# Generalizing the linear mixed model to factorial designs

Tanner Sorensen

University of Potsdam, Potsdam, Germany

Shravan Vasishth

University of Potsdam, Potsdam, Germany, and

School of Mathematics and Statistics, University of Sheffield, Sheffield, UK

July 22, 2015

The Gibson and Wu (2013) data-set has a two-condition design. This section presents a varying intercepts, varying slopes model for a $2 \times 2$ factorial design. Because of the more general matrix formulation we use here, the Stan code can be deployed with minimal changes for much more complex designs, including correlational studies.

Our example is the $2 \times 2$ repeated measures factorial design of Husain, Vasishth, and Srinivasan (2014, Experiment 1), also a self-paced reading study on relative clauses. The dependent variable was the reading time `rt` of the relative clause verb. The factors were relative clause type, which we code with the predictor `so` (`so` $= +1$ for object relatives and `so` $= -1$ for subject relatives) and distance between the head noun and the relative clause verb, which we code with the predictor `dist` (`dist` $= +1$ for far and `dist` $= -1$ for near). Their interaction is the product of the `dist` and `so` contrast vectors, and labeled as the predictor `int`. The 60 subjects were speakers of Hindi, an Indo-Aryan language spoken primarily in India. The 24 items were presented in a standard, fully balanced Latin square design. This resulted in a total of 1440 data points ($60 \times 24 = 1440$). The first few lines from the data frame are shown below.

The theoretical interest is in determining whether relative clause type and distance

| row | subj | item | so | dist | rt |
|---:|---:|---:|:---:|:---:|---:|
| 1 | 1 | 14 | s | n | 1561 |
| 2 | 1 | 16 | o | n | 959 |
| 3 | 1 | 15 | o | f | 582 |
| 4 | 1 | 18 | s | n | 294 |
| 5 | 1 | 4 | o | n | 438 |
| 6 | 1 | 17 | s | f | 286 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |
| 1440 | 9 | 13 | s | f | 516 |

Table 1

*The first six rows, and the last row, of the data-set of Husain et al. (2014, Experiment 1), as they appear in the data frame.*

influence reading time, and whether there is an interaction between these two factors. We use Stan to determine the posterior probability distribution of the fixed effect $\beta_1$ for relative clause type, the fixed effect $\beta_2$ for distance, and their interaction $\beta_3$.

We fit a varying intercepts, varying slopes model to this data-set. The grand mean $\beta_0$ of log `rt` is adjusted by subject and by item through the varying intercepts $u_0$ and $w_0$, which are unique values for each subject and item respectively. Likewise, the three fixed effects $\beta_1$, $\beta_2$, and $\beta_3$ which are associated with the predictors `so`, `dist`, and `int`, respectively, are adjusted by the by-subject varying slopes $u_1$, $u_2$, and $u_3$ and by-item varying slopes $w_1$, $w_2$, and $w_3$.

It is more convenient to represent this model in matrix form. We build up the model specification by first noting that, for each subject, the by-subject varying intercept $u_0$ and slopes $u_1$, $u_2$, and $u_3$ have a multivariate normal prior distribution with mean zero and covariance matrix $\Sigma_u$. Similarly, for each item, the by-item varying intercept $w_0$ and slopes $w_1$, $w_2$, and $w_3$ have a multivariate normal prior distribution with mean zero and covariance matrix $\Sigma_w$. We can write this as follows:

$$
\begin{pmatrix} u_0 \\ u_1 \\ u_2 \\ u_3 \end{pmatrix} \sim \mathrm{N} \left( \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \Sigma_u \right) \quad \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ w_3 \end{pmatrix} \sim \mathrm{N} \left( \begin{pmatrix} 0 \\ 0 \\ 0 \\ 0 \end{pmatrix}, \Sigma_w \right) \tag{1}
$$

```
1   rDat<-read.table("HusainEtAlexpt1data.txt",header=TRUE)
2   rDat$subj <- with(rDat,factor(subj))
3   rDat$item <- with(rDat,factor(item))
4
5   X <- unname(model.matrix(~1+so+dist+int, rDat))
6
7   stanDat <- within(list(),
8   {
9     N<-nrow(X)
10    P <- n_u <- n_w <- ncol(X)
11    X <- X
12    Z_u <- X
13    Z_w <- X
14    J <- length(levels(rDat$subj))
15    K <- length(levels(rDat$item))
16    rt <- rDat$rt
17    subj <- as.integer(rDat$subj)
18    item <- as.integer(rDat$item)
19  }
20  )
21  factorialFit <- stan(file="factorialModel.stan",
22                       data=stanDat,
23                       iter=2000, chains=4)
```

Listing 1: Preparation of data for analyzing the Husain et al. data-set, and running the model.

The error $\varepsilon$ is assumed to have a normal distribution with mean zero and standard deviation $\sigma_e$.

We proceed to implement the model in Stan. First we read in the data-set (see Listing 1). Instead of passing the predictors so, dist, and their interaction int to stan as vectors, as we did with so earlier, we make so, dist, and int into a design matrix X using the function model.matrix available in R.[1] The first column of the design matrix X consists of all ones. The second column is the predictor so which codes the factor for relative clause type. The third column the predictor dist which codes the factor for distance. The fourth column is the predictor int which codes the interaction between relative clause type and distance. The model matrix thus consists of a fully factorial $2 \times 2$ design, with blocks of this design repeated for each subject. For the full data-set, we could write it very compactly in matrix form as follows:

---

[1]Here, we would like to acknowledge the contribution of Douglas Bates in specifying the model in this general matrix form.

```
1   data {
2     int<lower=0> N;                  //no trials
3     int<lower=1> P;                  //no fixefs
4     int<lower=0> J;                  //no subjects
5     int<lower=1> n_u;                //no subj ranefs
6     int<lower=0> K;                  //no items
7     int<lower=1> n_w;                //no item ranefs
8     int<lower=1,upper=J> subj[N];    //subject indicator
9     int<lower=1,upper=K> item[N];    //item indicator
10    row_vector[P] X[N];              //fixef design matrix
11    row_vector[n_u] Z_u[N];          //subj ranef design matrix
12    row_vector[n_w] Z_w[N];          //item ranef design matrix
13    vector[N] rt;                    //reading time
14  }
15  parameters {
16    vector[P] beta;                  //fixef coefs
17    cholesky_factor_corr[n_u] L_u;   //cholesky factor of subj ranef corr matrix
18    cholesky_factor_corr[n_w] L_w;   //cholesky factor of item ranef corr matrix
19    vector<lower=0>[n_u] sigma_u;    //subj ranef std
20    vector<lower=0>[n_w] sigma_w;    //item ranef std
21    real<lower=0> sigma_e;           //residual std
22    vector[n_u] z_u[J];              //subj ranef
23    vector[n_w] z_w[K];              //item ranef
24  }
25  transformed parameters {
26    vector[n_u] u[J];                //subj ranefs
27    vector[n_w] w[K];                //item ranefs
28    {
29      matrix[n_u,n_u] Sigma_u;       //subj ranef cov matrix
30      matrix[n_w,n_w] Sigma_w;       //item ranef cov matrix
31      Sigma_u <- diag_pre_multiply(sigma_u,L_u);
32      Sigma_w <- diag_pre_multiply(sigma_w,L_w);
33      for(j in 1:J)
34        u[j] <- Sigma_u * z_u[j];
35      for(k in 1:K)
36        w[k] <- Sigma_w * z_w[k];
37    }
38  }
39  model {
40    //priors
41    L_u ~ lkj_corr_cholesky(2.0);
42    L_w ~ lkj_corr_cholesky(2.0);
43    for (j in 1:J)
44      z_u[j] ~ normal(0,1);
45    for (k in 1:K)
46      z_w[k] ~ normal(0,1);
47    //likelihood
48    for (i in 1:N)
49      rt[i] ~ lognormal(X[i] * beta +
50                        Z_u[i] * u[subj[i]] +
51                        Z_w[i] * w[item[i]],
52                        sigma_e);
53  }
```

Listing 2: Stan code for Husain et al data.

$$\log(\mathbf{rt}) = \mathbf{X}\beta + \mathbf{Z}_u \mathbf{u} + \mathbf{Z}_w \mathbf{w} + \varepsilon \tag{2}$$

Here, $\mathbf{X}$ is the $N \times P$ model matrix (with $N = 1440$, since we have 1440 data points; and $P = 4$ since we have the intercept plus three other fixed effects), $\beta$ is a $P \times 1$ vector of fixed effects parameters, $\mathbf{Z}_u$ and $\mathbf{Z}_w$ are the subject and item model matrices ($N \times P$), and $u$ and $w$ are the by-subject and by-item adjustments to the fixed effects estimates; these are identical to the design matrix $\mathbf{X}$ in the model with varying intercepts and varying slopes included. For more examples of similar model specifications in Stan, see the R package `RePsychLing` on github (https://github.com/dmbates/RePsychLing).

Having defined the model, we proceed to assemble the list `stanDat` of data, relying on the above matrix formulation; please refer to Listing 1. The number `N` of observations, the number `J` of subjects and `K` of items, the reading times `rt`, and the subject and item indicator variables `subj` and `item` are familiar from the previous models presented. The integer `P` is the number of fixed effects (four including the intercept). Model 2 includes a varying intercept $u_0$ and varying slopes $u_1$, $u_2$, $u_3$ for each subject, and so the number `n_u` of by-subject random effects equals `P`. Likewise, Model 2 includes a varying intercept $w_0$ and varying slopes $w_1$, $w_2$, $w_3$ for each item, and so the number `n_w` of by-item random effects also equals `P`. The data block contains the corresponding variables. We declare the fixed effects design matrix `X` as an array of `N` row vectors whose components are the predictors associated with the `N` reading times. Likewise for the subject and item random effects design matrices `Z_u` and `Z_w`, which correspond to $\mathbf{Z}_u$ and $\mathbf{Z}_w$ respectively in Model 2. The vector `beta` contains the fixed effects $\beta_0$, $\beta_1$, $\beta_2$, and $\beta_3$. The matrices `L_u`, `L_w` and the arrays `z_u`, `z_w` of vectors (not to be confused with the design matrices `Z_u` and `Z_w`) will generate the varying intercepts and slopes $u_0, \ldots, u_3$ and $w_0, \ldots, w_3$. The vector `sigma_u` contains the standard deviations of the by-subject varying intercepts and slopes $u_0, \ldots, u_3$, and the vector `sigma_w` contains the standard deviations of the by-item varying intercepts and slopes $w_0, \ldots, w_3$. The variable `sigma_e` is the standard deviation $\sigma_e$ of the error $\varepsilon$. The

transformed parameters block generates the by-subject intercepts and slopes $u_0, \ldots, u_3$ and the by-item intercepts and slopes $w_0, \ldots, w_3$.

We place lkj priors on the random effects correlation matrices through the `lkj_corr_cholesky(2.0)` priors on their Cholesky factors `L_u` and `L_w`. We implicitly place uniform priors on the fixed effects $\beta_0, \ldots, \beta_3$, the random effects standard deviations $\sigma_{u0}, \ldots, \sigma_{u3}$, and $\sigma_{w0}, \ldots, \sigma_{w3}$ and the error standard deviation $\sigma_e$ by omitting any prior specifications for them in the model block. We specify the likelihood with the probability statement that `rt[i]` is distributed log-normally with mean `X[i] * beta + Z_u[i] * u[subj[i]] + Z_w[i] * w[item[i]]` and standard deviation `sigma_e`. The next step towards model-fitting is to pass the list `stanDat` to `stan`, which compiles a C++ program to sample from the posterior distribution of the model parameters.

Figure 1 plots histograms of the marginal posterior distribution of the fixed effects. The HPD interval of the fixed effect $\hat{\beta}_1$ for relative clause type is entirely below zero. This is evidence that object relatives are read faster than subject relatives. The HPD interval of the fixed effect $\hat{\beta}_2$ for distance is also entirely below zero. This is evidence of a slowdown when the verb (where reading time was measured) is closer to the head noun of the relative clause. The HPD of the interaction $\hat{\beta}_3$ between relative clause type and distance is greater than zero, which is evidence for a greater slowdown on subject relatives when the distance between the verb and head noun is short.

A major advantage of the above matrix formulation is that we do not need to write a new Stan model for a future repeated measures factorial design. All we have to do now is define the design matrix $X$ appropriately, and include it (along with appropriately defined $Z_u$ and $Z_w$ for the subjects and items random effects) as part of the data specification that is passed to Stan.
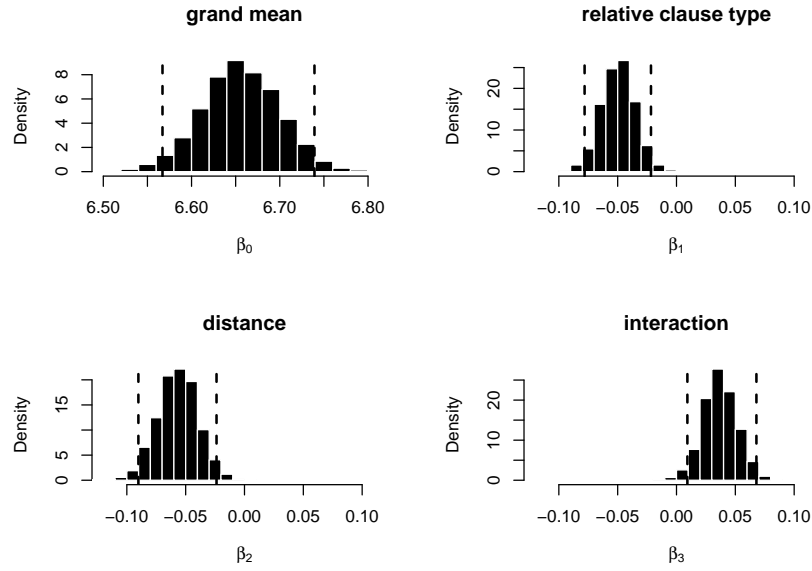
*Figure 1*. Marginal posterior distribution and HPD intervals of the fixed effects grand mean $\beta_0$, slope $\beta_1$ for relative clause type, slope $\beta_2$ for distance, and interaction $\beta_3$. All fixed effects are on the log-scale.

References

Gibson, E., & Wu, H.-H. I. (2013). Processing chinese relative clauses in context. *Language and Cognitive Processes*, *28*(1-2), 125–155.

Husain, S., Vasishth, S., & Srinivasan, N. (2014). Strong expectations cancel locality effects: Evidence from Hindi. *PLoS ONE*, *9*(7), 1–14.