

# Mixed-effects regression modeling

Roland Schäfer

Freie Universität Berlin

August 13, 2017

## 1 Introduction

## 2 Fundamentals

### 2.1 Introduction to random effects

*Generalized Linear Models* (GLMs), as discussed in the previous chapter, allow us to estimate the effects which various *predictors* or *regressors* (i. e., corpus linguistic variables) have on an *outcome* or *response* (i. e., another corpus linguistic variable).<sup>1</sup> Surely, the most typical application (in corpus linguistics) is the modeling of *alternations*, i. e., phenomena where the response variable of interest encodes a choice of forms or constructions, for example a case alternation (a binary or multi-valued categorical response, depending on the richness of the

---

<sup>1</sup>It should be noted that I use the term *variable* here in its neutral sense which is common in most strains of empirical and statistical research. The definition of a variable in Labovian variationist linguistics is more restricted and theoretically burdened.

language's case system), alternations of graphemic form such as contracted vs. non-contracted, ordering preferences such as the order of prenominal adjectives, or syntactic/constructional alternations such as the dative alternation.<sup>2</sup> The approach is called *generalized* in contrast to normal linear models because the response need not be numerical, and the *errors* or *residuals* do not have to be (approximately) normally distributed. First, this is achieved by allowing for different types of exponential distributions for the residuals, which requires the use of a more general estimator than least-squares, typically likelihood maximization. Second, *link functions* are introduced which relate the additive linear term that combines the predictors in a non-linear way to the response variable. Generalized Linear Mixed Models (GLMMs) are not much different. They add what are often called *random effects* and *mix* them with the normal predictors as used in GLMs. The latter one are called *fixed effects* in this terminology. Alternatively, statisticians speak of *multilevel models* or *hierarchical models* (Gelman & Hill 2006), a terminology to be explained in Section 2.3 and Section 3.

The purpose of including random effects is usually said to be the modeling of variance between groups of observations. A single observation (or *data point* or *measurement* or *unit*) is one atomic exemplar entering into the statistical analysis of the study. In corpus linguistics, single observations can be understood as a single line in a concordance, and they typically represent, for example (and with reference to the above examples), a clause or sentence in which one of the alternants of a case alternation occurs, an NP where two pre-nominal

---

<sup>2</sup>In this article, I restrict the discussion to GL(M)Ms with categorical responses, simply because the continuous responses in Linear (Mixed) Models – or LM(M)s – are not found very often in corpus linguistics. Also, an L(M)M can be understood as a GL(M)M with an identity link function and a Gaussian distribution for the residuals.

adjectives are used, a single occurrence of a contracted or non-contracted form, etc. When such observations are grouped, it is often plausible to assume that some variance in the choice of the alternating forms or constructions occurs at the group-level and not at the level of observations. Groups can be defined by any linguistically relevant grouping factor (a categorical variable, also called a nominal variable), such as the individual speakers (or authors, writers, etc.), their sex and gender, the regions where they were born or live, social groups with which they identify, but also time periods, genres, styles, etc. If the concordance in a study contains, say, ten exemplars each written by ten speakers, then the speaker grouping factor has ten levels and defines ten groups. We know that preferences vary between speakers, and it is therefore reasonable to take care of this variance in our statistics in some way. The same goes for the other possible groups just mentioned. Furthermore, it is known that specific lexemes often have idiosyncratic affinities towards alternants in alternating constructions. Therefore, exemplars containing specific lexemes can also be treated as groups with considerable between-group variance. As an example from outside corpus linguistics, variation between participants is standardly modeled by including a random effect for speaker in experimental settings.

While random effects are often presented like this using conceptual arguments, the crucial question in specifying models is not whether to include these grouping factors at all, but rather whether to include them as fixed effects or as random effects.<sup>3</sup> Random effects structures are very suitable for accounting for group-level variation in regression, but contrary to formulaic recommendations such as “Always include random effects for speaker and genre!”, the

---

<sup>3</sup>I assume that no corpus linguist would make the decision *not* to include relevant factors in their models when the corpus contains the corresponding meta data or these meta data can be annotated reliably with acceptable effort.

choice between fixed and random effects can and should be made based on an analysis and understanding of the data set at hand and the differences and similarities in the resulting estimates. Subsection 2.2, 2.3, and 2.4 introduce three important points to consider about the structure of the data typically used in mixed modeling. This is intended to show readers that mixed or multilevel/hierarchical modeling is simply a matter of doing justice to the structure of the data. Then, Section 2.5 provide a moderately technical introduction to the important technicalities in mixed modeling, including a discussion of when a factor should be included as a random effect and when as a fixed effect. Section 3 then shows how mixed models are specified and interpreted using R.

## 2.2 Crossed and nested effects

It was established in the previous section that random effects are a means of accounting for group-level variance in regression models. This section briefly introduces a distinction that plays a role in modeling when there is more than one grouping factor (to be used either as a fixed or random effect). When this is the case, each pair of grouping factors can be *nested* or *crossed*. By way of example, we can group exemplars (such as sentences) by the individual speakers who wrote or uttered them, and we can group speakers by their region of birth. Such a data set would intrinsically be *nested*, as Table 2 illustrates.

Since speakers have a unique region of birth, Tyneside is the unique region value for the speakers Bay and Riley, and Greater London is the unique region value for Dale and Hayden. There cannot be exemplars where, for example, the speaker is Bay and the region is Greater London (assuming that speakers are uniquely identified by the labels in the middle column). In this example,

<b>Exemplar</b>	<b>Speaker</b>	<b>Region</b>
1	Bay	Tyneside
2	Bay	Tyneside
3	Riley	Tyneside
4	Riley	Tyneside
5	Dale	Greater London
6	Dale	Greater London
7	Hayden	Greater London
8	Hayden	Greater London

Table 1: Illustration of nested factors

<b>Exemplar</b>	<b>Speaker</b>	<b>Mode</b>
1	Bay	Spoken
2	Bay	Written
3	Riley	Spoken
4	Riley	Spoken
5	Dale	Written
6	Dale	Written
7	Hayden	Spoken
8	Hayden	Written

Table 2: Illustration of crossed factors

the region factor nests the speaker factor.<sup>4</sup> This example was chosen because the nesting is conceptually necessary. However, even when a data set has a nested structure by accident, standard packages in R will treat them as nested, and a closer look at data sets should be part of any protocol for using GLMMs in corpus studies (see Section 3.1).

When the grouped entities do not uniquely belong to grouping factors, the factors are *crossed*. Continuing the example, crossed factors for speaker and mode are illustrated in Table 2. While there are only spoken sentences by Riley and only written sentences by Dale in the sample, there is one spoken and one

<sup>4</sup>The exemplar index should not be called not a grouping factor because its values are unique, and sentences are thus not “nested”. They represent the basic level of observations, and at that level, there is nothing to group.

written sentence each by Bay and Hayden. There is a many-to-many relation between speakers and modes, which is characteristic of crossed factors. In Table 1, the relation between speakers and regions was many-to-one, which is typical of nested factors. In experimental settings, the design often makes sure that the combinations of nested or crossed factors are represented by equal numbers of observations (such that, for example, there is an equal number of written and spoken sentences from each speaker). Contrarily, the situation in Table 2 is typical of corpus studies where pseudo-random sampling from a pre-compiled corpus such as the BNC was used. This does not affect the practical modeling procedures much, especially when random factors are used, as will be shown below. However, practitioners must be aware of it when interpreting the data.

Finally, it should be noted that grouping factors can form hierarchical structures. When grouping factors are nested, there can be more than one level of nesting. Mode could nest genre if genres are defined such that each genre is either exclusively spoken or written, and in a given corpus, speakers might be nested within genres because each of them only contributed material to one genre.<sup>5</sup> Similarly, we might want to describe – in a given study on adjectives – adjectives as being either intersective or non-intersective. Within the two groups, a finer-grained semantic classification might be nested, which itself nests single adjective lexemes. This gives rise to potentially complex hierarchical structures, which can often be modeled more effectively using random

---

<sup>5</sup>In this example, the second level of nesting is not a conceptual necessity. In fact, it would be quite surprising if the real world were shaped like this. However, standard corpus compilation techniques might easily lead to a situation where exactly this is the case, simply because it is often difficult to sample texts and utterances from single speakers across a wide range of genres.

effect structures compared to fixed effect structures.

## 2.3 Hierarchical or multilevel modeling

This section introduces the idea – often ignored in introductory texts written by practitioners and handbook articles – that so-called random effects actually introduce new levels of modeling, or *secondary models*. It is argued that this is, again, not a technical thing but required by the structure of certain data sets. Assume that we wanted to account for lexeme-specific variation in a study on an alternation phenomenon by specifying the lexeme as a random effect in the model. Additionally, we suspect or know that a lexeme’s overall frequency influences its preferences for occurring in the construction alternants. Now, we could simply quantize the frequency variable and turn it into an ordinal variable (for example in the form of frequency bands) and interpret it as a grouping factor which nests the lexeme grouping factor. However, frequency obviously is a numerical and not a categorical variable, and by using it as a grouping factor we would destroy valuable information that is encoded in the data. A similar situation would arise in a study of learner corpus data with a learner grouping factor if we also knew that the number of years learners have learned a language influences their performance with regard to a specific phenomenon. It should have become clear that in such cases, variables like *frequency* and *number of learning years* are constant for each level of the grouping factor (*lexeme* and *learner*, respectively), but we cannot treat them as nesting grouping factors themselves. In other words, each lexeme has exactly one overall frequency, and each learner has had a fixed number of years of learning the language.<sup>6</sup>

---

<sup>6</sup>In the given example, things would get more complicated if the corpus contained data by single learners from different points in time. We simplify the scenario for the sake of an

Level of observations			Group level	
Exemplar	Givenness	NP length	Verb	Verb frequency
1	New	8	give	6.99
2	Old	7	give	6.99
3	Old	5	give	6.99
4	Old	5	grant	5.97
5	New	9	grant	5.97
6	Old	6	grant	5.97
7	New	11	promise	5.86
8	New	10	promise	5.86
9	Old	9	promise	5.86

Table 3: Illustration of a data set which requires multilevel modeling; lemma frequencies are logarithmized frequencies per one million tokens taken from ENCOW14A

Such variables are thus reasonably interpretable only at the group-level. Table 3 illustrates such a data set (fictional in this case). It might be a small fraction of the data used to predict whether a dative NP is used in the dative shift construction or not. The exemplar indices, again, simply identify single sentences containing one of the constructions of interest. The discourse status obviously varies at the level observations, and so does the NP length in syllables. To capture verb lemma specific tendencies, a verb lemma grouping factor is added, and the verb lemma frequency necessarily varies at the group level because each lemma has a unique frequency. In such cases, an adequately specified multilevel model uses the group-level variables to partially predict the tendency of the grouping factor. Put differently, the idiosyncratic effect associated with a lexeme, speaker, genre, etc. is split up into a truly idiosyncratic preference and a preference predictable from group-level factors. This is achieved, in fact, by specifying another model (a linear model) that predicts the group-level random effect itself, and the second-level predictor is a fixed

---

easier-to-follow introduction.



effect in this model. Such second-level models can even contain modeled effects themselves, giving rise to third-level models, and so on. In a way, the data look similar to multilevel nesting, but (1) second-level models can account for continuous numerical predictors at the group-level, which nesting cannot, and (2) there might be situations where specifying even categorical second-level grouping factors as fixed effects in a second-level model is more appropriate than adding nested random effects (see Section 2.5).

As in the case of nested vs. crossed factors, standard packages in R often take care of hierarchical modeling automatically, given that the data are structured and are specified accordingly (see Section 3). This might, however, lead to situations where practitioners specify multilevel models without even knowing it, which in turn can lead to misinterpretations of the results. Therefore, multi-level modeling will be introduced as the more general framework for so-called mixed effects models in Sections 2.5 and 3.

## 2.4 Random slopes as interactions

Before moving on to the more technical discussion of hierarchical model specification in Section 2.5, one more basic concept will be discussed in this section, namely the data patterns that gives rise to *varying intercepts* and *varying slopes*. Varying intercepts are an adequate modeling tool when the overall tendency in the outcome variable changes with the levels of the grouping factor. It is shown that random slopes are just another way of modeling an interaction between influencing factors.

For the sake of the argument, we assume that we are looking at an alternation phenomenon like the dative alternation, wherein we are interested in the probability that, under given circumstances, the dative shift construction is chosen.

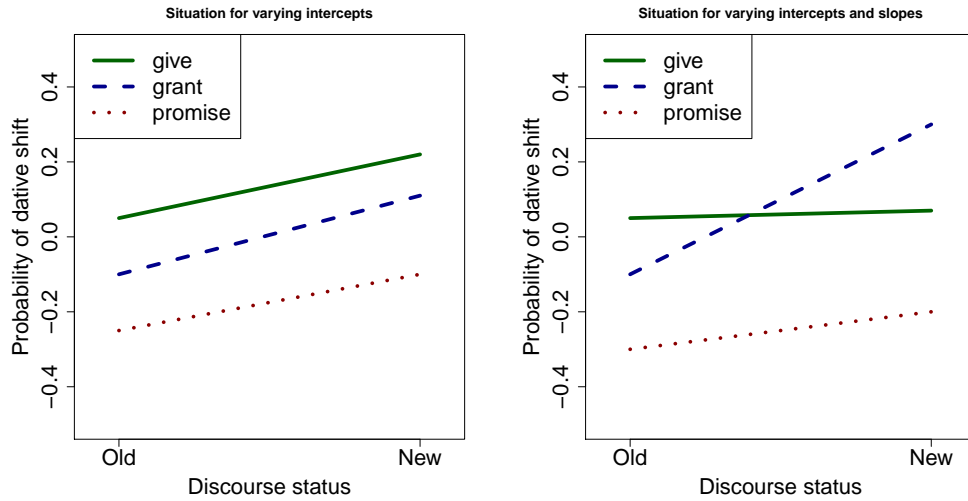


Figure 1: Illustration of data in situations for varying intercepts or varying intercepts and additional varying slopes

Looking at data set, it turns out that the probability of the dative shift changes for *old* and *new* dative NPs. The verb lemma – a typical candidate to be used as a random effect – also influences the probability of either variant being used. The situation can now be as in the left or the right panel of Figure 1. In the situation illustrated in the left panel, the overall level in probability changes with the verb lemma, but for each verb lemma, the values change roughly accordingly in exemplars with old and new dative NPs. Note that the lines are not perfectly parallel because the figure is supposed to be an illustration of a data set rather than a fitted model, and we always expect some chance variation in data sets. In the situation depicted in the right panel, however, the overall levels are different between lemmas, but the lemma-specific tendencies also vary between exemplars with old and new NPs. This is actually nothing but an interaction between two factors (verb lemma and discourse status). However, if the verb lemma factor is used as a random effect grouping factor, the interaction is modeled as a so-called *random slope*. In the next section, it is shown how all the different types of data sets discussed so far can be modeled using

fixed effects models or, alternatively, using mixed effects. Which one is more appropriate will be argued to be better understood as a technical rather than a conceptual question.

## 2.5 Model specification

In this section, it is discussed how the specification of mixed models differs from that of fixed effects models, and that for each model with random effects there is an alternative models with only fixed effects. A major focus is on the question of when to use fixed and random effects. The amount of technicality and notation is kept at the absolute minimum, but a few notational conventions are introduced as the absolute minimum required to understand both the output of `lme4` and other packages in R as well as the literature on mixed models. To make successful *practical* use of mixed models, some level of fundamental understanding is required. This section is based mostly on Part 2A of Gelman & Hill (2006) (pp. 235–342).

### 2.5.1 Simple random intercepts

Readers with experience in fixed effects modeling (see the previous chapter in this handbook) should see that a grouping factor encoding the verb lemma, the speaker, the mode, the genre corresponding to a corpus exemplar (and all the other grouping factors discussed in the previous sections) could be specified as a normal fixed effect in a GLM. In such a case, each of the  $m$  levels of the speaker factor is dummy-coded, and for all but one of these binary dummy variables, a coefficient is estimated. Logistic regression examples are used throughout this section, and we begin with the fictional corpus study of the dative alternation introduced in Sections 2.3 and 2.4. We begin with a minimal model with only

the dummies of the lemma grouping factor and one other (binary) predictor, namely discourse status. There are  $m$  verb lemmas (i. e., groups) and  $n$  observations. As index variables, we use  $j$  for groups and  $i$  for observations. In general,  $\alpha$  is used for intercepts and  $\beta$  for coefficients. A specification of such a model is given in (1).

$$Pr(y^i = 1) = \text{logit}^{-1}(\alpha_0 + \beta_d \cdot x_d^i + \beta_{l_1} \cdot x_{l_1}^i + \beta_{l_2} \cdot x_{l_2}^i + \cdots + \beta_{l_{m-1}} \cdot x_{l_{m-1}}^i) \quad (1)$$

This models the estimate of the probability ( $Pr$ ) that in observation  $i$ , the outcome variable  $y^i$  is 1, i. e., that dative shift occurs.  $\alpha_0$  is the intercept,  $\beta_d$  is the coefficient for the effect of discourse status.  $x_d^i$  is the value of the variable that encodes discourse status for exemplar  $i$  (0 for discourse-old NPs and 1 for discourse-new NPs). Furthermore,  $\beta_{l_j}$  are the coefficients for the lemma dummy variables. Finally,  $x_{l_j}^i$  is the value (0 or 1) for lemma  $j$  and observation  $i$ . If in exemplar 64, the lemma is *give* and *give* is encoded as group 12, then  $i = 64$ ,  $j = 12$ , and  $x_{l_{12}}^{64} = 1$ , whereas all  $x_{l_j}^{64} = 0$  with  $j \neq 12$ .<sup>7</sup> Because one verb lemma dummy variable is on the intercept  $\alpha_0$  and thus used as a reference, we only estimate  $m - 1$  instead of  $m$  coefficients, i. e.,  $j = 1, \dots, m - 1$ .<sup>8</sup> The function  $\text{logit}^{-1}$  is the *link function*, and its argument is the *linear term* of the model. It is obvious that in such a model, the effect of each verb lemma is treated as a

---

<sup>7</sup>It is unfortunate that multiple indexation is required to such an extent. However, any alternative notation has at least the same potential to confuse readers.

<sup>8</sup>It is often not explained in text books for practitioners that using one dummy as a reference level is necessary because otherwise infinitely many equivalent estimates of the model coefficients exist because one could simply add a constant to the intercept and subtract it from the dummies. However, the estimator works under the assumption that there is a unique maximum likelihood estimate for the coefficient matrix.

fixed population parameter, exactly the same as the effect of discourse status.

The coefficient  $\beta_d$  is estimated in exactly the same way as each  $\beta_{l_j}$ .

If we treat the same grouping factor as a random intercept, we let the intercept vary by group, and we give the varying intercepts a distribution instead of estimating  $m - 1$  coefficients. This is the only relevant difference between a fixed effect and a random effect. The model specification then looks like in (2).

$$P(y^i = 1) = \text{logit}^{-1}(\alpha_l^{j[i]} + \beta_d \cdot x_d^i) \quad (2)$$

We now have an intercept  $\alpha_l^{j[i]}$  that varies by group (instead of one term with its own coefficient per group). We use the notation  $\alpha_l^{j[i]}$  (borrowed in a modified form from Gelman & Hill 2006) to indicate that the correct  $j$ -th lemma intercept is chosen for the  $i$ -th observation. For example, if in exemplar 64, the lemma is *give*, which is group 12, then  $i = 64$  and  $j[64] = 12$  (i. e., the group appropriate for exemplar 64 is group 12), and  $\alpha_l^{j[64]} = \alpha_l^{12}$ . The term  $\beta_d \cdot x_d^i$  for the effect of discourse status remains unchanged when going from (1) to (2). Crucially, instead of estimating a batch of coefficients for the lemma effect,  $\alpha_l$  is itself modeled, and random terms are predicted for each level of the grouping factor. For this, the assumption in (3) is made.

$$\alpha_l^j \sim N(\mu_l, \sigma_l^2) \quad (3)$$

This is standard notation to indicate that the values of  $\alpha_l^j$  follow a normal distribution with mean  $\mu_l$  and a variance of  $\sigma_l^2$ . In fact, we can regard (3) as a minimal second-level model already, although one which simply predicts varying intercepts from a normal distribution. All more complex models to be discussed below are extensions of this approach. In the next section, the consequences

of going from a fixed effect to a random effect are discussed.

### 2.5.2 Choosing between random and fixed effects

There are primarily two points to consider which influence the decision to use random effects. First, the variance in the intercepts (and for random intercept-random slope models also the covariance between intercepts and slopes) needs to be estimated. Second, the random intercepts can be understood as a compromise between fitting separate models for each group of the grouping factor (*no pooling*) and fitting a model while ignoring the grouping factor altogether (*complete pooling*), see Gelman & Hill (2006). While all conditions which were discussed in the previous chapter (independence of observations, non-collinearity, etc.) must also be met by hierarchical models, these two points add some additional conditions.

As was stated above in (3), the random intercepts are assumed to follow a normal distribution, and the variance  $\sigma_I^2$  needs to be estimated with sufficient precision. From the estimated variance and the data, the estimator then predicts the *conditional modes* in GLMMs (*conditional means* in LMMs) for each group (see Bates 2010), which is the numerical value which software packages like lme4 produce, and these values are sometimes sloppily called “random effects” by practitioners. This procedure, however, requires that the number of groups must not be too low to effectively achieve this. As a rule of thumb, fewer than five levels means that a grouping factor should be included as a fixed effect, regardless of its conceptual nature. Even if there is a default recommendation to use a speaker grouping variable as a random effect, it is ill-advised to do so if there are exemplars from less than five speakers in the sample. Very often, the estimator will fail anyway under such conditions. At the same time, the

estimation of coefficients for a (dummy-coded) fixed effect becomes less precise and feasible with large numbers of levels, and at a certain point, a random effect might be the better or only option for technical reasons.

If, however, the number of groups is reasonably large, the next thing to consider is the number of observations per group. As mentioned at the outset of this section, alternatives to using a random effect would be to estimate a separate model for each level of the grouping factor, or to include it as a fixed effect. In both cases the effects are not treated as random variables, and fixed values per group are estimated without taking the between-group variance into account. The conditional modes are *shrunk* towards the overall intercept as a result (an effect called *shrinkage*), and they tend to be more leveled compared to fixed effects estimates. This effect becomes considerably weaker with larger per-group sample sizes. With relatively low numbers of observations per group, on the other hand, fixed effect estimates tend to become inexact and will probably be dismissed because of growing uncertainty in the estimate (large confidence intervals, non-significance). In this case, the conditional modes (or means) of random intercepts will just be shrunk more towards the overall intercept. This considered, including a grouping factor as a random effect might be a way to using it when the estimation of fixed effects fails.

This section closes with an illustration.<sup>9</sup> For this, 1,000 data sets were simulated which corresponded to the model in (4) and (5). We drop the subscripts on  $\alpha$ ,  $\beta$ , and  $\mu$  for convenience since there is only one random intercept.

---

<sup>9</sup>The code for these and other simulations is available under a Creative Commons Attribution license: <https://github.com/rsling/Rstuff/tree/master/simulations/glmm>

$$P(y^i = 1) = \text{logit}^{-1}(\alpha^{j[i]} + \beta_1 \cdot x_1^i + \beta_2 \cdot x_2^i) \quad (4)$$

$$\alpha^j \sim N(\mu, \sigma) \quad (5)$$

Again, this could be a model of a binary alternation.  $x_1$  was a binary variable (such as singular/plural) and  $x_2$  a continuous variable (such as the length of the construction). Since the data were simulated, the parameters to be estimated were known:  $\beta_1 = 0.8$ ,  $\beta_2 = -1.3$ ,  $\mu = 0$ ,  $\sigma = 1.5$ . The number of groups was set to 5 (the recommended minimum), the simulated values of the grouping factor were identical in each simulation, and there were 20 observations per group. Figure 2 shows the distribution of the group levels based on the conditional modes predicted of all but the first group in the 1,000 simulations. Figure 3 shows the estimated coefficients of a model where the grouping factor was added as a fixed effect.<sup>10</sup>

The per-group predictions lean slightly towards 0 in the GLMM (Figure 2), but the fixed effects estimates in the GLM are prone to massive misestimations (see the large spread and increased 95% intervals). Even with as few as five levels of a grouping factor, however, random and fixed effects lead to very similar results, albeit with different advantages and disadvantages. A few more differences will be discussed in Sections 2.5.3 and 2.5.4.

---

<sup>10</sup>The plots do not show the distribution of the raw conditional modes and coefficient estimates of the fixed effects. Rather, the overall intercept was taken into account, and the plots thus show the distribution of the per-group prediction of the models, taking only the grouping factor into account. This is what was pre-specified in the simulations.



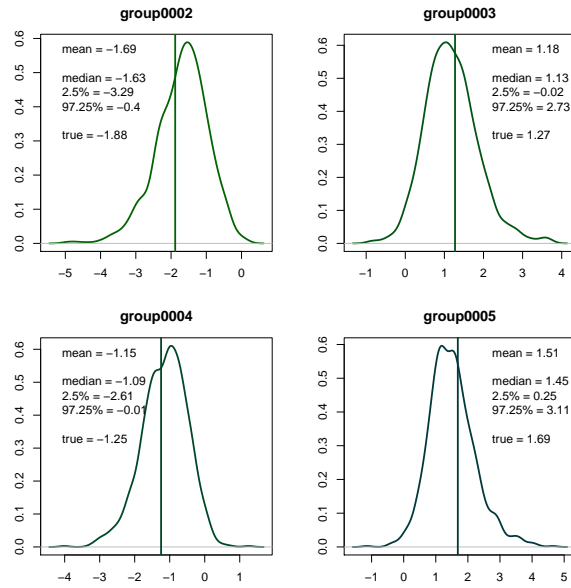


Figure 2: Group levels in sample GLMM based on predicted random effect (conditional mode); 5 groups; 20 observations per group; 1,000 simulations; the horizontal line marks the true value

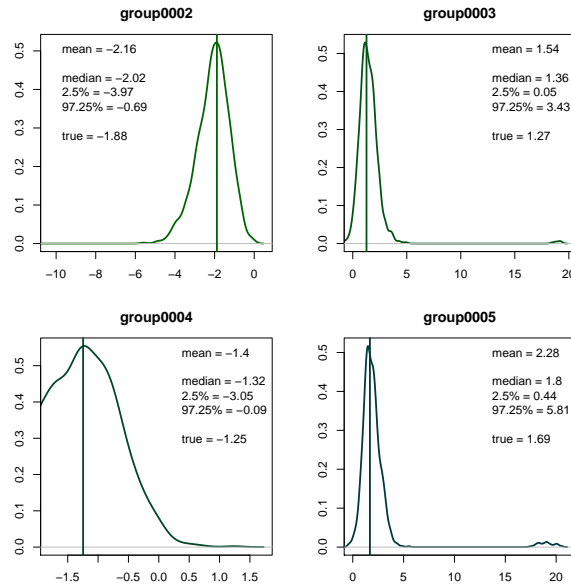


Figure 3: Estimated fixed effects for the grouping factor in sample GLM; 5 groups; 20 observations per level of the grouping factor; 1,000 simulations; the horizontal line marks the true value

### 2.5.3 Significance testing, model selection and coefficients of determination

One commonly given reason to use a random effect is that “the researchers are not interested in the individual levels of the random effect factor” (or variations thereof). Such recommendations should be taken with a grain of salt. Gelman & Hill (2006: 245–247) summarise the diverging recommendations and their various justifications, concluding that there is essentially no clean conceptual way of deciding what should be a random effect and what a fixed effect (see Section 2.5.2 for the more technical approach). Clearly, however, it is not adequate to do any kind of significance test on the levels of the random effect factor, as is customary for fixed effects. The conditional modes are not estimates, and only for estimates is significance testing a well-defined activity. There are ways of calculating confidence intervals (see Section 3.1), but they should not be misused for talking about significance. If this is desired, fixed effects are the way to go. Not doing it does, however, not mean that the researcher is not interested in the individual conditional modes of the random effect. Additionally, a random effect allows the researcher to quantify the between-group variance, which is not possible in the same way with fixed effects.

A related question is model selection, i. e., whether the inclusion of the random effect improves the model. First, the between-group variance should be checked. If it is close to 0, there is most likely not much going on between groups, or there simply was not enough data to estimate the variance. In LMMs, it is possible to compare the residual (observation-level) variance with the between-group variance to see which one is larger, and to which degree. If, for example, the residual variance is  $\sigma_\epsilon = 0.2$  and the between-group variance is  $\sigma_\alpha = 0.8$ , then we can say that the between-group variance is four

times larger than the residual variance, which indicates a high importance of the random effect. This comparison is not possible in most GLMMs, because the residuals do not have the same straightforward interpretation as in LMMs. Traditional model selection for random effects can be achieved through a likelihood ratio (LR) test comparing a model including the random effect and a model not including it. Such models, where one is strictly a simplification of the other, are called *nested model models* (not to be confused with *nested effects* discussed in Section 2.2). A much more robust option uses a parametric bootstrap replacements for the LR test (see Section 3.1). It is *not* appropriate to compare a model with a random effect and a model with the same factor as a fixed effect using *any* test or metric (including information criteria).

Coefficients of determination ( $R^2$ ) should not be used for model selection, because they usually do not penalize complexity enough (i. e., the coefficients rise with added complexity). They can still be used to give some idea of the model fit, however. For GLMMs, Nakagawa & Schielzeth (2013) have proposed a method that distinguishes between *marginal*  $R^2$  (fixed-effects-only) and *conditional*  $R^2$  (including random effects). This has become a de facto standard, and we now show its consistency with Nagelkerke's  $R^2$  for GLMs. Using the simulations described in the last section, Figures 4 and 5 show that the marginal  $R^2$  is roughly the same as Nagelkerke's  $R^2$  in a GLM which ignores the grouping factor, and the conditional  $R^2$  is roughly the same as Nagelkerke's  $R^2$  in a GLM which includes the grouping factor as a fixed effect, given that the estimators converge for all these three types of models.

Finally, it should be mentioned that adding a random intercept does not change the estimates nor the inferences (p-values) for fixed effects, at least if the grouping factor and the fixed effect are independent. Figure 6 compares the p-values

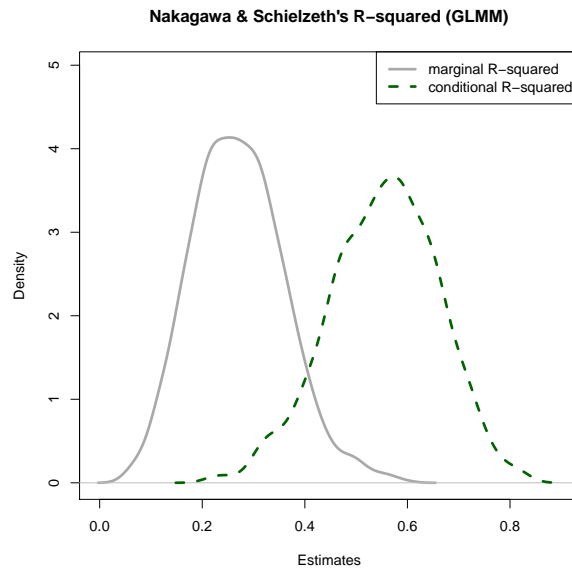


Figure 4: Distribution of Nakagawa & Schielzeth's  $R^2$  in the simulations described in Section 2.5.2

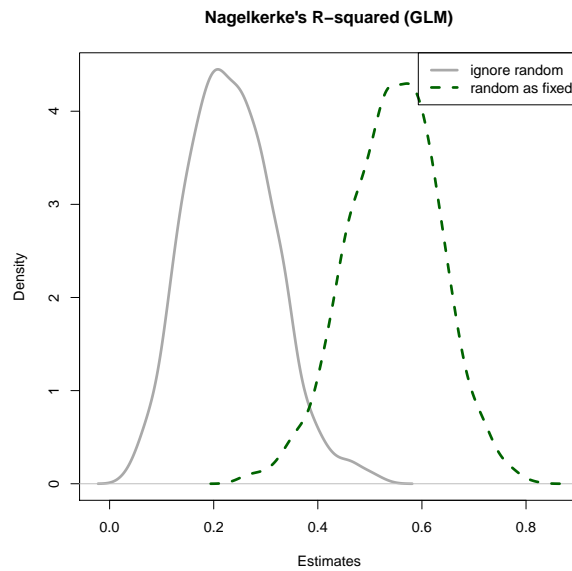


Figure 5: Nagelkerke's  $R^2$  in the simulations described in Section 2.5.2 for a GLM that ignores the grouping factor and a model that includes it as a fixed effect

obtained for the two fixed effects in the 1,000 simulations described above for the GLMM and the GLM ignoring the grouping factor. Notice that the  $\beta_2$  effect has the higher effect strength, and p-values are generally lower. If fixed effects and random effects are not independent, adding a varying slope should be considered, which is described in Section 2.5.4.

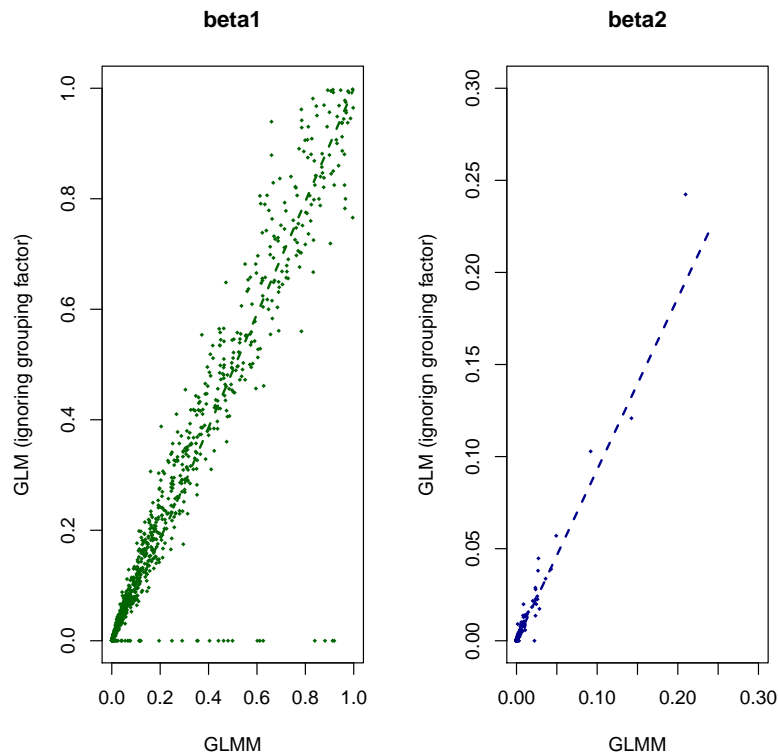


Figure 6: Comparison between the p-values of the fixed effects in the simulations described in Section 2.5.2 for the GLMM and the GLM ignoring the grouping factor; the line was drawn using LOWESS (locally weighted scatter-plot smoothing)

#### 2.5.4 More complex models

This section shows briefly how more complex models are specified before Section 3.1 demonstrated how all models discussed are estimated in R, including interpretations of the output.

**Varying intercepts and slopes** A varying slop is an interaction between a fixed effect and a random effect, see Section 2.4. While it is possible to have just a varying slope, this is rarely useful, and we discuss only varying-intercept and varying-slope (VIVS) models. We extend the simple model from (2), and the fixed effect coefficients for which a random slope is specified simply receive group indices; see (6). Put simply, instead of a fixed coefficient for the fixed effect, coefficients are predicted and assumed to come from a random (normal) distribution. We use  $\beta_{d:l}$  to denote the coefficient for discourse status varying by lemma.

$$P(y^i = 1) = \text{logit}^{-1}(\alpha_l^{j[i]} + \beta_{d:l}^{j[i]} \cdot x_d^i) \quad (6)$$

A source of problems in VIVS models is the fact that in addition to the variance in the intercepts and slopes, the covariance between them has to be estimated. If in groups with a higher-than-average intercept, the slope is also higher than average, they are positively correlated, and vice versa. These relations are captured in the covariance. Condition (7) is added, where the indices  $l$  and  $d : l$  are omitted for readability.

$$\begin{pmatrix} \alpha^j \\ \beta^j \end{pmatrix} \sim \left( \begin{pmatrix} \mu_\alpha \\ \mu_\beta \end{pmatrix}, \begin{pmatrix} \sigma_\alpha^2 & \rho\sigma_\alpha\sigma_\beta \\ \rho\sigma_\alpha\sigma_\beta & \sigma_\beta^2 \end{pmatrix} \right) \quad (7)$$

(7) says that the joint distribution of the intercepts  $\alpha^j$  and the slopes  $\beta^j$  follow a bivariate normal distribution with means  $\mu_\alpha$  and  $\mu_\beta$ . The variance in the intercepts is  $\sigma_\alpha$ , the variance in the slopes is  $\sigma_\beta$ , and the covariance between them is  $\rho$ . Figure 7 shows the bivariate density distributions for two negatively correlated, non-correlated, and positively correlated normally distributed variables.

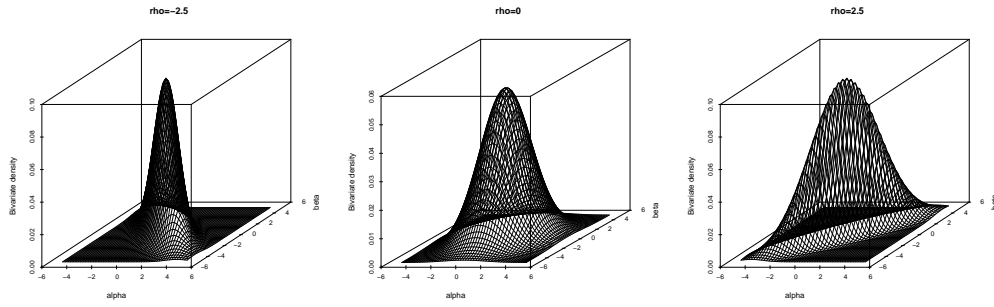


Figure 7: Bivariate normal density distribution with different correlations ( $\rho$ );  $\sigma_\alpha = \sigma_\beta = 3$ ,  $\mu_\alpha = \mu_\beta = 0$

Importantly, the number of parameters to estimate obviously increases with more complex model specifications, and the estimation of the parameters in the presence of complex variance-covariance matrices requires considerably more data than estimating a single variance parameter. The estimator might converge, but typically covariance estimates of  $-1$  or  $1$  indicate that the data was too sparse for a successful estimation of the parameter. In this case, the model is *over-parametrized* and needs to be simplified.

### **Nested and crossed random effects**

### **Second-level predictors**

### **Remarks about models for time series**

## 3 Packages and tools in R

### 3.1 Using lme4

### 3.2 Bootstrap methods for lme4

## 4 Further reading

## References

- Bates, Douglas M. 2010. *lme4: Mixed-effects modeling with R*. ().
- Gelman, Andrew & Jennifer Hill. 2006. *Data analysis using regression and multilevel/hierarchical models*. Cambridge: Cambridge University Press.
- Nakagawa, Shinichi & Holger Schielzeth. 2013. A general and simple method for obtaining  $R^2$  from generalized linear mixed-effects models. *Methods in Ecology and Evolution* 4(2). 133–142.