutorial we use a Normal likelihood, but other (better) alternatives would include using skew-normal or log-normal models, which are implemented in **brms** with the `skew_normal` and `lognormal` families. We provide examples in the supplementary materials.

```
prior1 <- c(  
  prior(normal(0, 10), class = Intercept),  
  prior(normal(0, 10), class = b, coef = gender),  
  prior(cauchy(0, 10), class = sigma)  
)
```

273 where a prior can be defined over a class of parameters (e.g., for all variance
274 components, using the `sd` class) or for a specific one, for instance as above by specifying the
275 coefficient (`coef`) to which the prior corresponds (here the slope of the constant effect of
276 gender).

277 The model can be fitted with `brms` with the following command:

```
library(brms)  
  
bmod1 <- brm(  
  distance ~ gender,  
  data = indo, family = gaussian(),  
  prior = prior1,  
  warmup = 2000, iter = 5000  
)
```

278 where `distance` is the distance from the centre of gravity. The `iter` argument serves
279 to specify the total number of iterations of the Markov Chain Monte Carlo (MCMC)
280 algorithm, and the `warmup` argument specifies the number of iterations that are run at the
281 beginning of the process to “calibrate” the MCMC, so that only `iter - warmup` iterations
282 are retained in the end to approximate the shape of the posterior distribution (for more
283 details, see McElreath, 2016).

284 Figure 3 depicts the estimations of this first model for the intercept α , the slope β , and
285 the residual standard deviation σ_e . The left part of the plot shows histograms of draws taken

286 from the posterior distribution, and from which several summaries can be computed (e.g.,
287 mean, mode, quantiles). The right part of Figure 3 shows the behaviour of the two
288 simulations (i.e., the two chains) used to approximate the posterior distribution, where the
289 x-axis represents the number of iterations and the y-axis the value of the parameter. This
290 plot reveals one important aspect of the simulations that should be checked, known as
291 *mixing*. A chain is considered well mixed if it explores many different values for the target
292 parameters and does not stay in the same region of the parameter space. This feature can be
293 evaluated by checking that these plots, usually referred to as *trace plots*, show random
294 scatter around a mean value (they look like a “fat hairy caterpillar”).

```
library(tidyverse)

bmod1 %>%
  plot(
    combo = c("hist", "trace"), widths = c(1, 1.5),
    theme = theme_bw(base_size = 10)
  )
```

295 The estimations obtained for this first model are summarised in Table 2, which
296 includes the mean, the standard error (SE), and the lower and upper bounds of the 95%
297 credible interval (CrI)⁵ of the posterior distribution for each parameter. As **gender** was
298 contrast-coded before the analysis (f = -0.5, m = 0.5), the intercept α corresponds to the
299 grand mean of the formant distance over all participants and has its mean around 0.16. The
300 estimate of the slope ($\beta = -0.04$) suggests that females are more variable than males in the
301 way they pronounce vowels, while the 95% CrI can be interpreted in a way that there is a

⁵ Where a credible interval is the Bayesian analogue of a classical confidence interval, except that probability statements can be made based upon it (e.g., “given the data and our prior assumptions, there is a 0.95 probability that this interval encompasses the population value θ ”).

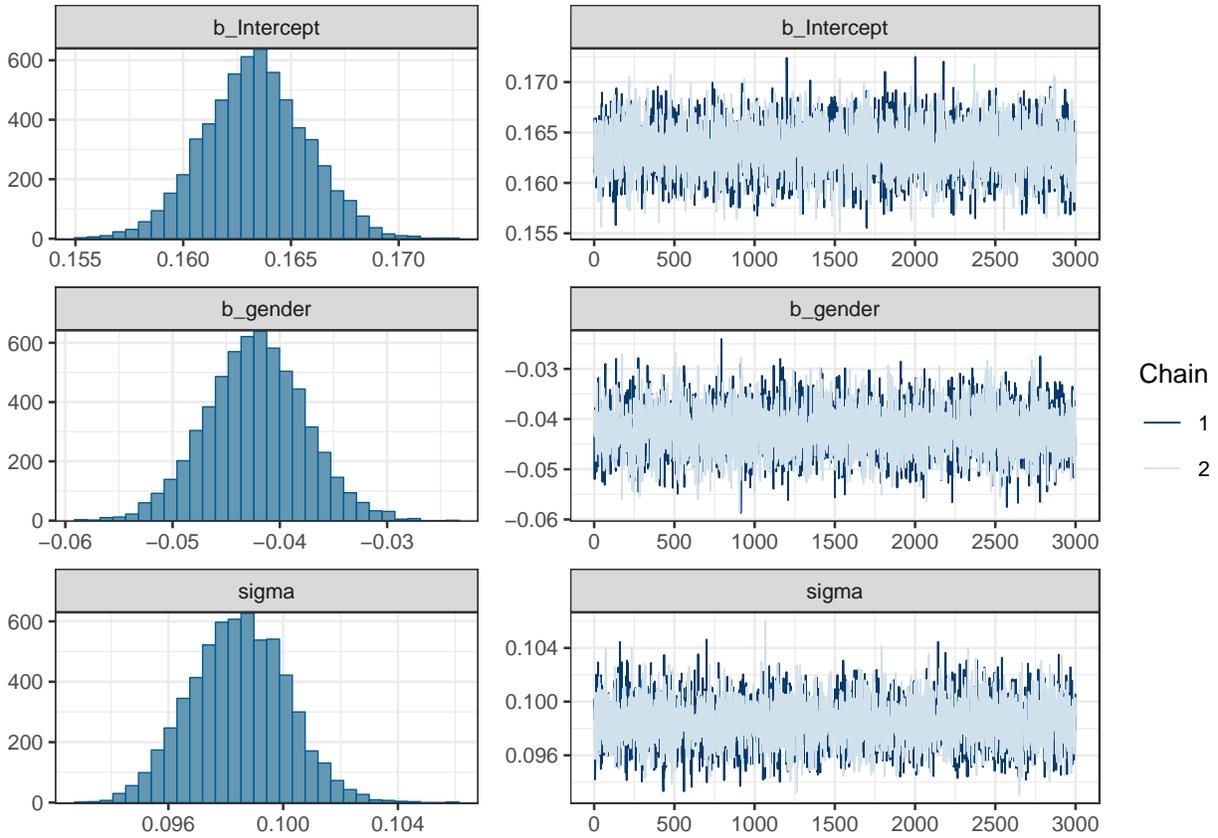


Figure 3. Histograms of posterior samples and trace plots of the intercept, the slope for gender and the standard deviation of the residuals of the constant effects model.

302 0.95 probability that the value of the intercept lies in the $[-0.05, -0.03]$ interval.

Table 2

Posterior mean, standard error, 95% credible interval and \hat{R} statistic for each parameter of the constant effect model *bmod1*.

parameter	mean	SE	lower bound	upper bound	Rhat
α	0.163	0.002	0.159	0.168	1.000
β	-0.042	0.005	-0.051	-0.033	1.000
σ_e	0.098	0.002	0.095	0.102	1.000

303 The Rhat value corresponds to the potential scale reduction factor \hat{R} (Gelman & Rubin,

1992), that provides information about the convergence of the algorithm. This index can be conceived as equivalent to the F-ratio in ANOVA. It compares the between-chains variability (i.e., the extent to which different chains differ one from each other) to the within-chain variability (i.e., how widely a chain explores the parameter space), and, as such, gives an index of the convergence of the chains. An overly large between-chains variance (as compared to the within-chain variability) would be a sign that chain-specific characteristics, like the starting value of the algorithm, have a strong influence on the final result. Ideally, the value of `Rhat` should be close to 1, and should not exceed 1.1. Otherwise, one might consider running more iterations or defining stronger priors (Bürkner, 2017b; Gelman et al., 2013).

2.3 Varying intercept model

The first model can be improved by taking into account the dependency between vowel formant measures for each participant. This is handled in MLMs by specifying unique intercepts $\alpha_{subject[i]}$ and by assigning them a common prior distribution. This strategy corresponds to the following by-subject varying-intercept model, `bmod2`:

$$\begin{aligned} \text{distance}_i &\sim \text{Normal}(\mu_i, \sigma_e) \\ \mu_i &= \alpha + \alpha_{subject[i]} + \beta \times \text{gender}_i \\ \alpha_{subject} &\sim \text{Normal}(0, \sigma_{subject}) \\ \alpha &\sim \text{Normal}(0, 10) \\ \beta &\sim \text{Normal}(0, 10) \\ \sigma_{subject} &\sim \text{HalfCauchy}(10) \\ \sigma_e &\sim \text{HalfCauchy}(10) \end{aligned}$$

This model can be fitted with `brms` with the following command (where we specify the HalfCauchy prior on $\sigma_{subject}$ by applying it on parameters of class `sd`):

```
prior2 <- c(
  prior(normal(0, 10), class = Intercept),
  prior(normal(0, 10), class = b, coef = gender),
  prior(cauchy(0, 10), class = sd),
  prior(cauchy(0, 10), class = sigma)
)

bmod2 <- brm(
  distance ~ gender + (1|subj),
  data = indo, family = gaussian(),
  prior = prior2,
  warmup = 2000, iter = 10000
)
```

320 As described in the first part of the present paper, we now have two sources of
321 variation in the model: the standard deviation of the residuals σ_e and the standard deviation
322 of the by-subject varying intercepts $\sigma_{subject}$. The latter represents the standard deviation of
323 the population of varying intercepts, and is also learned from the data. It means that the
324 estimation of each unique intercept will inform the estimation of the population of intercepts,
325 which, in return, will inform the estimation of the other intercepts. We call this sharing of
326 information between groups the *partial pooling* strategy, in comparison with the *no pooling*
327 strategy, where each intercept is estimated independently, and with the *complete pooling*
328 strategy, in which all intercepts are given the same value (Gelman et al., 2013; Gelman &
329 Hill, 2007; McElreath, 2016). This is one of the most essential features of MLMs, and what
330 leads to better estimations than single-level regression models for repeated measurements or
331 unbalanced sample sizes. This pooling of information is made apparent through the
332 phenomenon of *shrinkage*, which is illustrated in Figure 4, and later on, in Figure 6.

333 Figure 4 shows the posterior distribution as estimated by this second model for each

334 participant, in relation to the raw mean of its category (i.e., females or males), represented
335 by the vertical dashed lines. We can see for instance that participants M02 and F09 have
336 smaller average distance than the means of their groups, while participants M03 and F08
337 have larger ones. The arrows represent the amount of *shrinkage*, that is, the deviation
338 between the mean in the raw data (represented by a cross underneath each density) and the
339 estimated mean of the posterior distribution (represented by the peak of the arrow). As
340 shown in Figure 4, this *shrinkage* is always directed toward the mean of the considered group
341 (i.e., females or males) and the amount of *shrinkage* is determined by the deviation of the
342 individual mean from its group mean. This mechanism acts like a safeguard against
343 overfitting, preventing the model from overly trusting each individual datum.

344 The marginal posterior distribution of each parameter obtained with `bmod2` is
345 summarised in Table 3, where the `Rhat` values close to 1 suggest that the model has
346 converged. We see that the estimates of α and β are similar to the estimates of the first
347 model, except that the SE is now slightly larger. This result might seem surprising at first
348 sight, as we expected to improve the first model by adding a by-subject varying intercept. In
349 fact, it reveals an underestimation of the SE when using the first model. Indeed, the first
350 model assumes independence of observations, which is violated in our case. This highlights
351 the general need for careful consideration of the model’s assumptions when interpreting its
352 estimations. The first model seemingly gives highly certain estimates, but these estimations
353 are only valid in the “independence of observations” world (see also the distinction between
354 the *small world* and the *large world* in McElreath, 2016). Moreover, estimating an intercept
355 by subject (as in the second model) increases the precision of estimation, but it also makes
356 the average estimation less certain, thus resulting in a larger SE.

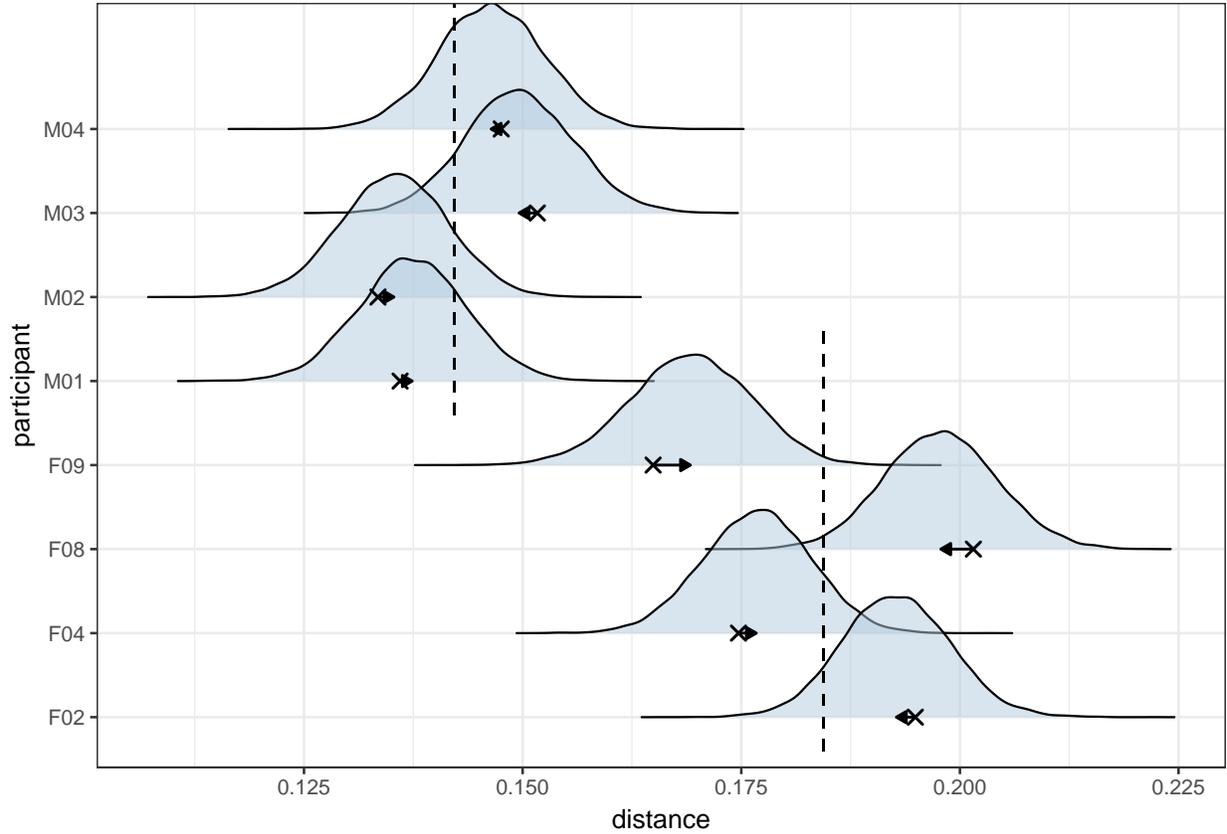


Figure 4. Posterior distributions by subject, as estimated by the `bmod2` model. The vertical dashed lines represent the means of the formant distances for the female and male groups. Crosses represent the mean of the raw data, for each participant. Arrows represent the amount of shrinkage, between the raw mean and the estimation of the model (the mean of the posterior distribution).

Table 3

Posterior mean, standard error, 95% credible interval and \hat{R} statistic for each parameter of model `bmod2` with a varying intercept by subject.

parameter	mean	SE	lower bound	upper bound	Rhat
α	0.163	0.006	0.150	0.176	1.001
β	-0.042	0.013	-0.068	-0.017	1.001
$\sigma_{subject}$	0.016	0.008	0.006	0.035	1.000
σ_e	0.098	0.002	0.095	0.101	1.000

357 This model (`bmod2`), however, is still not adequate to describe the data, as the
 358 dependency between repetitions of each vowel is not taken into account. In `bmod3`, we added
 359 a by-vowel varying intercept, thus also allowing each vowel to have a different general level of
 360 variability.

$$\begin{aligned} \text{distance}_i &\sim \text{Normal}(\mu_i, \sigma_e) \\ \mu_i &= \alpha + \alpha_{\text{subject}[i]} + \alpha_{\text{vowel}[i]} + \beta \times \text{gender}_i \\ \alpha_{\text{subj}} &\sim \text{Normal}(0, \sigma_{\text{subject}}) \\ \alpha_{\text{vowel}} &\sim \text{Normal}(0, \sigma_{\text{vowel}}) \\ \alpha &\sim \text{Normal}(0, 10) \\ \beta &\sim \text{Normal}(0, 10) \\ \sigma_e &\sim \text{HalfCauchy}(10) \\ \sigma_{\text{subject}} &\sim \text{HalfCauchy}(10) \\ \sigma_{\text{vowel}} &\sim \text{HalfCauchy}(10) \end{aligned}$$

361 This model can be fitted with `brms` with the following command:

```
prior3 <- c(
  prior(normal(0, 10), class = Intercept),
  prior(normal(0, 10), class = b, coef = gender),
  prior(cauchy(0, 10), class = sd),
  prior(cauchy(0, 10), class = sigma)
)

bmod3 <- brm(
  distance ~ gender + (1|subj) + (1|vowel),
  data = indo, family = gaussian(),
```

```
prior = prior3,
warmup = 2000, iter = 10000
)
```

362 where the same Half-Cauchy is specified for the two varying intercepts, by applying it
 363 directly to the `sd` class.

Table 4

Posterior mean, standard error, 95% credible interval and \hat{R} statistic for each parameter of model `bmod3` with a varying intercept by subject and by vowel.

parameter	mean	SE	lower bound	upper bound	Rhat
α	0.164	0.040	0.086	0.244	1.000
β	-0.042	0.013	-0.069	-0.014	1.000
$\sigma_{subject}$	0.017	0.008	0.007	0.036	1.000
σ_{vowel}	0.075	0.048	0.031	0.196	1.000
σ_e	0.088	0.002	0.085	0.091	1.000

364 The marginal posterior distribution of each parameter is summarised in Table 4. We
 365 can compute the intra-class correlation (ICC, see section 1.2) to estimate the relative
 366 variability associated with each varying effect: $ICC_{subject}$ is equal to 0.03 and ICC_{vowel} is
 367 equal to 0.42. The rather high ICC for vowels suggests that observations are highly correlated
 368 within each vowel, thus stressing the relevance of allocating a unique intercept by vowel⁶.

⁶ But please note that we do not mean to suggest that the varying intercept for subjects should be removed because its ICC is low.

369 2.4 Including a correlation between varying intercept and varying slope

370 One can legitimately question the assumption that the differences between male and
371 female productions are identical for each vowel. To explore this issue, we thus added a
372 varying slope for the effect of gender, allowing it to vary by vowel. Moreover, we can exploit
373 the correlation between the baseline level of variability by vowel, and the amplitude of the
374 difference between males and females in pronouncing them. For instance, we can observe
375 that the pronunciation of /a/ is more variable in general. We might want to know whether
376 females tend to pronounce vowels that are situated at a specific location in the F1-F2 plane
377 with less variability than males. In other words, we might be interested in knowing whether
378 the effect of **gender** is correlated with the baseline level of variability. This is equivalent to
379 investigating the *dependency*, or the correlation between the varying intercepts and the
380 varying slopes. We thus estimated this correlation by modelling α_{vowel} and β_{vowel} as issued
381 from the same multivariate normal distribution (a multivariate normal distribution is a
382 generalisation of the usual normal distribution to more than one dimension), centered on 0
383 and with some covariance matrix **S**, as specified on the third line of the following model:

$$\text{distance}_i \sim \text{Normal}(\mu_i, \sigma_e)$$

$$\mu_i = \alpha + \alpha_{\text{subject}[i]} + \alpha_{\text{vowel}[i]} + (\beta + \beta_{\text{vowel}[i]}) \times \text{gender}_i$$

$$\begin{bmatrix} \alpha_{\text{vowel}} \\ \beta_{\text{vowel}} \end{bmatrix} \sim \text{MVNormal}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \mathbf{S}\right)$$

$$\mathbf{S} = \begin{pmatrix} \sigma_{\alpha_{\text{vowel}}}^2 & \sigma_{\alpha_{\text{vowel}}}\sigma_{\beta_{\text{vowel}}}\rho \\ \sigma_{\alpha_{\text{vowel}}}\sigma_{\beta_{\text{vowel}}}\rho & \sigma_{\beta_{\text{vowel}}}^2 \end{pmatrix}$$

$$\alpha_{\text{subject}} \sim \text{Normal}(0, \sigma_{\text{subject}})$$

$$\alpha \sim \text{Normal}(0, 10)$$

$$\beta \sim \text{Normal}(0, 10)$$

$$\sigma_e \sim \text{HalfCauchy}(10)$$

$$\sigma_{\alpha_{\text{vowel}}} \sim \text{HalfCauchy}(10)$$

$$\sigma_{\beta_{\text{vowel}}} \sim \text{HalfCauchy}(10)$$

$$\sigma_{\text{subject}} \sim \text{HalfCauchy}(10)$$

$$\mathbf{R} \sim \text{LKJcorr}(2)$$

384 where \mathbf{R} is the correlation matrix $\mathbf{R} = \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}$ and ρ is the correlation between
 385 intercepts and slopes, used in the computation of \mathbf{S} . This matrix is given the
 386 LKJ-Correlation prior (Lewandowski, Kurowicka, & Joe, 2009) with a parameter ζ (zeta)
 387 that controls the strength of the correlation⁷. When $\zeta = 1$, the prior distribution on the
 388 correlation is uniform between -1 and 1 . When $\zeta > 1$, the prior distribution is peaked
 389 around a zero correlation, while lower values of ζ ($0 < \zeta < 1$) allocate more weight to
 390 extreme values (i.e., close to -1 and 1) of ρ (see Figure 5).

⁷ The LKJ prior is the default prior for correlation matrices in `brms`.

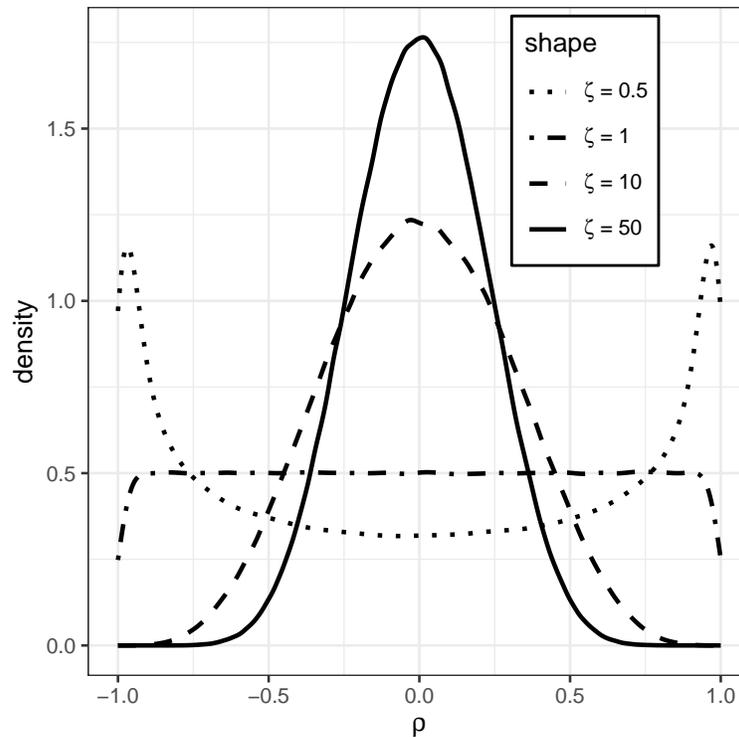


Figure 5. Visualisation of the LKJ prior for different values of the shape parameter ζ .

```
prior4 <- c(
  prior(normal(0, 10), class = Intercept),
  prior(normal(0, 10), class = b, coef = gender),
  prior(cauchy(0, 10), class = sd),
  prior(cauchy(0, 10), class = sigma),
  prior(lkj(2), class = cor)
)

bmod4 <- brm(
  distance ~ gender + (1|subj) + (1 + gender|vowel),
  data = indo, family = gaussian(),
  prior = prior4,
  warmup = 2000, iter = 10000
)
```

)

391 Estimates of this model are summarised in Table 5. This summary reveals a negative
 392 correlation between the intercepts and slopes for vowels, meaning that vowels with a large
 393 “baseline level of variability” (i.e., with a large average **distance** value) tend to be
 394 pronounced with more variability by females than by males. However, we notice that this
 395 model’s estimation of β is even more uncertain than that of the previous models, as shown
 396 by the associated standard error and the width of the credible interval.

Table 5

Posterior mean, standard error, 95% credible interval and \hat{R} statistic for each parameter of model `bmod4` with a varying intercept and varying slope by vowel.

parameter	mean	SE	lower bound	upper bound	Rhat
α	0.164	0.036	0.096	0.237	1.001
β	-0.042	0.030	-0.099	0.016	1.000
$\sigma_{subject}$	0.016	0.008	0.007	0.036	1.000
$\sigma_{\alpha_{vowel}}$	0.067	0.043	0.029	0.171	1.000
$\sigma_{\beta_{vowel}}$	0.052	0.031	0.022	0.132	1.000
ρ	-0.497	0.356	-0.951	0.371	1.001
σ_e	0.086	0.001	0.084	0.089	1.000

397 Figure 6 illustrates the negative correlation between the by-vowel intercepts and the
 398 by-vowel slopes, meaning that vowels that tend to have higher “baseline variability” (i.e.,
 399 /e/, /o/, /a/), tend to show a stronger effect of **gender**. This figure also illustrates the
 400 amount of shrinkage, here in the parameter space. We can see that the *partial pooling*
 401 estimate is shrunk somewhere between the *no pooling* estimate and the *complete pooling*
 402 estimate (i.e., the grand mean). This illustrates again the mechanism by which MLMs

403 balance the risk of overfitting and underfitting (McElreath, 2016).

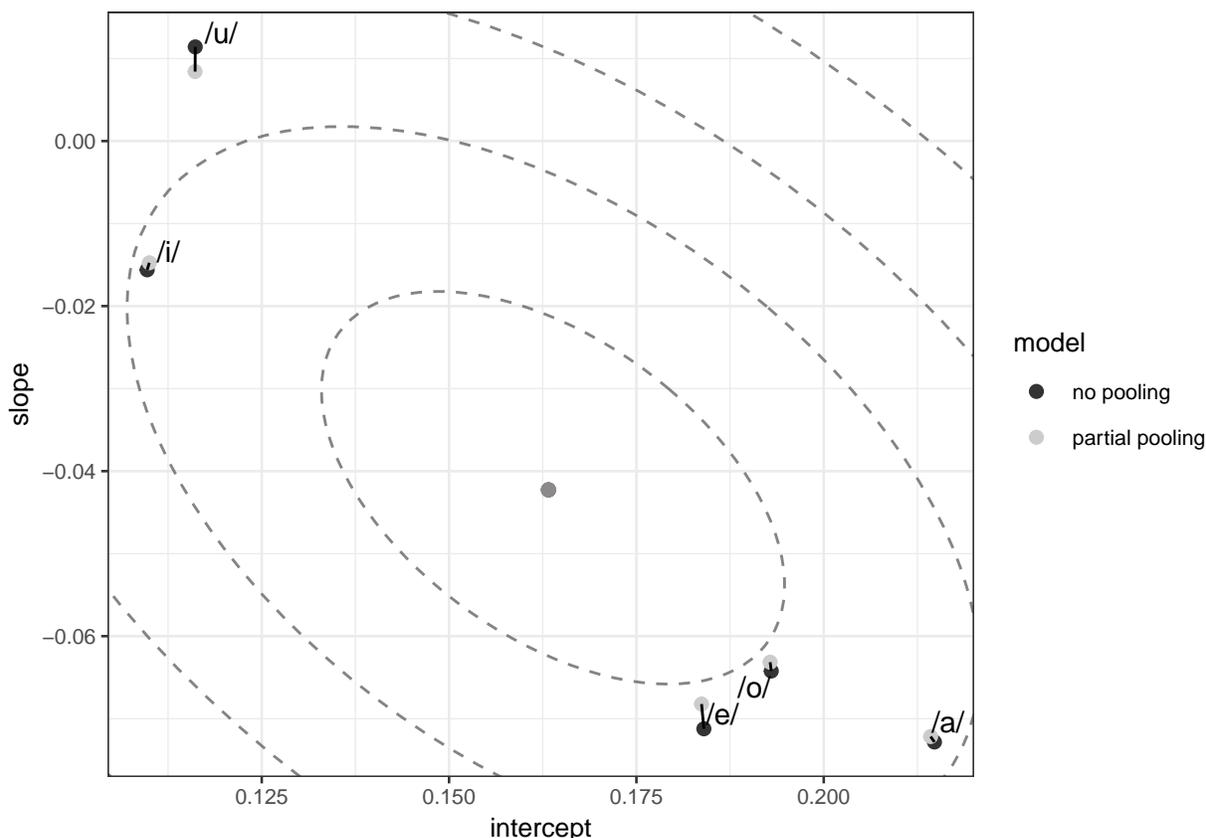


Figure 6. Shrinkage of estimates in the parameter space, due to the pooling of information between clusters (based on the `bmod4` model). The ellipses represent the contours of the bivariate distribution, at different degrees of confidence 0.1, 0.3, 0.5 and 0.7.

404 2.5 Varying intercept and varying slope model, interaction between subject 405 and vowel

406 So far, we modelled varying effects of subjects and vowels. In this study, these varying
407 factors were crossed, meaning that every subject had to pronounce every vowel. Let us now
408 imagine a situation in which Subject 4 systematically mispronounced the /i/ vowel. This
409 would be a source of systematic variation over replicates which is not considered in the
410 model (`bmod4`), because this model can only adjust parameters for either vowel or
411 participant, but not for a specific vowel for a specific participant.

412 In building the next model, we added a varying intercept for the interaction between
 413 subject and vowel, that is, we created an index variable that allocates a unique value at each
 414 crossing of the two variables (e.g., Subject1-vowel/a/, Subject1-vowel/i/, etc.), resulting in 8
 415 $\times 5 = 40$ intercepts to be estimated (for a review of multilevel modeling in various
 416 experimental designs, see Judd, Westfall, & Kenny, 2017). This varying intercept for the
 417 interaction between subject and vowel represents the systematic variation associated with a
 418 specific subject pronouncing a specific vowel. This model can be written as follows, for any
 419 observation i :

$$\text{distance}_i \sim \text{Normal}(\mu_i, \sigma_e)$$

$$\mu_i = \alpha + \alpha_{\text{subject}[i]} + \alpha_{\text{vowel}[i]} + \alpha_{\text{subject:vowel}[i]} + (\beta + \beta_{\text{vowel}[i]}) \times \text{gender}_i$$

$$\begin{bmatrix} \alpha_{\text{vowel}} \\ \beta_{\text{vowel}} \end{bmatrix} \sim \text{MVNormal} \left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \mathbf{S} \right)$$

$$\mathbf{S} = \begin{pmatrix} \sigma_{\alpha_{\text{vowel}}}^2 & \sigma_{\alpha_{\text{vowel}}} \sigma_{\beta_{\text{vowel}}} \rho \\ \sigma_{\alpha_{\text{vowel}}} \sigma_{\beta_{\text{vowel}}} \rho & \sigma_{\beta_{\text{vowel}}}^2 \end{pmatrix}$$

$$\alpha_{\text{subject}} \sim \text{Normal}(0, \sigma_{\text{subject}})$$

$$\alpha_{\text{subject:vowel}} \sim \text{Normal}(0, \sigma_{\text{subject:vowel}})$$

$$\alpha \sim \text{Normal}(0, 10)$$

$$\beta \sim \text{Normal}(0, 10)$$

$$\sigma_e \sim \text{HalfCauchy}(10)$$

$$\sigma_{\text{subject}} \sim \text{HalfCauchy}(10)$$

$$\sigma_{\text{subject:vowel}} \sim \text{HalfCauchy}(10)$$

$$\sigma_{\alpha_{\text{vowel}}} \sim \text{HalfCauchy}(10)$$

$$\sigma_{\beta_{\text{vowel}}} \sim \text{HalfCauchy}(10)$$

$$\mathbf{R} \sim \text{LKJcorr}(2)$$

421 This model can be fitted with the following command:

```
prior5 <- c(  
  prior(normal(0, 10), class = Intercept),  
  prior(normal(0, 10), class = b, coef = gender),  
  prior(cauchy(0, 10), class = sd),  
  prior(cauchy(0, 10), class = sigma),  
  prior(lkj(2), class = cor)  
)  
  
bmod5 <- brm(  
  distance ~ gender + (1|subj) + (1 + gender|vowel) + (1|subj:vowel),  
  data = indo, family = gaussian(),  
  prior = prior5,  
  warmup = 2000, iter = 10000  
)
```

422 Estimates of this model are summarised in Table 6. From this table, we first notice
423 that the more varying effects we add, the more the model is uncertain about the estimation
424 of α and β , which can be explained in the same way as in section 2.2. Second, we see the
425 opposite pattern for σ_e , the residuals standard deviation, which has decreased by a
426 considerable amount compared to the first model, indicating a better fit.

Table 6

Posterior mean, standard error, 95% credible interval and \hat{R} statistic for each parameter of model bmod5 with a varying intercept and a varying slope by vowel and a varying intercept for the interaction between subject and vowel.

parameter	mean	SE	lower bound	upper bound	Rhat
α	0.163	0.038	0.087	0.236	1.000
β	-0.042	0.030	-0.100	0.018	1.000
$\sigma_{subject}$	0.012	0.009	0.001	0.033	1.000
$\sigma_{subject:vowel}$	0.024	0.005	0.016	0.034	1.000
$\sigma_{\alpha_{vowel}}$	0.070	0.046	0.029	0.183	1.000
$\sigma_{\beta_{vowel}}$	0.050	0.038	0.013	0.144	1.000
ρ	-0.433	0.380	-0.946	0.454	1.000
σ_e	0.085	0.001	0.082	0.088	1.000

427

3 Model comparison

428 Once we have built a set of models, we need to know which model is the more accurate
 429 and should be used to draw conclusions. It might be a little tricky to select the model that
 430 has the better absolute fit on the actual data (using for instance R^2), as this model will not
 431 necessarily perform as well on new data. Instead, we might want to choose the model that
 432 has the best predictive abilities, that is, the model that performs the best when it comes to
 433 predicting data that have not yet been observed. We call this ability the out-of-sample
 434 predictive performance of the model (McElreath, 2016). When additional data is not
 435 available, cross-validation techniques can be used to obtain an approximation of the model's
 436 predictive abilities, among which the Bayesian leave-one-out-cross-validation (LOO-CV,
 437 Vehtari, Gelman, & Gabry, 2017). Another useful tool, and asymptotically equivalent to the

438 LOO-CV, is the Watanabe Akaike Information Criterion (WAIC, Watanabe, 2010), which
439 can be conceived as a generalisation of the Akaike Information Criterion (AIC, Akaike,
440 1974)⁸.

441 Both WAIC and LOO-CV indexes are easily computed in `brms` with the `WAIC` and the
442 `LOO` functions, where n models can be compared with the following call: `LOO(model1,`
443 `model2, ..., modeln)`. These functions also provide an estimate of the uncertainty
444 associated with these indexes (in the form of a SE), as well as a difference score ΔLOOIC ,
445 which is computed by taking the difference between each pair of information criteria. The
446 `WAIC` and the `LOO` functions also provide a SE for these delta values (ΔSE). A comparison of
447 the five models we fitted can be found in Table 7.

⁸ More details on model comparison using cross-validation techniques can be found in Nicenboim and Vasishth (2016). See also Gelman, Hwang, and Vehtari (2014) for a complete comparison of information criteria.

Table 7

Model comparison with LOOIC.

Model	LOOIC	SE	Δ LOOIC	Δ SE	right side of the formula
bmod5	-3600.29	68.26	0.00	0.00	gender + (1 subj) + (1 + gender vowel) + (1 subj:vowel)
bmod4	-3544.66	66.92	55.63	14.94	gender + (1 subj) + (1 + gender vowel)
bmod3	-3484.21	67.15	116.08	20.22	gender + (1 subj) + (1 vowel)
bmod2	-3119.41	65.32	480.88	39.50	gender + (1 subj)
bmod1	-3103.43	66.72	496.86	40.52	gender

449 We see from Table 7 that `bmod5` (i.e., the last model) is performing much better than
 450 the other models, as it has the lower LOOIC. We then based our conclusions (see last
 451 section) on the estimations of this model. We also notice that each addition to the initial
 452 model brought improvement in terms of predictive accuracy, as the set of models is ordered
 453 from the first to the last model. This should not be taken as a general rule though, as
 454 successive additions made to an original model could also lead to *overfitting*, corresponding
 455 to a situation in which the model is over-specified in regards to the data, which makes the
 456 model good to explain the data at hand, but very bad to predict non-observed data. In such
 457 cases, information criteria and indexes that rely exclusively on goodness-of-fit (such as R^2)
 458 would point to different conclusions.

459 4 Comparison of `brms` and `lme4` estimations

460 Figure 7 illustrates the comparison of `brms` (Bayesian approach) and `lme4` (frequentist
 461 approach) estimates for the last model (`bmod5`), fitted in `lme4` with the following command.

```
lmer_model <- lmer(
  distance ~ gender + (1|subj) + (1 + gender|vowel) + (1|subj:vowel),
  REML = FALSE, data = indo
)
```

462 Densities represent the posterior distribution as estimated by `brms` along with 95%
 463 credible intervals, while the crosses underneath represent the *maximum likelihood estimate*
 464 (MLE) from `lme4` along with 95% confidence intervals, obtained with parametric
 465 bootstrapping.

466 We can see that the estimations of `brms` and `lme4` are for the most part very similar.
 467 The differences we observe for $\sigma_{\alpha_{vowel}}$ and $\sigma_{\beta_{vowel}}$ might be explained by the skewness of the
 468 posterior distribution. Indeed, in these cases (i.e., when the distribution is not symmetric),
 469 the mode of the distribution would better coincide with the `lme4` estimate. This figure also

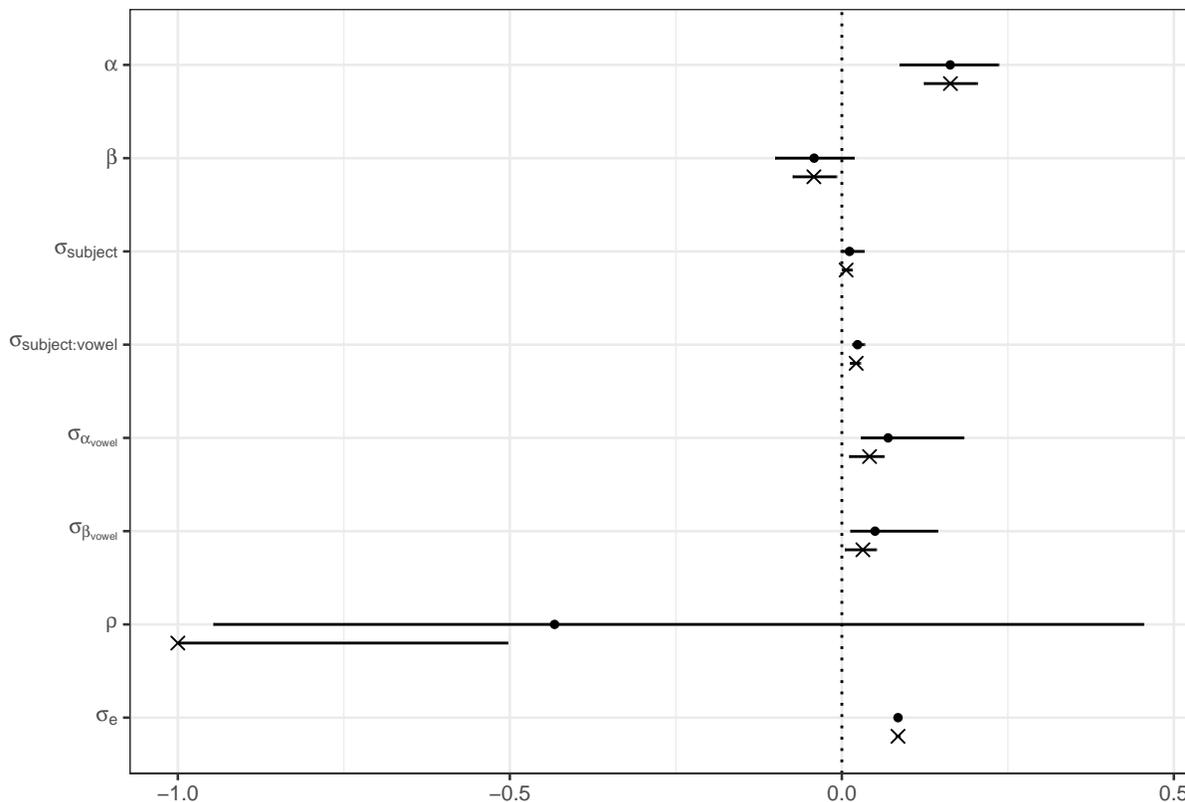


Figure 7. Comparison of estimations from `brms` and `lme4`. Dots represent means of posterior distribution along with 95% CrIs, as estimated by the `bmod5` model. Crosses represent estimations of `lme4` along with bootstrapped 95% CIs.

470 illustrates a limitation of frequentist MLMs that we discussed in the first part of the current
 471 paper. If we look closely at the estimates of `lme4`, we can notice that the MLE for the
 472 correlation ρ is at its boundary, as $\rho = -1$. This might be interpreted in (at least) two ways.
 473 The first interpretation is what Eager and Roy (2017) call the *parsimonious convergence*
 474 *hypothesis* (PCH) and consists in saying that this aberrant estimation is caused by the
 475 over-specification of the random structure (e.g., Bates et al., 2015a). In other words, this
 476 would correspond to a model that contains too many varying effects to be “supported” by a
 477 certain dataset (but this does not mean that with more data, this model would not be a
 478 correct model). However, the PCH has been questioned by Eager and Roy (2017), who have
 479 shown that under conditions of unbalanced datasets, non-linear models fitted with `lme4`

480 provided more prediction errors than Bayesian models fitted with **Stan**. The second
481 interpretation considers failures of convergence as a problem of frequentist MLMs *per se*,
482 which is resolved in the Bayesian framework by using weakly informative priors (i.e., the
483 LKJ prior) for the correlation between varying effects (e.g., Eager & Roy, 2017; Nicenboim &
484 Vasishth, 2016), and by using the full posterior for inference.

485 One feature of the Bayesian MLM in this kind of situation is to provide an estimate of
486 the correlation that incorporates the uncertainty caused by the weak amount of data (i.e., by
487 widening the posterior distribution). Thus, the **brms** estimate of the correlation coefficient
488 has its posterior mean at $\rho = -0.433$, but this estimate comes with a huge uncertainty, as
489 expressed by the width of the credible interval (95% CrI = $[-0.946, 0.454]$).

490 5 Inference and conclusions

491 Regarding our initial question, which was to know whether there is a gender effect on
492 vowel production variability in standard Indonesian, we can base our conclusions on several
493 parameters and indices. However, the discrepancies between the different models we fitted
494 deserve some discussion first. As already pointed out previously, if we had based our
495 conclusions on the results of the first model (i.e., the model with constant effects only), we
496 would have confidently concluded on a positive effect of gender. However, when we included
497 the appropriate error terms in the model to account for repeated measurements by subject
498 and by vowel, as well as for the by-vowel specific effect of gender, the large variability of this
499 effect among vowels lead the model to adjust its estimation of β , resulting in more
500 uncertainty about it. The last model then estimated a value of $\beta = -0.04$ with quite a large
501 uncertainty (95% CrI = $[-0.10, 0.02]$), and considering 0 as well as some positive values as
502 credible. This result alone makes it difficult to reach any definitive conclusion concerning the
503 presence or absence of a gender effect on the variability of vowels pronunciation in
504 Indonesian, and should be considered (at best) as suggestive.

505 Nevertheless, it is useful to recall that in the Bayesian framework, the results of our

506 analysis is a (posterior) probability distribution which can be, as such, summarised in
 507 multiple ways. This distribution is plotted in Figure 8, which also shows the mean and the
 508 95% CrI, as well as the proportion of the distribution below and above a particular value⁹.
 509 This figure reveals that 94.1% of the distribution is below 0, which can be interpreted as
 510 suggesting that there is a 0.94 probability that males have a lower mean formant distance
 511 than females (recall that female was coded as -0.5 and male as 0.5), given the data at hand,
 512 and the model.

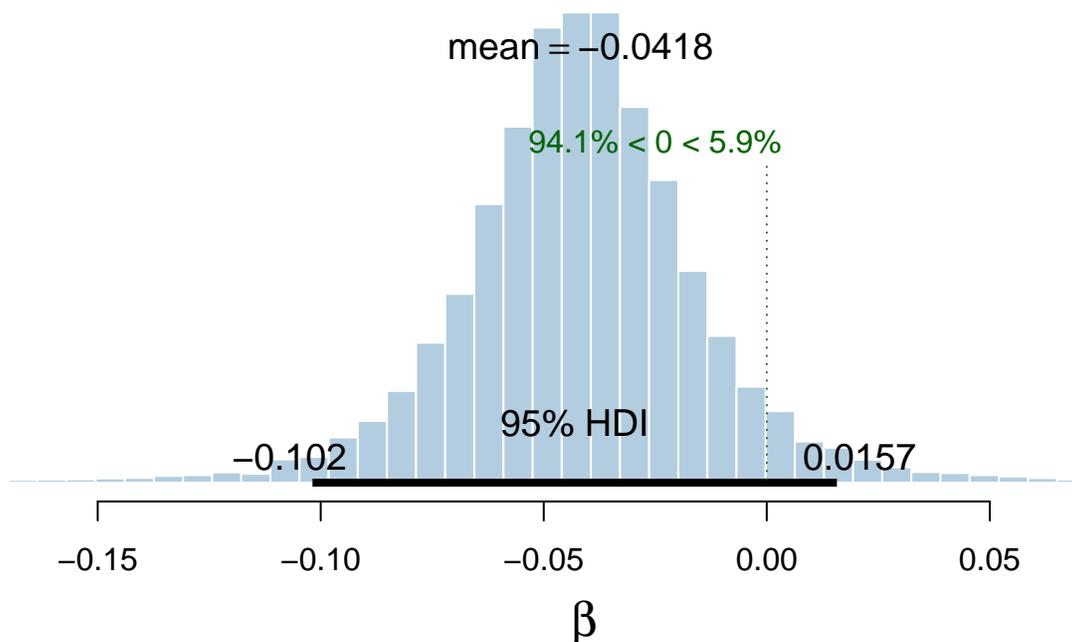


Figure 8. Histogram of posterior samples of the slope for `gender`, as estimated by the last model.

513 This quantity can be easily computed from the posterior samples:

```
post <- posterior_samples(bmod5) # extracting posterior samples
mean(post$b_gender < 0) # computing p(beta<0)
```

514 ## [1] 0.940625

⁹ We compare the distribution with 0 here, but it should be noted that this comparison could be made with whatever value.

515 Of course, this estimate can (and should) be refined using more data from several
 516 experiments, with more speakers. In this line, it should be pointed out that `brms` can easily
 517 be used to extend the multilevel strategy to meta-analyses (e.g., Bürkner et al., 2017;
 518 Williams & Bürkner, 2017). Its flexibility makes it possible to fit multilevel hierarchical
 519 Bayesian models at two, three, or more levels, enabling researchers to model the
 520 heterogeneity between studies as well as dependencies between experiments of the same
 521 study, or between studies carried out by the same research team. Such a modelling strategy
 522 is usually equivalent to the ordinary frequentist random-effect meta-analysis models, while
 523 offering all the benefits inherent to the Bayesian approach.

524 Another useful source of information comes from the examination of effects sizes. One
 525 of the most used criteria is Cohen’s d standardized effect size, that expresses the difference
 526 between two groups in terms of their pooled standard deviation:

$$\text{Cohen's } d = \frac{\mu_1 - \mu_2}{\sigma_{pooled}} = \frac{\mu_1 - \mu_2}{\sqrt{\frac{\sigma_1^2 + \sigma_2^2}{2}}}$$

527 However, as the total variance is partitioned into multiple sources of variation in
 528 MLMs, there is no unique way of computing a standardised effect size. While several
 529 approaches have been suggested (e.g., dividing the mean difference by the standard deviation
 530 of the residuals), the more consensual one involves taking into account all of the variance
 531 sources of the model (Hedges, 2007). One such index is called the δ_t (where the t stands for
 532 “total”), and is given by the estimated difference between group means, divided by the
 533 square root of the sum of all variance components:

$$\delta_t = \frac{\beta}{\sqrt{\sigma_{subject}^2 + \sigma_{subject:vowel}^2 + \sigma_{\alpha_{vowel}}^2 + \sigma_{\beta_{vowel}}^2 + \sigma^2}}$$

534 As this effect size is dependent on the parameters estimated by the model, one can
 535 derive a probability distribution for this index as well. This is easily done in R, computing it

536 from the posterior samples:

```
delta_t <-
  # extracting posterior samples from bmod5
  posterior_samples(bmod5, pars = c("^b_", "sd_", "sigma") ) %>%
  # taking the square of each variance component
  mutate_at(.vars = 3:7, .funs = funs(.^2) ) %>%
  # dividing the slope estimate by the square root of the sum of
  # all variance components
  mutate(delta = b_gender / sqrt(rowSums(.[3:7]) ) )
```

537 This distribution is plotted in Figure 9, and reveals the large uncertainty associated
538 with the estimation of δ_t .

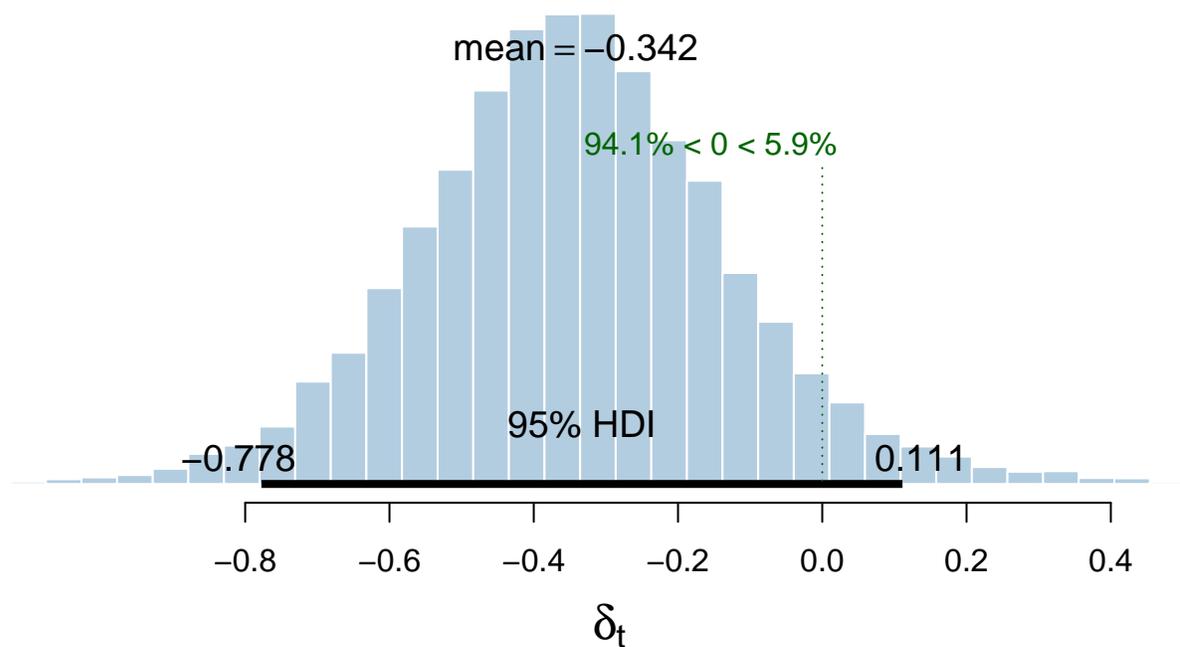


Figure 9. Posterior distribution of δ_t .

539 In the same fashion, undirected effect sizes (e.g., R^2) can be computed directly from
540 the posterior samples, or included in the model specification as a parameter of the model, in
541 a way that at each iteration of the MCMC, a value of the effect size is sampled, resulting in

542 an estimation of its full posterior distribution (see for instance Gelman & Pardoe, 2006 for
 543 measures of explained variance in MLMs and @Marsman2017 for calculations in ANOVA
 544 designs). A Bayesian version of the R^2 is also available in `brms` using the `bayes_R2` method,
 545 for which the calculations are based on Gelman, Goodrich, Gabry, and Ali (2017).

```
bayes_R2(bmod5)
```

```
546 ##      Estimate  Est.Error  2.5%ile  97.5%ile
547 ## R2  0.295614  0.01589917  0.2635006  0.3262617
```

548 In brief, we found a weak effect of gender on vowel production variability in Indonesian
 549 ($\beta = -0.04$, 95% CrI = $[-0.10, 0.02]$, $\delta_t = -0.34$, 95% CrI = $[-0.78, 0.11]$), this effect being
 550 associated with a large uncertainty (as expressed by the width of the credible interval). This
 551 result seems to show that females tend to pronounce vowels with more variability than males,
 552 while the variation observed across vowels (as suggested by $\sigma_{\beta_{vowel}}$) suggests that there might
 553 exist substantial inter-vowel variability, that should be subsequently properly studied. A
 554 follow-up analysis specifically designed to test the effect of gender on each vowel should help
 555 better describe inter-vowel variability (we give an example of such an analysis in the
 556 supplementary materials).

557 To sum up, we hope that this introductory tutorial has helped the reader to understand
 558 the foundational ideas of Bayesian MLMs, and to appreciate how straightforward the
 559 interpretation of the results is. Moreover, we hope to have demonstrated that although
 560 Bayesian data analysis may still sometimes (wrongfully) sound difficult to grasp and to use,
 561 the development of recent tools like `brms` helps to build and fit Bayesian MLMs in an
 562 intuitive way. We believe that this shift in practice will allow more reliable statistical
 563 inferences to be drawn from empirical research.

6 Supplementary materials

Supplementary materials, reproducible code and figures are available at: osf.io/dpzcb. A lot of useful packages have been used for the writing of this paper, among which the `papaja` and `knitr` packages for writing and formatting (Aust & Barth, 2017; Xie, 2015), the `ggplot2`, `viridis`, `ellipse`, `BEST`, and `ggridges` packages for plotting (Garnier, 2017; Kruschke & Meredith, 2017; Murdoch & Chow, 2013; Wickham, 2009; Wilke, 2017), as well as the `tidyverse` and `broom` packages for code writing and formatting (Robinson, 2017; Wickham, 2017).

Acknowledgements

We thank Brice Beffara for helpful comments on a previous version of this manuscript, as well as Shravan Vasishth and one anonymous reviewer for insightful suggestions during the review process.

References

- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6), 716–723. <https://doi.org/10.1109/tac.1974.1100705>
- Aust, F., & Barth, M. (2017). *papaja: Create APA manuscripts with R Markdown*. Retrieved from <https://github.com/crsh/papaja>
- Bakan, D. (1966). The test of significance in psychological research. *Psychological Bulletin*, 66(6), 423–437. <https://doi.org/10.1037/h0020412>
- Barr, D. J., Levy, R., Scheepers, C., & Tily, H. J. (2013). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*, 68(3), 255–278. <https://doi.org/10.1016/j.jml.2012.11.001>
- Bates, D., Kliegl, R., Vasishth, S., & Baayen, R. H. (2015a). *Parsimonious mixed models*. Retrieved from <https://arxiv.org/pdf/1506.04967.pdf>

- 588 Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015b). Fitting linear mixed-effects models
589 using lme4. *Journal of Statistical Software*, *67*(1), 1–48.
590 <https://doi.org/10.18637/jss.v067.i01>
- 591 Bürkner, P.-C. (2017a). *Advanced bayesian multilevel modeling with the R package brms*.
592 Retrieved from <https://arxiv.org/pdf/1705.11123>
- 593 Bürkner, P.-C. (2017b). brms: An R package for bayesian multilevel models using Stan.
594 *Journal of Statistical Software*, *80*(1), 1–28. <https://doi.org/10.18637/jss.v080.i01>
- 595 Bürkner, P.-C., Williams, D. R., Simmons, T. C., & Woolley, J. D. (2017). Intranasal
596 oxytocin may improve high-level social cognition in schizophrenia, but not social
597 cognition or neurocognition in general: A multilevel Bayesian meta-analysis.
598 *Schizophrenia Bulletin*, *43*(6), 1291–1303. <https://doi.org/10.1093/schbul/sbx053>
- 599 Cumming, G. (2012). *Understanding the new statistics: Effect sizes, confidence intervals,*
600 *and meta-analysis*. New York: Routledge.
- 601 Cumming, G. (2014). The new statistics: Why and how. *Psychological Science*, *25*(1), 7–29.
602 <https://doi.org/10.1177/0956797613504966>
- 603 Dienes, Z. (2011). Bayesian versus orthodox statistics: Which side are you on? *Perspectives*
604 *on Psychological Science*, *6*(3), 274–290. <https://doi.org/10.1177/1745691611406920>
- 605 Eager, C., & Roy, J. (2017). *Mixed effects models are sometimes terrible*. Retrieved from
606 <https://arxiv.org/pdf/1701.04858.pdf>
- 607 Garnier, S. (2017). *viridis: Default color maps from 'matplotlib'*. Retrieved from
608 <https://CRAN.R-project.org/package=viridis>
- 609 Gelman, A. (2005). Analysis of variance — why it is more important than ever. *The Annals*
610 *of Statistics*, *33*(1), 1–53. <https://doi.org/10.1214/009053604000001048>
- 611 Gelman, A. (2006). Prior distributions for variance parameter in hierarchical models.

- 612 *Bayesian Analysis*, 1(3), 515–534. <https://doi.org/10.1214/06-ba117a>
- 613 Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2013).
614 *Bayesian data analysis, third edition*. CRC Press.
- 615 Gelman, A., Goodrich, B., Gabry, J., & Ali, I. (2017). *R-squared for Bayesian regression*
616 *models*. Retrieved from
617 https://github.com/jgabry/bayes_R2/blob/master/bayes_R2.pdf
- 618 Gelman, A., & Hill, J. (2007). *Data analysis using regression and multilevel/hierarchical*
619 *models*. Cambridge University Press, New York.
- 620 Gelman, A., Hill, J., & Yajima, M. (2012). Why we (usually) don't have to worry about
621 multiple comparisons. *Journal of Research on Educational Effectiveness*, 5, 189–211.
622 <https://doi.org/10.1080/19345747.2011.618213>
- 623 Gelman, A., Hwang, J., & Vehtari, A. (2014). Understanding predictive information criteria
624 for Bayesian models. *Statistics and Computing*, 24(6), 997–1016.
625 <https://doi.org/10.1007/s11222-013-9416-2>
- 626 Gelman, A., & Pardoe, I. (2006). Bayesian measures of explained variance and pooling in
627 multilevel (hierarchical) models. *Technometrics*, 48(2), 241–251.
628 <https://doi.org/10.1198/004017005000000517>
- 629 Gelman, A., & Rubin, D. B. (1992). Inference from Iterative Simulation Using Multiple
630 Sequences. *Statistical Science*, 7(4), 457–472. <https://doi.org/10.1214/ss/1177011136>
- 631 Gigerenzer, G., Krauss, S., & Vitouch, O. (2004). The null ritual: What you always wanted
632 to know about significance testing but were afraid to ask. *The Sage Handbook of*
633 *Methodology for the Social Sciences*, 391–408.
634 <https://doi.org/10.4135/9781412986311.n21>
- 635 Hedges, L. V. (2007). Effect sizes in cluster-randomized designs. *Journal of Educational and*

- 636 *Behavioral Statistics*, 32(4), 341–370. <https://doi.org/10.3102/1076998606298043>
- 637 Hoekstra, R., Morey, R. D., Rouder, J. N., & Wagenmakers, E.-J. (2014). Robust
638 misinterpretation of confidence intervals. *Psychonomic Bulletin & Review*, 21(5),
639 1157–1164. <https://doi.org/10.3758/s13423-013-0572-3>
- 640 Janssen, D. P. (2012). Twice random, once mixed: Applying mixed models to simultaneously
641 analyze random effects of language and participants. *Behavior Research Methods*,
642 44(1), 232–247. <https://doi.org/10.3758/s13428-011-0145-1>
- 643 Judd, C. M., Westfall, J., & Kenny, D. A. (2017). Experiments with more than one random
644 factor: Designs, analytic models, and statistical power. *Annual Review of Psychology*,
645 68, 601–625. <https://doi.org/10.1146/annurev-psych-122414-033702>
- 646 Kline, R. B. (2004). *Beyond significance testing: Reforming data analysis methods in*
647 *behavioral research*. (pp. 61–91). Washington, DC: APA Book.
648 <https://doi.org/10.1037/10693-003>
- 649 Kruschke, J. K. (2015). *Doing Bayesian data analysis, Second Edition: A tutorial with R,*
650 *JAGS, and Stan*. Burlington, MA: Academic Press / Elsevier.
- 651 Kruschke, J. K., & Liddell, T. M. (2017a). Bayesian data analysis for newcomers.
652 *Psychonomic Bulletin & Review*, 1–23. <https://doi.org/10.3758/s13423-017-1272-1>
- 653 Kruschke, J. K., & Liddell, T. M. (2017b). The Bayesian new statistics: Hypothesis testing,
654 estimation, meta-analysis, and power analysis from a Bayesian perspective.
655 *Psychonomic Bulletin & Review*, 1–29. <https://doi.org/10.3758/s13423-016-1221-4>
- 656 Kruschke, J. K., & Meredith, M. (2017). *BEST: Bayesian estimation supersedes the t-test*.
657 Retrieved from <https://CRAN.R-project.org/package=BEST>
- 658 Lambdin, C. (2012). Significance tests as sorcery: Science is empirical—significance tests are
659 not. *Theory & Psychology*, 22(1), 67–90. <https://doi.org/10.1177/0959354311429854>

- 660 Lewandowski, D., Kurowicka, D., & Joe, H. (2009). Generating random correlation matrices
661 based on vines and extended onion method. *Journal of Multivariate Analysis*, 100(9),
662 1989–2001. <https://doi.org/10.1016/j.jmva.2009.04.008>
- 663 Marsman, M., Waldorp, L., Dablander, F., & Wagenmakers, E.-J. (2017). *Bayesian*
664 *estimation of explained variance in ANOVA designs*. Retrieved from
665 http://maartenmarsman.com/wp-content/uploads/2017/04/MarsmanEtAl_R2.pdf
- 666 McCloy, D. R. (2014). Phonetic effects of morphological structure in Indonesian vowel
667 reduction. In *Proceedings of meetings on acoustics* (Vol. 12, pp. 1–14).
668 <https://doi.org/10.1121/1.4870068>
- 669 McCloy, D. R. (2016). *phonR: Tools for phoneticians and phonologists*. Retrieved from
670 <https://cran.r-project.org/web/packages/phonR/>
- 671 McElreath, R. (2016). *Statistical Rethinking*. Chapman; Hall/CRC.
- 672 Morey, R. D., Hoekstra, R., Rouder, J. N., Lee, M. D., & Wagenmakers, E.-J. (2015). The
673 fallacy of placing confidence in confidence intervals. *Psychonomic Bulletin & Review*,
674 23, 103–123. <https://doi.org/10.3758/s13423-015-0947-8>
- 675 Murdoch, D., & Chow, E. D. (2013). *ellipse: Functions for drawing ellipses and ellipse-like*
676 *confidence regions*. Retrieved from <https://CRAN.R-project.org/package=ellipse>
- 677 Nicenboim, B., & Vasishth, S. (2016). Statistical methods for linguistic research:
678 Foundational Ideas – Part II. *Language and Linguistics Compass*, 10(11), 591–613.
679 <https://doi.org/10.1111/lnc3.12207>
- 680 Polson, N. G., & Scott, J. G. (2012). On the half-cauchy prior for a global scale parameter.
681 *Bayesian Analysis*, 7(4), 887–902. <https://doi.org/10.1214/12-BA730>
- 682 R Core Team. (2017). *R: A language and environment for statistical computing*. Vienna,
683 Austria: R Foundation for Statistical Computing. Retrieved from

684 <https://www.R-project.org/>

685 Robinson, D. (2017). *broom: Convert statistical analysis objects into tidy data frames.*

686 Retrieved from <https://CRAN.R-project.org/package=broom>

687 Scott, J. G., & Berger, J. O. (2010). Bayes and empirical-Bayes multiplicity adjustment in
688 the variable-selection problem. *The Annals of Statistics*, *38*(5), 2587–2619.

689 <https://doi.org/10.1214/10-AOS792>

690 Sorensen, T., Hohenstein, S., & Vasishth, S. (2016). Bayesian linear mixed models using

691 Stan: A tutorial for psychologists, linguists, and cognitive scientists. *The Quantitative*

692 *Methods for Psychology*, *12*(3), 175–200. <https://doi.org/10.20982/tqmp.12.3.p175>

693 Stan Development Team. (2016). Stan modeling language users guide and reference manual.

694 Retrieved from <http://mc-stan.org>

695 Trafimow, D., Amrhein, V., Areshenkoff, C. N., Barrera-Causil, C. J., Beh, E. J., Bilgiç, Y.

696 K., ... Marmolejo-Ramos, F. (2018). Manipulating the alpha level cannot cure
697 significance testing. *Frontiers in Psychology*, *9*.

698 <https://doi.org/10.3389/fpsyg.2018.00699>

699 Vehtari, A., Gelman, A., & Gabry, J. (2017). Practical Bayesian model evaluation using

700 leave-one-out cross-validation and WAIC. *Statistics and Computing*, *27*(5), 1413–1432.

701 <https://doi.org/10.1007/s11222-016-9696-4>

702 Watanabe, S. (2010). Asymptotic equivalence of Bayes cross validation and widely

703 applicable information criterion in singular learning theory. *Journal of Machine*

704 *Learning Research*, *11*, 3571–3594.

705 Watt, D., & Fabricius, A. (2002). Evaluation of a technique for improving the mapping of

706 multiple speakers' vowel spaces in the F1~F2 plane. *Leeds Working Papers in*

707 *Linguistics and Phonetics*, *9*(9), 159–173.

- 708 Wickham, H. (2009). *ggplot2: Elegant graphics for data analysis*. Springer-Verlag New York.
709 Retrieved from <http://ggplot2.org>
- 710 Wickham, H. (2017). *tidyverse: Easily install and load 'tidyverse' packages*. Retrieved from
711 <https://CRAN.R-project.org/package=tidyverse>
- 712 Wilke, C. O. (2017). *ggridges: Ridgeline plots in 'ggplot2'*. Retrieved from
713 <https://CRAN.R-project.org/package=ggridges>
- 714 Williams, D. R., & Bürkner, P.-C. (2017). Psychoneuroendocrinology Effects of intranasal
715 oxytocin on symptoms of schizophrenia: A multivariate Bayesian meta-analysis.
716 *Psychoneuroendocrinology*, *75*, 141–151.
717 <https://doi.org/10.1016/j.psyneuen.2016.10.013>
- 718 Xie, Y. (2015). *Dynamic documents with R and knitr* (2nd ed.). Boca Raton, Florida:
719 Chapman; Hall/CRC. Retrieved from <https://yihui.name/knitr/>