

DRAFT of November 6, 2017

# Mixed-effects regression modeling

to appear in: St. Gries & M. Paquot (eds.), Practical Handbook of Corpus Linguistics

Roland Schäfer

Freie Universität Berlin

November 6, 2017

## 1 Introduction

Mixed effects modeling – alternatively called *hierarchical* or *multilevel modeling* is a straightforward extension of (generalised) linear modeling as discussed in the previous chapter. A common characterisation of mixed-effects modeling is that it accounts for situations where observations are *clustered* or *come in groups*. In corpus linguistics, there could be clusters of observations defined by individual speakers, registers, genres, modes, lemmas, etc. Instead of estimating coefficients for each level of such a grouping factor (so-called *fixed effects*), in a mixed model they can be modeled as normally distributed random variables (so-called *random effects*) with predictions being made for each group. This chapter introduces readers to the situations where mixed-effects modeling is useful or necessary. The proper specification of models is discussed, as

well as some model diagnostics and ways of interpreting the output. Readers are assumed to be familiar with the concepts covered in the previous chapter.

## 2 Fundamentals

### 2.1 When are random effects useful?

#### 2.1.1 Introduction to random effects

(Generalised) Linear Mixed Models (GLMMs) are an extension of (Generalised) Linear Models. They add what are often called *random effects* and *mix* them with the normal predictors (*fixed effects*) as used in GLMs. Alternatively, statisticians speak of *multilevel models* or *hierarchical models* (Gelman & Hill 2006), a terminology to be explained in Section 2.1.3.

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 a study. In corpus linguistics, single observations can be understood as single lines in a concordance, and they represent, for example, clauses or sentences in which one of the alternants of a morpho-syntactic alternation occurs, NPs where two pre-nominal adjectives are used, single occurrences of a contracted or non-contracted form, etc. When such observations are grouped, it is often plausible to assume that there is some variance in the choice of the alternating forms or constructions at the group-level. If this is the case and the grouping factor is not included in the model, the error terms within the groups will be correlated. Since the estimators used for estimating the parameters of GLMs work under the assumption of non-correlated errors, standard errors for model coefficients will typically be estimated as smaller than they

nominally are, leading to increased Type I error rates in inferences about the coefficients.<sup>1</sup> This gets even worse when there are within-group tendencies regarding the direction and strength of the influence of the other regressors, i. e., when there is an interaction between them and the grouping factor (e. g., Schielzeth & Forstmeier 2009). This is why known variation by group should always be accounted for in the model, and random effects are often a convenient way to do so.

Groups can be defined by any linguistically relevant grouping factor, such as the individual speakers (or authors, writers, etc.), the regions where they were born or live, social groups with which they identify, but also time periods, genres, styles, etc. Furthermore, it is known that specific lexemes often have idiosyncratic affinities towards alternants in alternating constructions. Therefore, exemplars containing specific lexemes also constitute groups.

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. Random effects structures are very suitable for accounting for group-level variation in regression, but while formulaic recommendations such as “Always include random effects for speaker and genre!” provide useful guidance for beginners, the 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 models. The remainder of Section 2.1 introduces three important points to consider about the structure of the data typically used in mixed modeling. Then, Section 2.2 provides a moderately technical introduction to modeling. Section 2.3 shows how mixed models are specified using the `lme4` package in R, and Section 2.4 deals with

---

<sup>1</sup>Type I errors occur when the null hypothesis is rejected although it is true.

Exemplar	Speaker	Region
1	Daryl	Tyneside
2	Daryl	Tyneside
3	Riley	Tyneside
4	Riley	Tyneside
5	Dale	Greater London
6	Dale	Greater London
7	Reed	Greater London
8	Reed	Greater London

Table 1: Illustration of nested factors

the interpretation of the output.

### 2.1.2 Crossed and nested effects

This section discusses a distinction that arises when there is more than one grouping factor. 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 Daryl and Riley, and Greater London is the unique *region* value for Dale and Reed. In this example, the region factor nests the speaker factor. 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 also treat them as nested (see Section 2.3).

When the grouped entities (themselves groups) do not uniquely belong to levels of the grouping factor, 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

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

Table 2: Illustration of crossed factors

sample, there is one spoken and one written sentence each by Daryl and Reed. 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 is many-to-one, which is typical of nested factors.

With more than two grouping factors, 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. 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. However, not all of these structures should be modeled as nested random effects. In the latter case, for example, the low number of levels in one factor (intersectivity with just two levels) predestines it as a second-level predictor rather than a nesting factor; see Section 2.1.3.

### 2.1.3 Hierarchical or multilevel modeling

This section describes the types of data to be used in true multilevel models. Let us 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

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 fictional data set which requires multilevel modeling; lemma frequencies are logarithm-transformed frequencies per one million tokens taken from ENCOW14A

in the model. Additionally, we suspect or know that a lexeme’s overall frequency influences its preferences for occurring in the construction alternants. 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. In such cases, variables like *frequency* and *number of learning years* are constant for each level of the grouping factor (*lexeme* and *learner*, respectively). 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>2</sup>

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 ditransitive verb is used in the dative shift construction or not. The discourse status and the NP length status obviously vary at the level observations. To capture verb lemma specific tendencies,

<sup>2</sup>In the given example, things would get more complicated if the corpus contained observations of single learners at different points in time. We simplify the scenario for the sake of an easier-to-follow introduction. See also the last subsection of Section 2.2.4.

a verb lemma grouping factor is added. 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 variables. This is achieved by specifying a second (linear) model which predicts the group-level random effect itself. Such second-level models can even contain modeled effects themselves, giving rise to third-level models, and so on. 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 2.2).

As in the case of nested vs. crossed factors, standard packages in R usually take care of hierarchical modeling automatically, given that the data are structured and are specified accordingly. 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.

#### 2.1.4 Random slopes as interactions

This section introduces 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.

We assume that we are looking at an alternation phenomenon like the dative

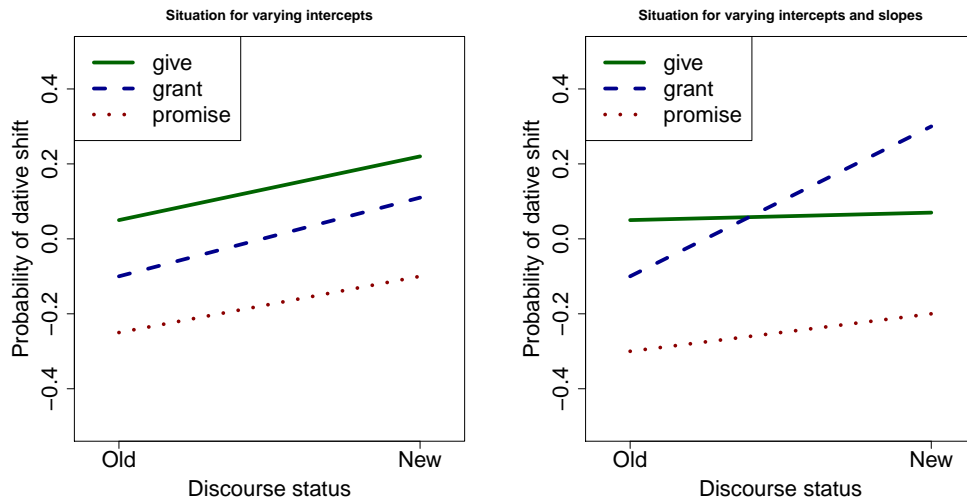


Figure 1: Illustration of fictional data in situations for varying intercepts or varying intercepts and additional varying slopes

alternation, wherein we are interested in the probability that, under given circumstances, the dative shift construction is chosen. In the examination of the data, it turns out that the probability of the dative shift changes for *old* and *new* dative NPs. The verb lemma 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 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 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 models. Which one is more appropriate will be argued to be better understood as a technical rather than a conceptual question.

## 2.2 Model specification and modeling assumptions

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. It is based mostly on Part 2A of Gelman & Hill (2006) (pp. 235–342).

### 2.2.1 Simple random intercepts

Readers with experience in fixed effects modeling should be able to see that a grouping factor encoding the verb lemma 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 2.1.3 and 2.1.4. We first specify 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).

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

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 the 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  $j$  in 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  $j \neq 12$ . 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>3</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 fixed population parameter, exactly the same as the effect of discourse status. 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 relevant difference between a fixed effect and a random effect. The model specification then looks like in (2).

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

---

<sup>3</sup>Picking one dummy as a reference level is necessary because otherwise infinitely many equivalent estimates of the model coefficients exist because one could simply add any arbitrary constant to the intercept. However, the estimator works under the assumption that there is a unique maximum likelihood estimate.

We now have an intercept  $\alpha_l^{j[i]}$  which 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 random effect. 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.

### 2.2.2 Choosing between random and fixed effects

There are primarily two points to consider which influence the decision whether 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: Ch. 12).

As was stated above in (3), the random intercepts are assumed to follow a normal distribution, and the variance  $\sigma_l^2$  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: Ch. 1), which is the numerical value which software packages like lme4 produce for each level of the grouping factor. 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. Along the same lines, mode (typically spoken vs. written) is no suitable grouping factor for use as a random effect.

If, however, the number of groups is reasonably large, the next thing to consider is the number of observations per group. 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 coefficients per group are estimated without taking the between-group variance into account. With a random effect, however, the conditional modes/means are pulled (*shrunk*) towards the overall intercept (*shrinkage*). When there the number of observations in a group is low, the conditional mode/mean is simply shrunk more strongly, predicting only a small deviation from the overall tendency. Fixed effect estimates, on the other hand, become inexact and will probably be dismissed because of growing uncertainty in the estimate (large confidence intervals, non-significance) when

the number of observations in a levels is low. Put differently, low numbers of observations in all or some groups are often detrimental for using fixed effects grouping factors. Random effects can deal with situations like this much better because of shrinkage. On the downside, a conditional mode that was strongly shrunken (due to a low number of observations) cannot be distinguished straightforwardly from a conditional mode of a group which simply does not deviate a lot from the average tendency. For fixed effects, we have both a parameter estimate and a possible significance test, but for random effects, we only have the prediction of the conditional mode/mean. Still, including a grouping factor as a random effect might be the only way of using it at all when the estimation as a fixed effect fails.

This section closes with an illustration.<sup>4</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.

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

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

Again, this could be a model of a binary alternation.  $x_1$  is a binary variable (such as given/new) and  $x_2$  a continuous variable (such as NP length). Since the data were simulated, the parameters to be estimated are known:  $\beta_1 = 0.8$ ,  $\beta_2 = -1.3$ ,  $\mu = 0$ ,  $\sigma = 1.5$ . The number of groups was set to 5, 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 lev-

---

<sup>4</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>

els based on the conditional modes predicted of all but the first group in the 1,000 simulations. Figure 3 shows the group estimates from a model where the grouping factor was added as a fixed effect.<sup>5</sup>

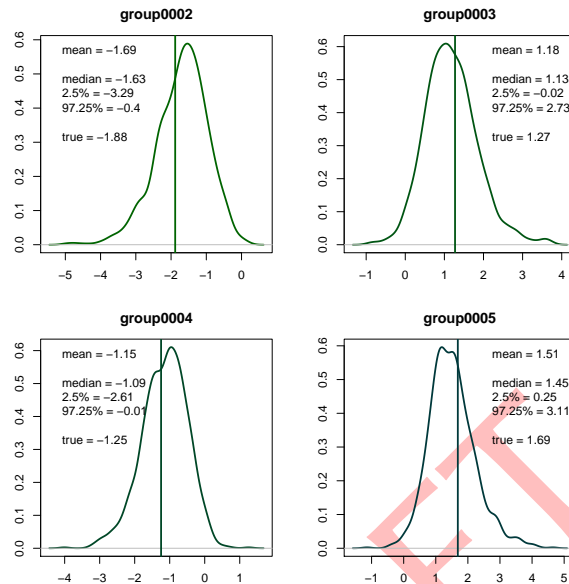


Figure 2: Group levels in sample GLMM based on predicted random effect (conditional mode); 5 groups; 20 observations per group; 1,000 simulations; the horizontal line marks the true value

The per-group predictions lean slightly towards 0 in the GLMM (Figure 2), but the fixed effects estimates in the GLM are prone to occasional misestimations. While the overall spread is roughly the same for both approaches, the GLM estimates have massive outliers (down to approximately  $-10$  and up to  $20$ ), which leads to slightly larger 95% percentile intervals. In this scenario (crafted to work reasonably well with a fixed or a random effect), however, random and fixed effects lead to very similar results, albeit with different advantages and

<sup>5</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, all other things being equal. This is what was pre-specified in the simulations.

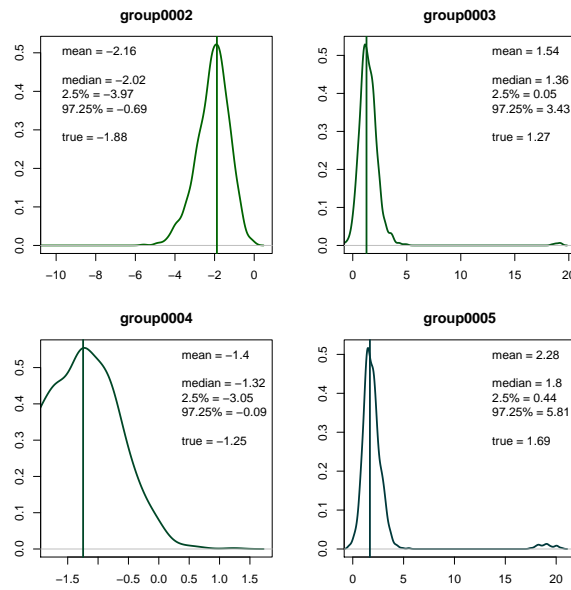


Figure 3: Estimated fixed effects for the grouping factor in sample GLM; 5 groups; 20 observations per level of the grouping factor; 1,000 simulations; the horizontal line marks the true value

disadvantages. More differences will be discussed in Sections 2.2.3 and 2.2.4.

### 2.2.3 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 and partially contradicting recommendations for what should be a random effect along with their motivations. They conclude that there is essentially no principled conceptual or mathematical way of deciding what should be a random effect and what a fixed effect. In this chapter, a more technical approach (which favours the solution that leads to the more robust model estimates) was therefore suggested. However, it is not adequate to do any kind of significance test on the levels of

the random effect because they are not estimates in the conceptual and technical sense.<sup>6</sup> There are ways of calculating *prediction intervals* (which are not the same as confidence intervals) for conditional modes in order to specify the quality of the fit (see Section 2.3), but they should not be misused for talking about significance. Not doing significance tests for single levels of the grouping factor does, however, not mean that the researcher is not interested in the individual conditional modes, which is proven by the fact that they are often reproduced in research papers, for example in the form of a dot plot. Also, the simulation in Section 2.2.2 shows that we can use a random effect and still get a good idea of the per-group tendencies. 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 quality. It is recommended here to include all conceptually necessary random effects and only remove them if they have no effect. To check whether this is the case, the estimated between-group variance is the first thing to look at. 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 would indicate a high importance of the random effect. This comparison is impossi-

---

<sup>6</sup>Essentially, we do not assume them to be fixed population parameters, which would be the case for estimates such as fixed effects coefficients.



ble in GLMMs because their (several types of) residuals do not have the same straightforward interpretation as in LMMs. Furthermore, likelihood ratio (LR) tests are available for comparing a model including the random effect and a model not including it. Such pairs of models, where one is strictly a simplification of the other, are called *nested models* (not to be confused with *nested effects* discussed in Section 2.1.2). An alternative option are parametric bootstrap replacements for the LR test. It is not appropriate to compare a GLMM with a random effect and a GLM with the same factor as a fixed effect using any test or metric (including so-called information criteria such as Akaike's or Bayes').

Coefficients of determination (pseudo- $R^2$ ) can be used to give some idea of the overall model fit. For GLMMs, Nakagawa & Schielzeth (2013) have proposed a method that distinguishes between *marginal  $R^2$*  (only fixed effects) and *conditional  $R^2$*  (fixed and random effects). This has become a de facto standard, and we now show its consistency with Nagelkerke's  $R^2$  for GLMs. Using the simulated data described in the last section, Figures 4 and 5 show that the marginal  $R^2$  for a GLMM estimate is roughly the same as Nagelkerke's  $R^2$  for a GLM estimate where the grouping factor is ignored. Also, the conditional  $R^2$  for a GLMM estimate is roughly the same as Nagelkerke's  $R^2$  for a GLM estimate which includes the grouping factor as a fixed effect. It should be noted that the simulations were explicitly designed such that the grouping factor (five levels with enough observations per level) could be used successfully as a fixed effect or a random effect, which is usually not the case with real data.

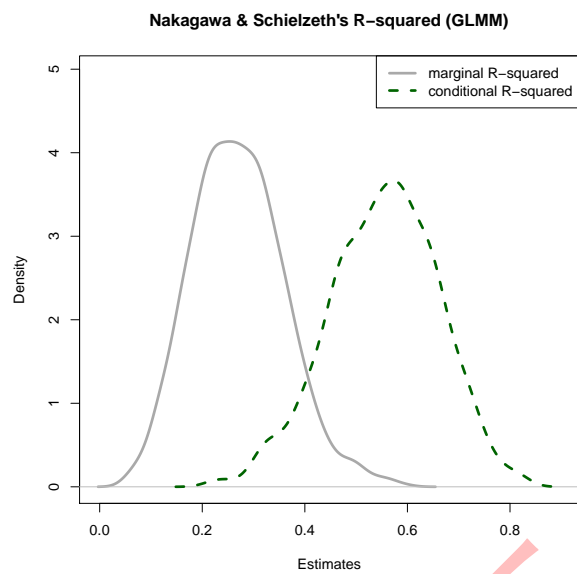


Figure 4: Distribution of Nakagawa & Schielzeth's  $R^2$  in the simulations described in Section 2.2.2

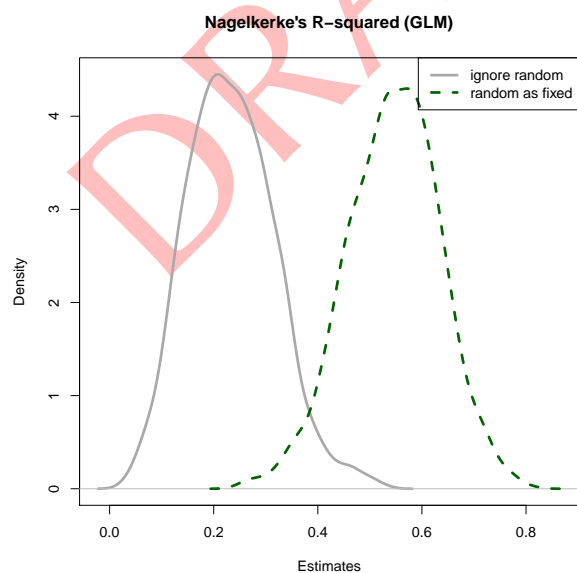


Figure 5: Nagelkerke's  $R^2$  in the simulations described in Section 2.2.2 for a GLM that ignores the grouping factor and a model that includes it as a fixed effect

### 2.2.4 More complex models

**Varying intercepts and slopes** 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). Instead of estimating a fixed coefficient, 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  $\beta^j$  follows 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 coefficient for the covariance between them is  $\rho$ . Figure 6 shows the bivariate density distributions for two (1) negatively correlated, (2) non-correlated, and (3) positively correlated normally distributed variables.

The number of variance parameters to be estimated thus obviously increases

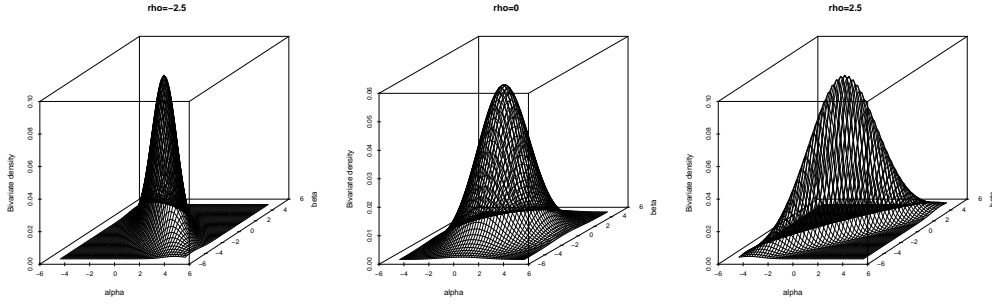


Figure 6: Bivariate normal density distribution with different correlation coefficients  $\rho$ ;  $\sigma_\alpha = \sigma_\beta = 3$ ;  $\mu_\alpha = \mu_\beta = 0$

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-parametrised* and needs to be simplified.

**Nested and crossed random effects** As it was explained in Section 2.1.2, nested random effects are adequate when grouping factors are nested within other grouping factors. Technically, while varying slopes can be understood as interactions between a fixed and a random effect, nested random intercepts can be understood interactions between two or more random effects. Crossed random effects are just several unrelated random effects. (8) shows the model specification, extending (2) with a varying intercept  $\alpha_s$ . This could be for example semantic classes which nest individual lemmas. It could also be another grouping factor for speaker, completely unrelated to the lemmas.

$$P(y^i = 1) = \text{logit}^{-1}(\alpha_s^{k[i]} + \alpha_l^{j[i]} + \beta_d \cdot x_d^i) \quad (8)$$

The difference is that in the nested case,  $k[i] = k[j]$ , i. e., the level of the nesting factor can be selected based on the nested factor as well as based on the single observation. As was mentioned in Section 2.1.2, the question is rather one of how the way the data are organised.

**Second-level predictors** In Section 2.1.3, situations were introduced where the random effects themselves can be partially predicted from fixed-effects. In this case, an additional linear model is specified for the random effect instead of the simple normal distribution predictor. We extend (2) by a predictor  $\gamma_f$  for the lemma frequency. The lemma frequencies themselves we denote by  $u_f$ , and we index them with  $j$ , just like the verb lemmas. This is reasonable because for each verb lemma, there is exactly one frequency. The first-level model specification remains the same, namely (9).

$$P(y^i = 1) = \text{logit}^{-1}(\alpha_i^{j[i]} + \beta_d \cdot x_d^i) \quad (9)$$

However, instead of (3), the varying intercept is predicted from (10).

$$\alpha_i^j \sim N(\gamma_0 + \gamma_f \cdot u_f^j, \sigma_l^2) \quad (10)$$

Instead of just the mean of the  $\alpha_j$  values, the model in (10) specifies a second-level intercept  $\gamma_0$  and a second-level fixed coefficient  $\gamma_f$ .

**Remarks on models for longitudinal studies** A longitudinal study is one wherein single subjects (usually speakers) are observed at different points in time, for example second-language learners after different years of learning a second language. First, the observations are obviously grouped by the individual speakers. It is thus recommended to include the speaker grouping factor,

typically as a random effect. Second, there might be time parameters such as years of learning.

We now assume we have a (fictional) sample from a learner corpus of German as a second language. In a logistic regression, we examine whether Swedish and English learners use the weak or the strong forms of attributive adjectives in NPs with a determiner. The crucial morpho-syntactic variable is whether the determiner has itself a strong ending or not, and we include an appropriate first-level term  $\beta_d \cdot x_d^i$  in the model. A random effect for individual learners is also included as  $\alpha_l^{j[i]}$ . Additionally, we add a term for the number of years the learners have learned German. However, since learners progress with different speed, a random slope is indicated, and the term becomes  $\beta_{y:l}^{j[i]} \cdot x_y^i$ . The first-level model is thus (11).

$$P(y^i = 1) = \text{logit}^{-1}(\alpha_l^{j[i]} + \beta_d \cdot x_d^i + \beta_{y:l}^{j[i]} \cdot x_y^i) \quad (11)$$

The main purpose of our study might be to find out whether Swedish and English learners differ with respect to the phenomenon at hand. Therefore, the first language is added as a second-level predictor with the term  $\gamma_f \cdot u_f^j$ , where  $u_f^j$  is 1 if the language of learner  $j$  is Swedish, and 0 if it is English. Since we have a random intercept and a random slope we need to distinguish between  $\gamma_f^\alpha \cdot u_f^j$  for the random intercept and  $\gamma_f^\beta \cdot u_f^j$  for the random slope. The second level model to go with (11) is therefore (12).

$$\begin{pmatrix} \alpha^j \\ \beta^j \end{pmatrix} \sim \left( \begin{pmatrix} \gamma_0^\alpha + \gamma_f^\alpha \cdot u_f^j \\ \gamma_0^\beta + \gamma_f^\beta \cdot u_f^j \end{pmatrix}, \begin{pmatrix} \sigma_\alpha^2 & \rho\sigma_\alpha\sigma_\beta \\ \rho\sigma_\alpha\sigma_\beta & \sigma_\beta^2 \end{pmatrix} \right) \quad (12)$$

Although this is a quite complex model already, there might be more problems. The performance of learners after  $n + 1$  years of learning is usually correlated

to a high degree with their performance after  $n$  years. This can lead to *auto-correlation* in the errors, violating basic modeling assumptions. To take care of this, model which allow for explicitly specified error structures must be used. Anyone wishing to use such models should consult further references (e. g., Fox 2016 and Zuur et al. 2009).

### 2.3 Specifying models using lme4 in R

This section and the next focus on lme4, the standard package to do multilevel modeling in R with maximum likelihood methods (Bates et al. 2015).

**Varying intercepts** The functions `lmer` and `glmer` extend the syntax of `lm` and `glm`. The varying intercept model in (2) is specified as follows in R (using informative variable names instead of Greek letters).

```
glmer(formula = construction ~ discourse + (1 | lemma),
      family = binominal(link=logit), data = my.data)
```

The pipe operator `x1 | x2` can be read as *x1 varies by x2*. The intercept is denoted by 1, and hence `(1 | lemma)` simply says that the intercept varies by lemma.

**Varying intercepts and slopes** The VIVS model in (6) is specified as follows (only the formula).

```
construction ~ discourse + (1 + discourse | lemma)
```

Before the pipe, the part of the model is repeated that should be modeled as varying by the grouping factor after the pipe. If a varying slope is specified, a

varying intercept is silently assumed. The last formula can therefore be abbreviated to the following equivalent one.

```
construction ~ discourse + (discourse | lemma)
```

In order to let *only* the slope vary, the intercept has to be removed explicitly from the random part of the formula.

```
construction ~ discourse + (discourse - 1 | lemma)
```

**Multiple random effects** When there is more than one random effect, several bracketed terms are added. The following is the recommended specification for models like (8), regardless of whether the effects are nested or crossed.

```
construction ~ discourse + (1 + | lemma) +
                    (1 | semantics)
```

Sometimes the following notation is used for nested random effects, where semantics nests lemma.

```
construction ~ discourse + (1 | semantics / lemma)
```

lme4 expands this to the following underlying syntax, which shows more clearly that nesting is handled as a kind of interaction.

```
construction ~ discourse + (1 | semantics) +
                    (1 | semantics : lemma)
```



There is a random intercept for semantics and one for each combination of semantics and lemma. While these notations are seemingly very explicit about the nesting structure, they are not necessary under normal circumstances. If the grouping factor lemma is nested within semantics (see Table 1 for a similar situation), lme4 automatically treats it as nested, and the results are exactly the same with all the three aforementioned notations.

However, the following specification is *not* equivalent and leads to problematic results.

```

construction ~ discourse + (1 | semantics) +
                (1 | lemma) +
                (1 | semantics : lemma)

```

This instructs lme4 to estimate the variance of lemma not just restricted to the permutations of the levels of lemma and semantics (i.e., semantics:lemma), but also outside of these specific permutations. In the nested case, there are no occurrences outside of these permutations, however, and the variance for lemma alone will be estimated close (but not exactly equal) to 0. To compensate for the spurious estimate for lemma, the variance estimate for semantics:lemma will be shifted unpredictably.

There is one situation where the explicit notation for nested factors is necessary. This is when the data are stored in a suboptimal way. Such a suboptimal version of Table 1 would look something like Table 4. Here, the speaker factor is encoded as the initial letter of the name only. Hence, Daryl and Dale (coming from two different regions) cannot be distinguished from each other, and Riley and Reed cannot, either. This leaves lme4 no way of recognising that the data structure is nested, and the user has to explicitly provide that information.

Exemplar	Speaker	Region
1	D	Tyneside
2	D	Tyneside
3	R	Tyneside
4	R	Tyneside
5	D	Greater London
6	D	Greater London
7	R	Greater London
8	R	Greater London

Table 4: Illustration of nested factors, organised suboptimally

It would, of course, be better *not* to organise data that way.

**Second-level predictors** (9) and (10) have the following lme4 syntax.

```
construction ~ discourse + frequency + (1 | lemma)
```

If the data is organised as shown in Table 3 – i. e., , with the second-level regressor not having any variance within the levels of the grouping factor –, lme4 will detect this and treat frequency as a second-level effect. However, second-level predictors for random slopes are more tricky to specify (see Gelman & Hill 2006: 280-282). Assuming that the effect of discourse status varies with the lemma, which itself comes with a second-level model including frequency as a regressor, the specification looks as follows.

```
construction ~ discourse + frequency +
              discourse : frequency +
              (1 + discourse | lemma)
```

A second-level regressor on a varying slope is thus an interaction between a first-level and a second-level fixed effect.

## 2.4 After fitting models with lme4

**Basics and varying intercepts** The output for GLMMs in lme4 can be understood straightforwardly after what was said in Sections 2. Here is a possible output of the summary function for a fit of model (2) and (3), repeated here as (13) and (14). Artificial data were used.

$$P(y^i = 1) = \text{logit}^{-1}(\alpha_l^{j[i]} + \beta_d \cdot x_d^i) \quad (13)$$

$$\alpha_l^j \sim N(\mu_l, \sigma_l^2) \quad (14)$$

```

Generalized linear mixed model fit by maximum likelihood
Family: binomial ( logit )
Formula: construction ~ discourse + (1 | lemma)
Data: observations
Random effects:
Groups Name          Variance Std.Dev.
lemma (Intercept) 1.29      1.136
Number of obs: 250, groups: lemma, 5

Fixed effects:
              Estimate Std. Error z value Pr(>|z|)
(Intercept)   0.7638     0.5513   1.385   0.166
discourse1    1.5064     0.3626   4.154 3.26e-05 ***

```

Some clutter as well as information which we do not interpret here (AIC, BIC, and information about the residuals) have been removed. In this output, the (Intercept) estimate (0.7638) is  $\hat{\mu}_l$ , and the Variance estimate for the lemma

random intercept (1.29) is  $\hat{\sigma}_l^2$ .<sup>7</sup> The estimate for `discourse1` (1.5064) corresponds to  $\hat{\beta}_d$ . Finally, we learn that there were five different lemmas and 250 observations in total.

To see whether the random intercept has a considerable influence, we should first look at the variance estimate. Here, it is larger than 1, which would be surprising if there were nothing going on in terms of between-lemma variation. It is possible to compute confidence intervals for the variance estimate using the `confint` function. Assuming the original model was stored in `alternation`, the following two alternatives work.

```
confint(alternation, parm="theta_", method = "profile")
confint(alternation, parm="theta_", method = "boot",
        nsim = 250)
```

The profile method uses LR tests and the bootstrap method uses a parametric bootstrap. For this model (where the variance estimate was 1.29 and the true value used to generate the data was 1.5), the profile method gives 0.5808 ... 2.6433 and the bootstrap with 250 simulations gives  $3.9665 \cdot 10^{-6}$  ... 1.8023 as the 95% confidence interval. Since the bootstrap (especially with smaller original sample sizes as in this case) typically tends to run into replications where the estimation of the variance fails and is thus returned as 0, the bootstrap interval is skewed to the left, while the profile confidence interval frames the true value symmetrically. The bootstrap is thus not always more robust or intrinsically better.

Although the authors of the `lme4` package advise against it, a significance test on the deviances of a simple GLM and a GLMM with an added single random

---

<sup>7</sup>The notation  $\hat{v}$  is used to denote an estimate of or a prediction for the variable  $v$ .

effect can be performed with the `anova` function.

```
alternation.0 <- glm(construction ~ discourse,
                    data = observations,
                    family = binomial(link=logit))
anova(alternation, alternation.0)
```

The GLMM object must be the first argument to `anova`. In this case, the output looks like this, indicating a significant effect, although the p-values should not be considered highly reliable.

```
Data: observations
Models:
alternation.0: construction ~ discourse
alternation:   construction ~ discourse + (1 | lemma)
              Df  logLik deviance  Chisq Df Pr(>Chisq)
alternation.0 2 -134.52  269.05
alternation   3 -119.12  238.24 30.801  1 2.859e-08 ***
```

If the nested (simpler) model still contains a random effect, the `update` function can be used to build the simpler model. Also, in addition to the `anova` command, there is a drop-in replacement for LR tests using bootstrap methods in the `pbkrtest` package (Halekoh & Højsgaard 2014).

```
alternation.1 <- glmer(construction ~ discourse +
                      (1 + | lemma) + (1 | semantics),
                      family = binomial(link=logit),
                      data = my.data)
alternation.2 <- update(object = alternation,
```

```

                                formula = construction ~
                                discourse +
                                (1 + | lemma))

require(pbkrtest)

PBmodcomp(alternation.1, alternation.2, nsim = 250)

```

The coefficient of determination (pseudo- $R^2$ ) according to Nakagawa & Schielzeth (2013) can be computed using the function `r.squaredGLMM` from the `MuMIn` package (Bartoń 2016).

```

require(MuMIn)
r.squaredGLMM(alternation)

```

In this case, it gives us  $R^2_m = 0.1101$  and  $R^2_c = 0.3608$ , so there is a considerable difference between the marginal  $R^2$  (without random effects) and conditional  $R^2$  (with random effects).

To inspect the conditional modes, the `ranef` function can be used, and it can also output standard errors for them.

```

ranef(alternation, drop = T, condVar = T)

```

**Varying intercepts and slopes** In the summary output for a VIVS model such as (6), the columns `Variance` and `Corr` can be regarded as specifying the lower triangle of the variance-covariance matrix.<sup>8</sup>

---

<sup>8</sup>Alternatively, the `VarCorr` function could be used to extract this information.

```
Random effects:
```

Groups Name	Variance	Std.Dev.	Corr
lemma (Intercept)	0.6103	0.7812	
discourse	0.8944	0.9457	-0.39

```
Number of obs: 2500, groups: lemma, 50
```

```
Fixed effects:
```

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	0.6332	0.1226	5.163	2.43e-07 ***
discourse	-1.0428	0.1492	-6.990	2.75e-12 ***

This output tells us that the estimated variance in the intercepts is  $\hat{\sigma}_\alpha^2 = 0.6103$ , the estimated variance in the slopes is  $\hat{\sigma}_\beta^2 = 0.8944$ , and the covariance coefficient estimate is  $\hat{\rho} = -0.39$  (a healthy value). The means are estimated as  $\hat{\mu}_\alpha = 0.6332$  and  $\hat{\mu}_\beta = -1.0428$ . Compare this to Section 2.2.4, especially (7). It is possible to reconstruct group-wise models from this output and a lookup of the group-specific predictions using the `ranef` function. For the first group, for example, the following can be done.

```
ranef(alternation)$lemma[1,]
```

The output is as follows.

```
(Intercept) discourse
0001      0.4351156 -1.227842
```

This means that for the first lemma, actual predictions for the outcome of the alternation can be made using (15), where values are rounded to two decimal digits. Compare this to (6).

$$Pr(y^i = 1) = \text{logit}^{-1}([0.63 + 0.44] + [-1.04 - 1.23] \cdot x_d^i) \quad (15)$$

### 3 Representative studies

#### **Wolk et al. (2013)**

**Research questions** The authors aim to achieve two things. First, they want to compare changes in two word order-related alternations in the history of English between 1650 and 1999: the dative alternation and the genitive alternation. They look for influencing features shared in both cases as well as construction-specific features. Second, they aim to show that historical data fits well into a probabilistic, cognitively oriented view of language.

**Data** The authors use the ARCHER corpus, which contains texts from various registers from 1650 to 1999. For both constructions, carefully designed sampling protocols were used (see their Section 4). For the annotation of the data, both available corpus meta data were used (text ID, register, time in fractions of centuries, centered at 1800) as well as a large number of manually coded variables (constituent length, animacy, definiteness, etc.). Furthermore, the possessor head lemma (genitive alternation) and the verb lemma (dative alternation) were coded.

**Method** Two mixed effects logistic regression models are estimated. For the genitive alternation, the text ID and the possessor head lemma are used as crossed random effects. The authors state on p. 399 that they collapsed



all head noun lemmas with less than four occurrences into one category because otherwise “difficulties” would arise. However, it is the advantage of random effects modeling that it can deal with a situation where categories have low numbers of observations (see *shrinkage*, Section 2.2.2). For the dative alternation, the model includes the text ID, the register (which nests the text ID) as well as the lemma of the theme argument and the verb.

**Results** It is found that many factors have a shared importance in both alternations, e. g., definiteness, animacy, construction length. It is also argued that the observed tendencies – such as *short-before-long* and *animate referents first* – are in line with synchronic corpus based and experimental findings about general cognitive principles underlying the framework of probabilistic grammar. These principles remain in effect, but the strength of their influence changes over time.

### Gries (2015)

**Research questions** The paper is programmatic in nature. The author re-analyses data from a previously published study on verb particle placement in English. He uses a GLMM instead of a fixed-effects logistic regression to show that including random effects in order to account for variation related to mode, register, and subregister increases the quality and predictive power of the model. He also argues that not doing so, corpus linguists risk violating fundamental assumptions about the independence of the error terms in models.

**Data** The data are 2,321 instances of particle verbs showing either verb–direct object–particle or verb–particle–direct object order, taken from the ICE-GB. The grouping factors derived from the structure of the corpus are mode (only two levels), register (five levels), and subregister (13 levels). They are nested: mode nests register, which nests subregister. Additionally, verb and particle lemma grouping factors are annotated. Finally, two fixed effects candidates are annotated (the type of the head of direct object and the logarithmised length of the direct object in words).

**Method** The author uses the model selection protocol described in Zuur et al. (2009) to first find the optimal random effects structure using ANOVAs and AIC comparisons as well as analyses of the estimated variance for single random effects. He then goes on to find the optimal fixed effects structure. Additionally, he compares the pseudo- $R^2$  measure of the resulting mixed models.

**Results** It is found that the verb and particle lemma play and the subregister play a significant role. Notably, the variance estimate for mode is close to 0 from the beginning of the model selection procedure. This is not surprising, as two levels are not nearly enough in order for the variance to be reliably estimated, and it could maybe be used as a second-level predictor instead (see Section 2.2.2). The  $R^2$  values of the final model are very high, with a considerable difference between marginal  $R_m^2 = 0.57$  and conditional  $R_c^2 = 0.748$ , which indicates that the random effects do in fact improve the model fit. It is also shown that the classification accuracy is considerably improved over that of a GLM without random effects, but dif-

ferently for different lexical groups and subregisters. The paper thus shows that it is not appropriate to ignore lexical grouping factors and grouping factors derived from the corpus structure, especially as both are easy to annotate automatically.

## 4 Further reading

Chapters 1–15 and Chapters 20–24 of Gelman & Hill (2006) are a highly recommended read, especially for R and lme4 users. Similarly, Zuur et al. (2009) has a reputation among R users of mixed effects models in many fields. The companion to lme4, Bates (2010) and the overview in Bates et al. (2015) are obligatory reads for users of lme4.

## References

- Bartoń, Kamil. 2016. *MuMIn: Multi-Model Inference*. R package version 1.15.6. <https://CRAN.R-project.org/package=MUMIn>.
- Bates, Douglas M. 2010. Lme4: mixed-effects modeling with R. <http://lme4.r-forge.r-project.org/IMMwR/lrgprt.pdf>.
- Bates, Douglas, Martin Mächler, Ben Bolker & Steve Walker. 2015. Fitting linear mixed-effects models using lme4. *Journal of Statistical Software* 67(1). 1–48.
- Fox, John. 2016. *Applied regression analysis & generalized linear models*. 3rd edn. London: Sage Publications.
- Gelman, Andrew & Jennifer Hill. 2006. *Data analysis using regression and multilevel/hierarchical models*. Cambridge: Cambridge University Press.
- Gries, Stefan Th. 2015. The most underused statistical method in corpus linguistics: multi-level (and mixed-effects) models. *Corpora* 10(1). 95–126.

- Halekoh, Ulrich & Søren Højsgaard. 2014. A Kenward-Roger approximation and parametric bootstrap methods for tests in linear mixed models – the R package pbrtest. *Journal of Statistical Software* 59(9). 1–30.
- 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.
- Schielzeth, Holger & Wolfgang Forstmeier. 2009. Conclusions beyond support: overconfident estimates in mixed models. *Behavioral Ecology* 20(2). 416–420.
- Wolk, Christoph, Joan Bresnan, Anette Rosenbach & Benedikt Szmrecsanyi. 2013. Dative and genitive variability in Late Modern English: Exploring cross-constructural variation and change. *Diachronica* 30(3). 382–419.
- Zuur, Alain F., Elena N. Ieno, Neil Walker, Anatoly A. Saveliev & Graham M. Smith. 2009. *Mixed effects models and extensions in ecology with R*. Berlin etc.: Springer.