

Detecting Treatment-Subgroup Interactions in Clustered Data with Generalized Linear Mixed-Effects Model Trees

M. Fokkema¹, N. Smits², A. Zeileis³, T. Hothorn⁴, H. Kelderman⁵

¹Universiteit Leiden, ²Universiteit van Amsterdam, ³Universität Innsbruck, ⁴Universität Zürich, ⁵Universiteit Leiden and Vrije Universiteit, Amsterdam

Abstract

Identification of subgroups of patients for which treatment A is more effective than treatment B, and vice versa, is of key importance to the development of personalized medicine. Tree-based algorithms are helpful tools for the detection of such interactions, but none of the available tree-based algorithms allow for taking into account a clustered or nested structure. Therefore, we propose the generalized linear mixed-effects model trees (GLMM tree) algorithm, which allows for the detection of treatment-subgroup interactions, while accounting for the clustered structure of a dataset. The algorithm uses model-based recursive partitioning (MOB) to detect treatment-subgroup interactions, and a GLMM for the estimation of random-effects parameters. In a simulation study, we evaluate the performance of GLMM trees and compare it with that of trees without random effects, and GLMMs with pre-specified interactions. GLMM tree was found to accurately recover treatment-subgroup interactions in 90% of the simulated datasets, whereas trees without random effects only did so in 61% of datasets. GLMM tree also outperformed trees without random effects in terms of predictive accuracy and Type-I error rates. Furthermore, compared to GLMMs with pre-specified interaction effects, GLMM tree showed better predictive accuracy. We illustrate the application of GLMM trees on an individual patient-level data meta-analysis on treatments for depression. We conclude that GLMM tree is a promising exploratory tool for the detection of treatment-subgroup interactions in clustered datasets.

The authors would like to thank Prof. Pim Cuijpers, Prof. Jeanne Miranda, Dr. Boadie Dunlop, Prof. Rob DeRubeis, Prof. Zindel Segal, Dr. Sona Dimidjian, Prof. Steve Hollon and Erica Weitz for granting access to the dataset for the application. The work for this paper was partially done while MF, AZ and TH were visiting the Institute for Mathematical Sciences, National University of Singapore in 2014. The visit was supported by the Institute.

Keywords: model-based recursive partitioning, treatment-subgroup interactions, random effects, generalized linear mixed-effects model, classification and regression trees

Introduction

In research on the efficacy of treatments for somatic and psychological disorders, the one-size-fits-all paradigm is slowly losing ground, and stratified (or personalized) medicine is becoming increasingly important. Stratified medicine presents the challenge of finding which patients respond best to which treatments. This can be referred to as the detection of treatment-subgroup interactions (e.g., Doove, Dusseldorp, Van Deun, & Van Mechelen, 2014). Often, treatment-subgroup interactions are studied using linear models, such as factorial analysis of variance techniques, in which potential moderators have to be specified a-priori, have to be checked one at a time, and continuous moderator variables have to be discretized. This may hamper identification of which treatment works best for whom, especially when there are no a-priori hypotheses about treatment-subgroup interactions. As noted by Kraemer, Frank, and Kupfer (2006), there is a need for methods that generate instead of test such hypotheses.

Tree-based methods are such hypothesis-generating methods, as they can automatically detect subgroups which differ in the expected outcomes for one or more treatments. Due to their flexibility, tree-based methods are particularly useful for exploratory purposes, as they can handle many potential predictor variables at once and can automatically detect (higher order) interactions between predictor variables (Strobl, Malley, & Tutz, 2009). As such, tree-based methods are preeminently suited to the detection of treatment-subgroup interactions. Several tree-based algorithms for the detection of treatment-subgroup interactions have already been developed (Dusseldorp, Doove, & Van Mechelen, 2015; Dusseldorp & Meulman, 2004; Su, Tsai, Wang, Nickerson, & Li, 2009; Foster, Taylor, & Ruberg, 2011; Lipkovich, Dmitrienko, Denne, & Enas, 2011; Zeileis, Hothorn, & Hornik, 2008; Seibold, Zeileis, & Hothorn, 2016; see Doove et al., 2014 for an overview). Also, Zhang, Tsiatis, Laber, and Davidian (2012); Zhang, Tsiatis, Davidian, Zhang, and Laber (2012) have developed a flexible classification-based approach, allowing users to select from a range of statistical methods, including trees.

In many instances, researchers may want to detect treatment-subgroup interactions in a generalized linear mixed-effects (GLMM) type model. For example, in individual-level patient data meta-analysis, where datasets of multiple clinical trials on the same treatments are pooled (e.g., Koopman, Van der Heijden, Glasziou, Grobbee, & Rovers, 2007). In such analyses, the nested or clustered structure of the dataset should be taken into account by including study-specific random effects in the model, prompting the need for

1 a mixed-effects model (e.g., Cooper & Patall, 2009; Higgins, Whitehead, Turner, Omar, &
 2 Thompson, 2001). In linear models, ignoring the clustered structure may lead, for example,
 3 to biased inference due to underestimated standard errors in linear models (e.g., Bryk &
 4 Raudenbush, 1992; Van den Noortgate, Opdenakker, & Onghena, 2005). For tree-based
 5 methods, ignoring the clustered structure has been found to result in the detection of
 6 spurious subgroups and inaccurate predictor variable selection (e.g., Sela & Simonoff, 2012;
 7 Martin, 2015). However, none of the purely tree-based methods for treatment-subgroup
 8 interaction detection allow for taking into account the clustered structure of a dataset.
 9 Therefore, in the current paper, we present a tree-based algorithm which can be used
 10 for the detection of (treatment-subgroup) interactions and non-linearities in GLMM type
 11 models: generalized linear mixed-effects model tree, or GLMM tree.

12 The GLMM tree algorithm builds on model-based recursive partitioning (MOB,
 13 Zeileis et al., 2008), which offers a flexible framework for subgroup detection. For ex-
 14 ample, GLM-based MOB has been applied to detect treatment-subgroup interactions for
 15 the treatment of depression (Driessen et al., 2016) and amyotrophic lateral sclerosis (Seibold
 16 et al., 2016). In contrast to other purely tree-based methods (e.g., Zeileis et al., 2008; Su et
 17 al., 2009; Dusseldorp & Van Mechelen, 2014)), GLMM tree allows for taking into account
 18 the clustered structure of datasets. In contrast to previously suggested regression trees with
 19 random effects, GLMM tree allows for treatment effect estimation, with continuous as well
 20 as non-continuous response variables (e.g., Hajjem, Bellavance, & Larocque, 2011; Sela &
 21 Simonoff, 2012).

22 In what follows, we first introduce an artificial motivating dataset, which will be
 23 used to illustrate and explain treatment-subgroup interaction detection with GLM-based
 24 MOB and GLMM tree. In the Simulation-based Evaluation, we evaluate the performance
 25 of GLMM tree in simulated datasets, and compare it with that of GLM-based MOB and
 26 GLMMs with pre-specified interaction effects. The comparison with GLM-based MOB al-
 27 lows for assessing the extent to which including random-effects estimation improves the
 28 performance of MOB trees. The comparison with GLMMs with pre-specified interactions
 29 will allow for assessing the extent to which using a tree-based method will improve detection
 30 of treatment-subgroup interactions. In the Application, we apply the algorithm to an exist-
 31 ing dataset of a patient-level meta-analysis on the effects of psycho- and pharmacotherapy
 32 for depression. Finally, in the Discussion we summarize the results and point out some
 33 directions for future research.

34 *Artificial motivating dataset*

35 We created a dataset representing observations on 150 participants in a randomized
 36 clinical trial. Every participant was randomly assigned to Treatment 1 or Treatment 2, and
 37 has a value for the response variable, with which the effect of treatment is assessed: the pos-

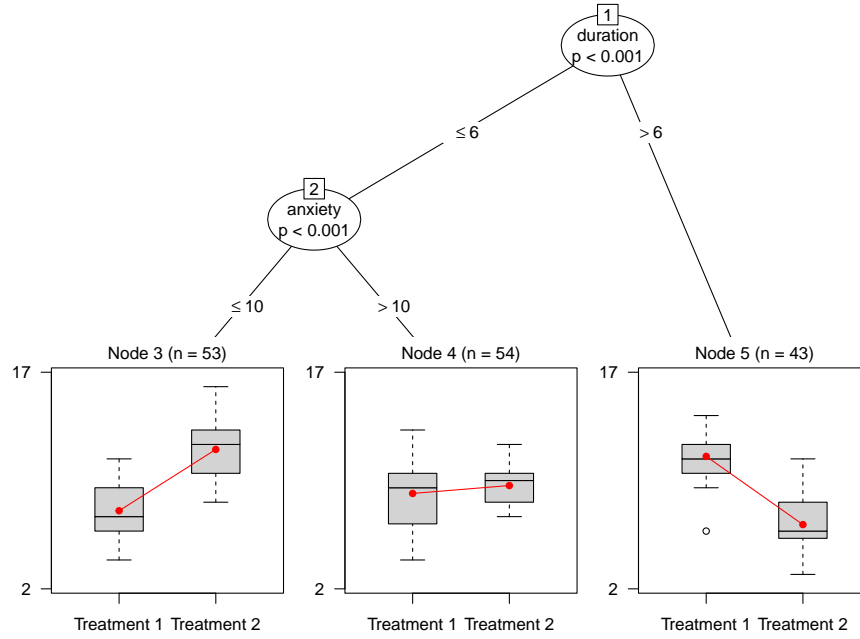


Figure 1. GLMM tree of the motivating example dataset. Three covariates (anxiety questionnaire score, duration of depressive symptoms at baseline in months and age) were used as potential splitting variables, and the clustering structure was taken into account by estimating random intercepts.

1 ttreatment total score on a depression inventory. For all participants, three covariate values
 2 are available: duration of depressive symptoms prior to treatment in months (duration,
 3 range 0–15); age in years (age, range 18–75); anxiety inventory total score (anxiety, range
 4 3–18).

5 The simulated dataset has 3 subgroups with differential treatment effectiveness. The
 6 first subgroup consists of observations with duration ≤ 6 and anxiety ≤ 10 . In this subgroup,
 7 Treatment 1 is more effective than Treatment 2: the mean of the response variable is 7 for
 8 Treatment 1, and 11 for Treatment 2. The second subgroup consists of observations with
 9 duration ≤ 6 and anxiety > 10 . In this subgroup, both therapies are equally effective: the
 10 mean value of the response variable is 9 for Treatment 1, and 9 for Treatment 2. The third
 11 subgroup consists of observations with duration > 6 . In this subgroup, Treatment 2 is more
 12 effective than Treatment 1: the mean value of the response variable is 12 for Treatment 1,
 13 and 7 for Treatment 2.

14 Participants were part of one of ten clusters, each with a different value for the random
 15 intercept. Data were generated such that covariates and cluster-specific intercepts were
 16 uncorrelated. 43% of variance in post-treatment depression scores was due to treatment-

subgroup interactions, while 8% of variance was due to cluster-specific variation.

See Figure 1 for the GLMM tree fitted to the simulated data set which recovers the true underlying structure.

Model-based recursive partitioning

The rationale behind MOB is that a global parametric model may not describe the data well, and when additional covariates are available it may be possible to partition the dataset with respect to these covariates, and find a better model in each cell of the partition (Zeileis et al., 2008). This is reminiscent of the classification and regression tree (CART) algorithm of Breiman, Friedman, Olshen, and Stone (1984), which splits the dataset into subsets, for which the distributions of the outcome variable are most different. However, CART trees detect differences in the mean (or median) value across terminal nodes, whereas MOB trees detect differences in parameters of more complex models across terminal nodes.

For example, let us take a global GLM to estimate the overall treatment effect in the motivating dataset. The expectation μ_i of outcome y_i given the treatment regressors x_i is modeled through a linear predictor and suitable link function:

$$E[y_i|x_i] = \mu_i, \quad (1)$$

$$g(\mu_i) = x_i^\top \beta, \quad (2)$$

where $x_i^\top \beta$ is the linear predictor for observation i and g is the link function. β is a vector of fixed-effects regression coefficients, the first element representing the intercept, corresponding to the mean value of the linear predictor in the first treatment group, and the second element representing the slope, which is the mean difference in the linear predictor between the first and second treatment groups. Thus, for simplicity we assume x_i and β to have length 2 in the current paper. Also, for the continuous response variable in the motivating data set, we employ a Gaussian distribution with identity link and denote the error by $\epsilon_i = y_i - \mu_i$ with variance σ_ϵ^2 . However, the model can easily accommodate additional treatment conditions and covariates, and binary or count outcome variables.

The result of fitting a global GLM to the motivating dataset is depicted in Figure 2; the boxplots show the distribution of the posttreatment depression scores in both treatment groups. The global model does not describe the data well: there is substantial residual variance and the slope of the regression line is nearly zero. This does not necessarily mean that posttreatment depression score and treatment type are unrelated, as the effect of treatment may be moderated by variables not yet included in the model.

The MOB algorithm can be used to detect such moderation, by testing for parameter stability over a set of auxiliary covariates, also known as *partitioning variables*. When the partitioning is based on a GLM, instabilities are differences in $\hat{\beta}$ across partitions of

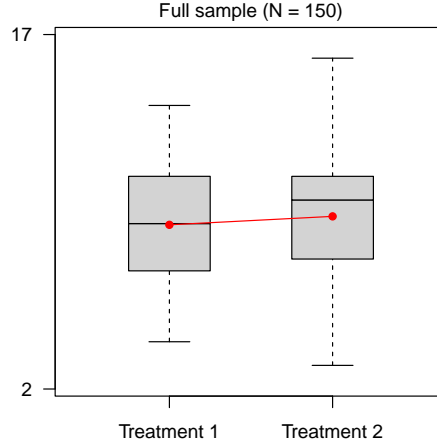


Figure 2. Example of a global GLM for treatment outcomes, based on the artificial motivating dataset ($N = 150$). The dot for Treatment 1 represents the first, and the slope of the regression line represents the second element of β .

the dataset, which are defined by one or more auxiliary covariates not included in the linear predictor. To find these partitions, the MOB algorithm cycles iteratively through the following steps (Zeileis et al., 2008): (1) fit the parametric model to the dataset, (2) test for parameter instability over a set of partitioning variables, (3) if there is some overall parameter instability, split the dataset with respect to the variable associated with the highest instability, (4) repeat the procedure in each of the resulting subgroups.

More specifically, in step (2), to test for parameter instability, the so-called *scores* are computed, using the score function. When the model is correctly specified, the expected value of the score for each observation equals zero. Therefore, under the null hypothesis of parameter stability, the scores do not systematically deviate from the expected value of zero, when observations are ordered by the values of a potential partitioning variable U_k (c.f., Merkle & Zeileis, 2013). To statistically test whether the scores systematically deviate from zero with respect to variable U_k , the class of generalized M-fluctuation tests is used (Zeileis, 2005; Zeileis & Hornik, 2007).

If the null hypothesis of parameter stability in step (2) can be rejected, that is, if at least one of the partitioning variables U_k yields a p -value for the M-fluctuation test below the pre-specified significance level α , the dataset is partitioned into two subsets in step (3). This partition is created using U_{k^*} , the variable with the minimal p -value in step (2). The split point for U_{k^*} is selected, by taking the value that minimizes the sum of the values of

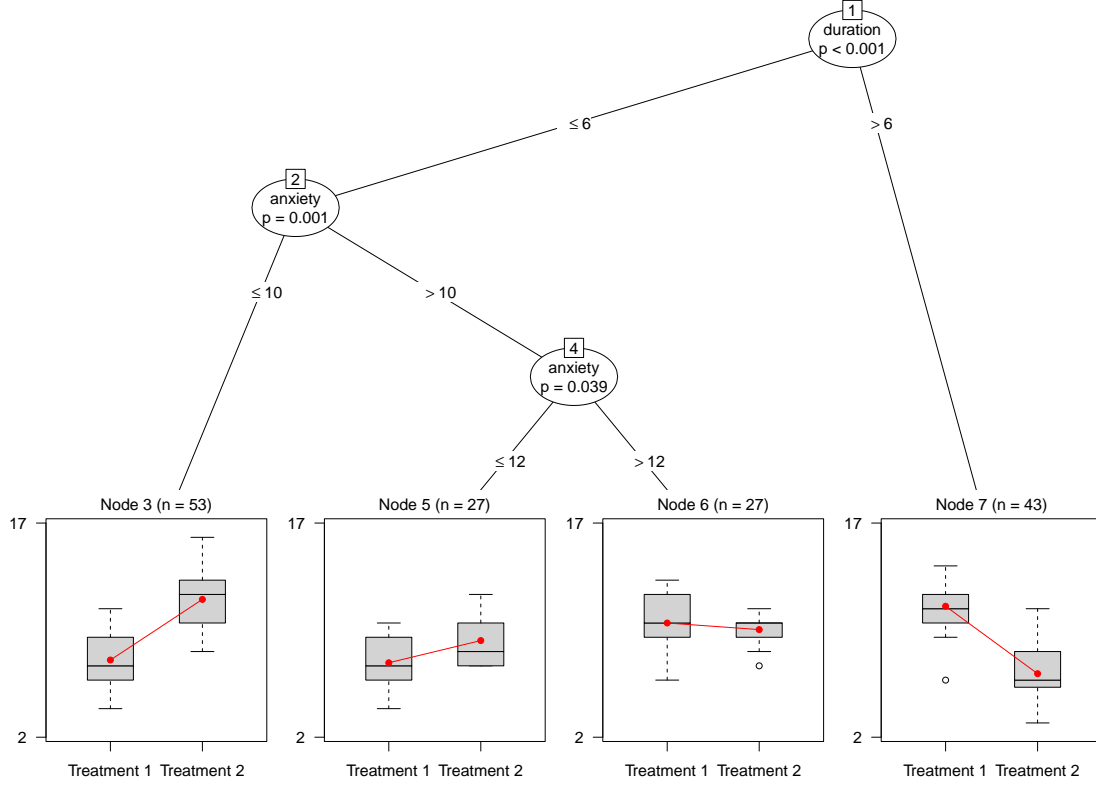


Figure 3. Example of a tree representation of model-based recursive partition, based on the artificial motivating dataset. Three additional covariates (anxiety questionnaire score, duration of depressive symptoms at baseline in months and age) were used as potential splitting variables.

1 the objective function in both partitions (Zeileis et al., 2008). In step (4), steps (1) through
 2 (3) are repeated in each partition, until the null hypothesis of parameter stability can no
 3 longer be rejected (or the subsets become too small).

4 Due to the binary recursive nature of MOB, the resulting partition can be represented
 5 as a binary tree. If the partitioning is based on a GLM, the result is a GLM tree, with a
 6 local fixed-effects regression model in every j -th ($j = 1, \dots, J$) terminal node or subgroup:

$$g(\mu_{ij}) = x_i^\top \beta_j \quad (3)$$

7 Note that, if the recursive subgroup structure (i.e., the partition) were known, the tree could
 8 be estimated as a single GLM where all x -variables interact with the subgroup indicator.
 9 Somewhat more formally, the model could then be written: $g(\mu_i) = x_i^{*\top} \beta^*$, where x_i^* are
 10 the values of the $2J$ interactions between the subgroups from the tree, and the elements
 11 of x_i . The corresponding β^* would have length $2J$ as well, containing subgroup-specific

fixed-effects coefficients.

Figure 3 shows the GLM tree grown on the motivating dataset. By using the three auxiliary covariates (anxiety, duration and age), MOB partitioned the observations into four subgroups, each with a different estimate for β_j . Age was correctly not detected as a partitioning variable, and the left- and rightmost subgroups are in accordance with the true treatment-subgroup interactions described above. However, the two subgroups in the middle do not represent true subgroups, which is due to the clustered structure of the dataset not being taken into account.

Generalized linear mixed-effects model trees

For datasets containing observations from multiple clusters (e.g., trials or research centers), application of a GLMM would be more appropriate. The model in Equation 2 is then extended to include cluster-specific, or random effects:

$$g(\mu_i) = x_i^\top \beta + z_i^\top b \quad (4)$$

For a simple random intercept z_i is a unit vector of length M , of which the m -th element takes a value of 1, and all other elements take a value of 0; m ($m = 1, \dots, M$) denotes the cluster which observation i is part of. Further, b is a random vector of length M , each m -th element representing the random intercept for cluster m . For simplicity we only employ a cluster-specific intercept here, but further random effects can be easily included in z_i . Furthermore, within the GLMM it is assumed that b is normally distributed, with mean zero and variance σ_b^2 . The parameters of the GLMM can be estimated with, for example, maximum likelihood (ML) and restricted ML (REML).

Note that, if the random-effects coefficients were known, the model could be estimated by a simple GLM as in Equation 2 where $z_i^\top b$ would only be added as an offset (i.e., a variable with a fixed coefficient of 1) to the linear predictor.

Although the random-effects part of the GLMM in Equation 4 accounts for the nested structure of the dataset, the global fixed-effects part may not describe the data well. Therefore, we propose the GLMM tree model, in which the fixed-effects part may be partitioned as in Equation 3 while still adjusting for random effects:

$$g(\mu_i) = x_i^\top \beta_j + z_i^\top b \quad (5)$$

To estimate the parameters of this model, we take an approach similar to that of the mixed-effects regression tree (MERT) approach of Hajjem et al. (2011) and Sela and Simonoff (2012). In the MERT approach, the fixed-effects part of a GLMM is replaced by a CART tree with constant fits in the nodes, and the random-effects part is estimated as usual. To estimate a MERT, an iterative approach is taken, alternating between (1) assuming random

effects known, allowing for estimation of the CART tree, and (2) assuming the CART tree known, allowing for estimation of the random effects.

For estimating GLMM trees, we take this approach two steps further: (1) Instead of a CART tree with constant fits to estimate the fixed-effects part of the GLMM, we use a GLM tree. This allows not only for detection of differences in intercepts across terminal nodes, but also for detection of differences in slopes such as treatment effects. (2) By using generalized linear (mixed) models, the response may also be a binary or count variable instead of a continuous variable. The GLMM tree algorithm takes the following steps to estimate the model in Equation 5:

Step 0: Initialize by setting r and all values $\hat{b}_{(r)}$ to 0.

Step 1: Set $r = r + 1$. Estimate GLM tree $(x_i^\top \hat{\beta}_{j(r)})$, with $z_i^\top \hat{b}_{(r-1)}$ as an offset.

Step 2: Estimate random effects in the mixed-effects model $x_i^\top \hat{\beta}_{j(r)} + z_i^\top \hat{b}_{(r)}$ with subgroups $j(r)$ from the GLM tree.

Step 3: Repeat Steps 1 and 2 until convergence.

The algorithm initializes by setting b to 0, since the random effects are initially unknown. In every iteration, the GLM tree and coefficients $\beta_{j(r)}$ and $b_{(r)}$ are re-estimated. The GLM tree is estimated, given the estimated random effects from the last iteration, and the random effects are estimated, given the estimated GLM tree from the current iteration. Iterations are continued until convergence, which is monitored by computing the log-likelihood criterion of the mixed-effects model in Equation 5. Typically, this converges if the tree does not change from one iteration to the next.

In Figure 1, the result of applying the GLMM tree algorithm to the motivating dataset is presented. By taking into account the clustering of observations, the spurious split involving the anxiety variable no longer appears in the tree and the true treatment subgroups have been recovered.

Simulation-based evaluation

To assess the performance of GLMM tree, we carried out three simulation studies: In Studies I and II, we compared the performance of GLMM tree with that of GLM tree. In Study I, we assessed the accuracy of both tree-based algorithms in datasets with treatment-subgroup interactions. In Study II, we assessed the Type-I error of both tree-based algorithms in datasets without treatment-subgroup interactions. In Study III, we compared the performance of GLMM tree in detecting treatment interactions with that of a more classical approach of interaction detection: a GLMM with pre-specified interactions.

General simulation design

In all simulation studies, we varied the following data-generating parameters:

1. Three levels for sample size: $N = 200$, $N = 500$, $N = 1000$.
2. Two levels for the number of potential partitioning covariates U_1 through U_K : $K = 5$ and $K = 15$.
3. Two levels for the intercorrelation between the potential partitioning covariates U_1 through U_K : $\rho_{U_k, U_{k'}} = 0.0$, $\rho_{U_k, U_{k'}} = 0.3$.
4. Three levels for the number of clusters: $M = 5$, $M = 10$, $M = 25$.
5. Three levels for the population standard deviation of the normal distribution from which the cluster specific intercepts were drawn: $\sigma_b = 0$, $\sigma_b = 5$, $\sigma_b = 10$.
6. Two levels for the intercorrelation between b and one of the U_k variables: b and all U_k covariates uncorrelated, b correlated with one of the U_k covariates.

In addition, each simulation study required a different design for main and interaction effects, which will be described in their respective sections below.

Following the approach of Dusseldorp and Van Mechelen (2014), all partitioning covariates U_1 through U_K were drawn from a multivariate normal distribution with means $\mu_{U_1} = 10$, $\mu_{U_2} = 30$, $\mu_{U_4} = -40$, and $\mu_{U_5} = 70$. Means for the other covariates (i.e., μ_{U_3} , and μ_{U_6} through $\mu_{U_{15}}$) were drawn from a discrete uniform distribution on the interval $[-70, 70]$. All covariates U_1 through U_{15} had the same standard deviation: $\sigma_{U_k} = 10$. Correlations between U variables vary according to the third facet of the data-generating design described above.

To generate the cluster-specific intercepts b_m , we partitioned the sample into M equally-sized clusters, conditional on one of the variables U_1 through U_5 , producing the correlations in the sixth facet of the simulation design. For each cluster, a single value b_m was drawn from a normal distribution with mean 0 and the value of σ_b given by the fifth facet of the simulation design. When b was correlated with one of the potential partitioning variables, this variable was randomly selected.

For every observation, we generated a binomial variable (with probability .5) as an indicator for treatment type. Random errors ϵ were drawn from a normal distribution with $\mu_\epsilon = 0$ and $\sigma_\epsilon = 5$. The value of the outcome variable y_i was calculated as the sum of the random intercept, (node-specific) fixed effects and the random error term.

Software

R (R Core Team, 2016) was used for the generation and analysis of all datasets. The **partykit** package (version 1.0-2; Hothorn & Zeileis, 2015) was employed for estimating GLM trees, using the **lmtree** function for normal linear regressions and **glmtree** for other response distributions. For estimating GLMMs the **lmer** function (or **glmer** function, re-

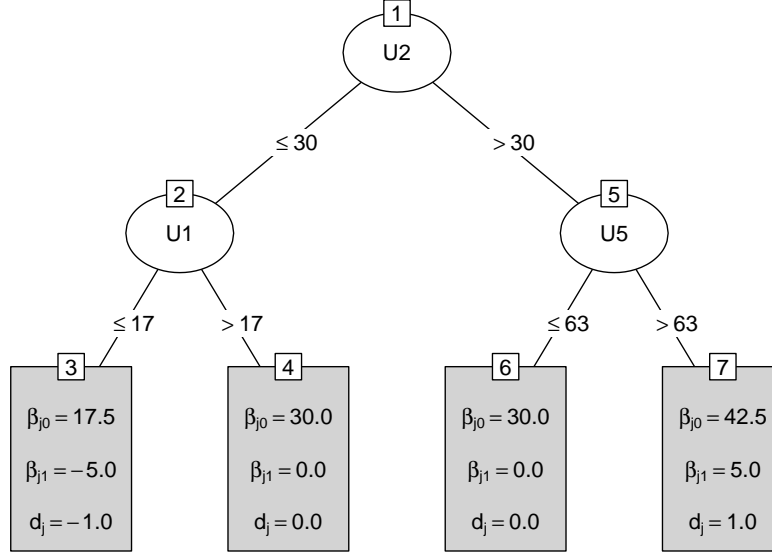


Figure 4. Data-generating model for treatment-subgroup interactions. Parameter d denotes the standardized mean difference between the outcomes of Treatment 1 and 2 (i.e., $\beta_{j1}/\sigma_\epsilon$).

spectively) from the `lme4` package (version 1.1-7; Bates, Mächler, Bolker, & Walker, 2015) was employed, using restricted maximum likelihood (REML) estimation.

For estimation of GLMM trees the former two packages were combined in a new package `glmertree` (version 0.1-0; Fokkema & Zeileis, 2016; available from R-Forge). This package provides functions `lmertree` and `glmertree` that iterate between estimation of the `lmtree`/`glmmtree` model and the `lmer`/`glmer` model.

The significance level α for the parameter instability tests in the trees was set to .05, with a Bonferroni correction applied for multiple testing, in all simulations. The minimum number of observations per node in the tree was set to 20 and the maximum tree depth was set to four, thus limiting the number of potential subgroups to eight.

In Study III, the `lmerTest` package (version 2.0-32; ?, ?) was used to assess statistical significance of fixed-effects predictors in GLMMs. The `lmerTest` package calculates effective degrees of freedom and p -values based on Satterthwaite approximations.

Study I: Method

Treatment-subgroup interaction design. For generating datasets with treatment-subgroup (or piecewise) interactions, we used a design from Dusseldorp and Van Mechelen (2014) which is depicted in Figure 4. Figure 4 shows four subgroups, characterized by val-

ues of the (true) partitioning variables U_2 , and U_1 or U_5 . Two of the subgroups have mean differences in treatment outcome, indicated by a non-zero value of β_{j1} , and two subgroups do not have mean differences in treatment outcome, indicated by a β_{j1} value of 0.

In this simulation study, some of the potential partitioning covariates were true partitioning covariates (U_1 , U_2 and U_5), whereas others were not (e.g., U_3 and U_4). Therefore, an extra level was added to the sixth facet of the *General simulation design* described above. Also, to assess the effect of the size of the treatment-effect differences in nodes 3 and 7, an additional facet was added:

6. Three levels for the intercorrelation between b and one of the U_k variables: b and all U_k covariates uncorrelated, b correlated with one of the true partitioning covariates (U_1 , U_2 or U_5), b correlated with one of the noise variables (U_3 or U_4).

7. Two levels for β_{j1} , the unstandardized mean difference in treatment outcomes. The absolute value of the treatment-effect difference was varied to be $|\beta_{j1}| = 2.5$ (corresponding to a medium effect size, Cohen's $d = 0.5$; Cohen, 1992) and $|\beta_{j1}| = 5.0$ (corresponding to a large effect size; Cohen's $d = 1.0$).

Again, 50 datasets were generated for each cell of the design. As a result, $50 \times 3 \times 2 \times 2 \times 3 \times 3 \times 3 \times 2 = 32,400$ datasets were generated. In every dataset, the outcome variable was calculated as $y_i = x_i^\top \beta_j + z_i^\top b_m + \epsilon_i$.

Assessment of performance. Performance of the algorithms was assessed by means of the total number of nodes in estimated GLM and GLMM trees, the accuracy of the resulting trees and predictive accuracy. An accurately recovered tree was defined as a tree with (1) a total of seven nodes (2) the first split involving variable U_2 and a value of 30 ± 5 , (3) the next split on the left involving variable U_1 and a value of 17 ± 5 , and (4) the next split on the right involving variable U_5 and a value of 63 ± 5 . The allowance of ± 5 equals plus or minus half the population standard deviation of the partitioning variable (σ_{U_k}).

Predictive accuracy of each method was assessed by calculating the correlation between true and predicted treatment-effect differences in every dataset. To prevent overly optimistic estimates of predictive accuracy (Hastie, Tibshirani, & Friedman, 2009), predictive accuracy was assessed using test datasets. Test datasets were generated from the same population as training datasets, but test observations were not drawn from the same clusters as the training observations, but from 'new' clusters.

To identify the most important predictors of tree size, tree accuracy and predictive accuracy, we used ANOVAs and GLMs with algorithm type and the parameters of the data-generating process as independent variables. First-order interactions between algorithm type and each of the data-generating parameters were also included. The most important predictors were further assessed with graphical displays.

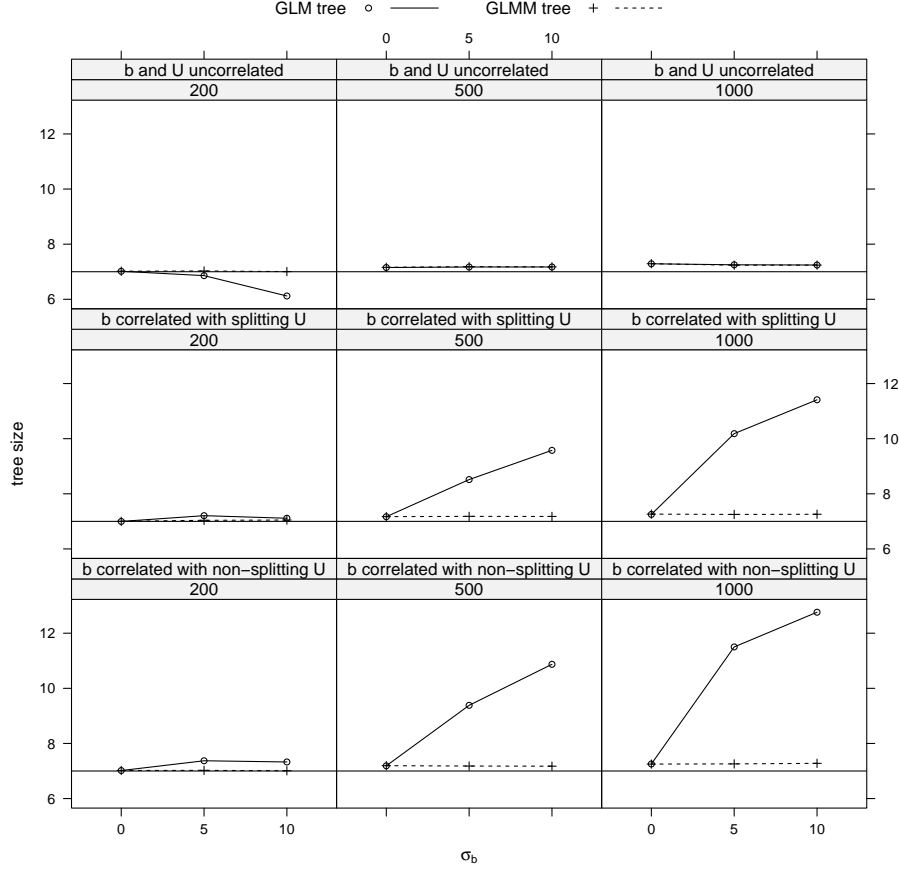


Figure 5. Average tree size of GLM and GLMM trees for datasets with treatment-subgroup interactions. Rows represent levels of dependence between random effects (b) and one of the partitioning variables U_k ; columns represent levels of sample size.

1 Study I: Results

2 *Tree size.* GLMM tree yielded trees with an average size of 7.15 nodes (SD = 0.61),
3 whereas GLM tree yielded an trees with an average size of 8.15 nodes (SD = 2.05), indicating
4 that GLM trees involved more spurious splits than GLMM trees on average. The estimated
5 probability that a dataset was erroneously not partitioned (the Type-II error) was 0 for
6 both algorithms.

7 An ANOVA indicated that the most important predictors of tree size were algorithm
8 type, sample size, variance of the random intercept, and the correlation between partitioning
9 variables. The graphical display in Figure 5 indicates that GLMM tree grows trees of about
10 the true tree size in all conditions. In the absence of random effects (i.e., $\sigma_b = 0$), GLM
11 and GLMM trees grow trees of about the same size. However, clear differences in tree size

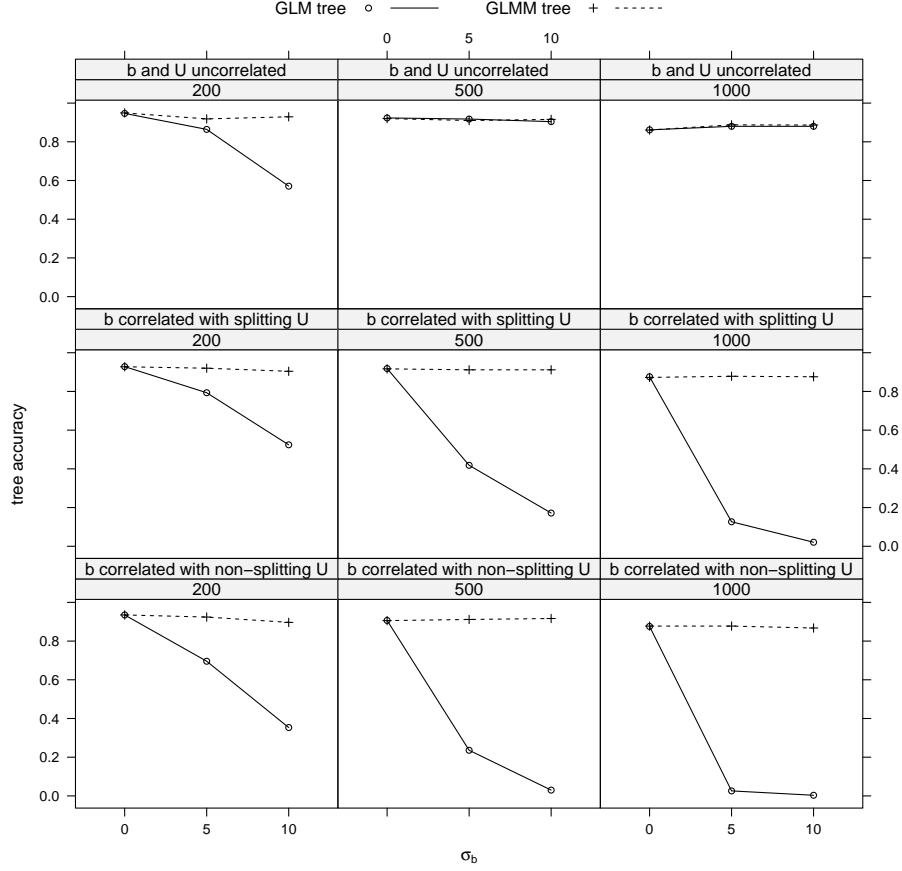


Figure 6. Tree accuracy of GLM and GLMM tree. Tree accuracy is defined as the proportion of datasets in which the true tree was accurately recovered. Rows represent levels of dependence between random effects (b) and one of the partitioning variables U_k , columns represent levels of sample size.

1 were observed when $\sigma_b > 0$. When sample size is small and the random intercept is not
 2 correlated with a partitioning covariate, GLM tree lacks power and grows trees that are
 3 too small, on average. When sample size is larger and random intercepts are not correlated
 4 with a partitioning covariate, GLM and GLMM trees are about equally sized. However,
 5 when random intercepts are correlated with a partitioning covariate, GLM starts to create
 6 spurious splits, especially with larger sample sizes and larger σ_b values.

7 *Accuracy of recovered trees.* To assess the accuracy of trees recovered by both algo-
 8 rithms, we inspected variables and values that were selected for partitioning. For the first
 9 split, GLMM tree selected the true partitioning variable (U_2) in all datasets, and GLM tree
 10 in all but one datasets. The splitting value for U_2 selected for the first split was 29.94 for

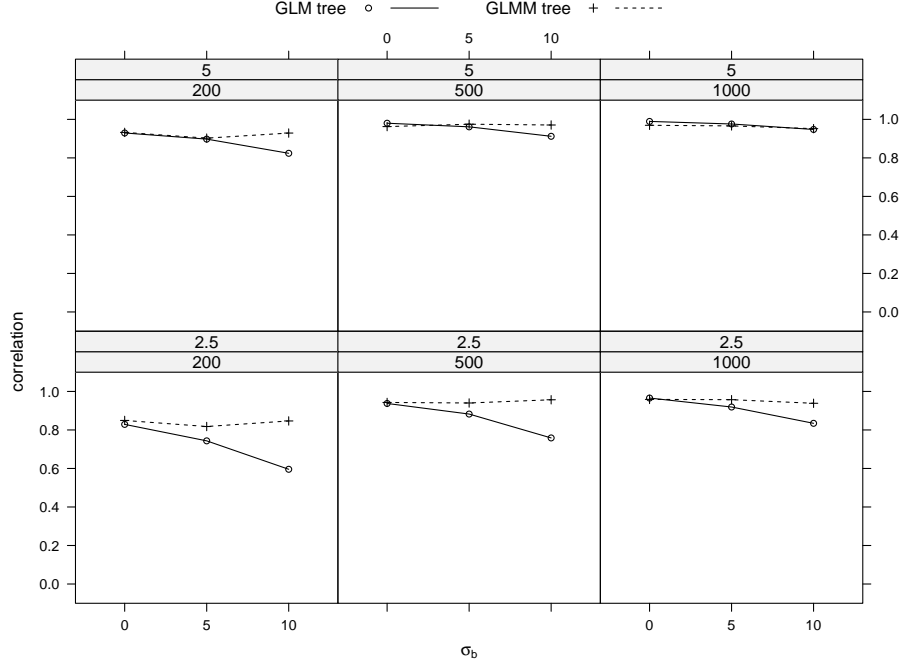


Figure 7. Average predictive accuracy of GLM and GLMM trees. Predictive accuracy of trees is defined as the correlation between the true and predicted outcome differences between Treatment 1 and 2. Rows represent the levels of the absolute value of the unstandardized treatment-effect difference in subgroups with treatment-effect differences, columns represent levels of sample size.

both GLM and GLMM tree, which is very close to the true splitting value of 30 (Figure 4).

Further splits were more accurately recovered by GLMM tree, which accurately recovered the partition in 90.40% of datasets, and GLM tree in only 61.44% of datasets. A GLM with tree accuracy as the outcome variable indicates that the most important predictors of tree accuracy were algorithm type, sample size, variance of the random intercept, and level of dependence between the partitioning variables and the random intercept. The graphical display in Figure 6 indicates that in the absence of random effects, the trees recovered by GLM and GLMM tree were about equally accurate. In the presence of random effects, GLM trees were much less accurate than GLMM trees when random effects were correlated with a partitioning covariate. When random intercepts were not correlated with one of the U_k variables, GLMM tree outperformed GLM tree when sample size was small (i.e., $N = 200$).

Predictive accuracy. The predicted treatment-effect differences of GLMM tree were closer to the true differences than those of GLM tree, with a mean correlation of .93 (SD = 0.12) for GLMM tree and .88 (SD = 0.19) for GLM tree. An ANOVA indicated that the most important predictors of predictive accuracy were algorithm type, sample size, variance

of the random intercept, and the effect size of the treatment difference. The graphical display in Figure 7 shows clear main effects of sample size and treatment-effect differences for both algorithms. In the absence of random effects (i.e., $\sigma_b = 0$), predictions of GLM and GLMM tree were about equally accurate. In the presence of random effects, GLM tree predictions were clearly less accurate than those of GLMM tree when treatment-effect differences were smaller (i.e., Cohen's $d = .5$). This effect was also observed, but much less pronounced, when treatment-effect differences were larger (i.e., Cohen's $d = 1.0$). In other words, GLMM tree outperforms GLM tree in the presence of random effects, especially when sample size and treatment-effect differences are smaller.

Study II: Method

Design. In the second simulation study we assessed the Type-I error of GLM and GLMM tree. Type-I error was defined as the proportion of datasets without treatment-subgroup interactions which were erroneously partitioned by the algorithm. In the simulated datasets in this study there was only a main effect of treatment in the population. Put differently, there was only a single global value of $\beta_j = \beta$ in every dataset.

To assess the effect of the treatment-effect difference β , an additional facet was added to the *General simulation design* described above:

7. Two levels for β , the unstandardized mean difference in treatment outcomes. The absolute value of the treatment-effect difference was varied to be $|\beta_{j1}| = 2.5$ (corresponding to a medium effect size, Cohen's $d = 0.5$) and $|\beta_{j1}| = 5.0$ (corresponding to a large effect size; Cohen's $d = 1.0$).

For each cell of the *General simulation design* described above, 50 datasets were generated. As a result, $50 \times 3 \times 2 \times 2 \times 3 \times 3 \times 2 \times 2 = 21,600$ datasets without treatment-subgroup interactions were generated. In every dataset, the outcome variable was calculated as $y_i = x_i^\top \beta + z_i^\top b_m + \epsilon_i$.

Assessment of performance. The total number of nodes in estimated GLM and GLMM tree were calculated. A tree of size > 1 was classified as a Type-I error. To identify the most important predictors of Type-I error, we used a GLMs with algorithm type and the parameters of the data-generating process as independent variables. First-order interactions between algorithm type and each of the data-generating parameters were also included. The most important predictors were further assessed with a graphical display.

Study II: Results

In datasets without treatment-subgroup interactions, average tree size was 1.09 (SD = 0.44) for GLMM tree, and 2.02 (SD = 1.68) for GLM tree. The average Type-I error rate was only .04 for GLMM tree, and .33 for GLM tree. Main predictors of type-I error were

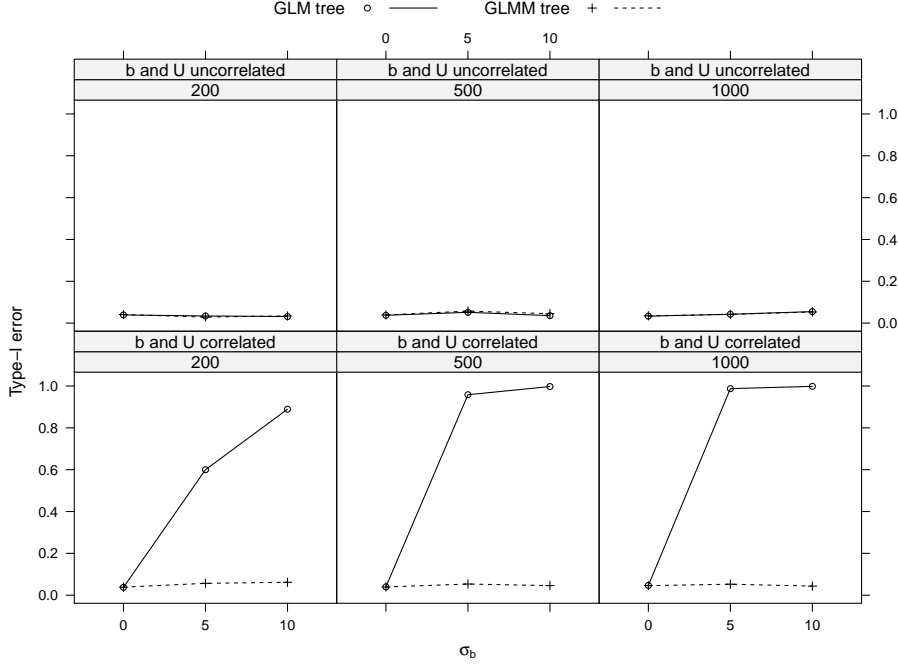


Figure 8. Type-I error of GLM and GLMM trees. Rows represent dependence between random effects (b) and one of the partitioning variables U_k ; columns represent sample size.

found to be sample size, variance of the random intercept, and dependence between the random intercept and one of the partitioning variables.

The graphical display in Figure 8 indicates that GLMM tree has a Type-I error rate somewhat below the pre-specified α level under all circumstances. The same goes for GLM tree, whenever the random intercept has 0 variance, or the random intercept is not correlated to one of the partitioning covariates. However, when the random intercept is correlated to one of the potential partitioning covariates, the type-I error rapidly increases for GLM tree. With large sample sizes and large variance of the random intercept, GLM tree will almost certainly create spurious splits.

Study III: Method

Interaction design. The treatment-subgroup interactions in Study I (Figure 4) can be referred to as piecewise interactions, as their effect is a stepwise function of the moderator (partitioning) variables U_1 , U_2 and U_5 . Trees are pre-eminently suited for recovering such piecewise interactions, but may have difficulty when the true interactions are continuous functions of moderator variables (for example, $U_1 \cdot U_2$). At the same time, linear regression models with pre-specified interaction terms may perform well in recovering continuous in-

Table 1: Fixed-effects terms in simulations with continuous and combined continuous and piecewise interaction designs.

term	β_3	β_4	β_6	β_7	β
intercept	27	27	27	27	27
U_2	.1	.1	.1	.1	.1
$U_2 \cdot U_1$	-.357	0	0	0	-.357
$U_2 \cdot U_5$	0	0	0	.357	.357
$U_2 \cdot U_1 \cdot treatment$	-.151	0	0	0	-.151
$U_2 \cdot U_5 \cdot treatment$	0	0	0	.151	.151

Note. Subscripted β values refer to the combined piecewise and continuous interaction design, with the values of the subscripts denoting the terminal nodes in Figure 4. β without a subscript refers to the global coefficients in the purely continuous interaction design.

teractions, but may have difficulty in recovering piecewise interactions. Therefore, in the third simulation study, we added a seventh facet to the *General simulation design* described above:

7. Three levels for interaction type: continuous, piecewise and a combination of both.

For datasets with purely piecewise interactions, the same partition as in Study II (depicted in Figure 4) was used. In other words, the outcome variable in this design was calculated as $y_i = x_i^\top \beta_j + z_i^\top b + \epsilon_i$, with the value of β_j depending on the values of U_2 , U_1 and U_5 .

For the datasets with both piecewise and continuous interactions, the partition as depicted in Figure 4 was used. However, the fixed-effects part $x_i^\top \beta_j$ in each of the terminal nodes comprised continuous main and (treatment) interaction effects of the partitioning variables. These node-specific effects are described in Table 1. β_j values were chosen to yield the same treatment-subgroup means as in Figure 4. Continuous interactions were created using centered U_k variables, calculated by subtracting the variable mean.

For datasets with purely continuous interactions, the outcome variable was calculated as $y_i = x_i^\top \beta + z_i^\top b + \epsilon_i$. Here, β has a global value and no subscript, and comprised purely continuous main and (treatment) interaction effects of the moderator variables U , as described in Table 1.

Furthermore, in this simulation study, the number of cells in the design was reduced by limiting the fourth facet of the data-generating design to a single level ($M = 25$ clusters), as Study I and II indicated no effects of the number of clusters. The fifth facet of the

1 data-generating design was limited to two levels ($\sigma_b = 2.5$ and $\sigma_b = 7.5$). As a result,
 2 $50 \times 3 \times 2 \times 2 \times 1 \times 2 \times 2 \times 3 = 7,200$ training datasets were generated for this study.

3 *GLMMs with pre-specified interactions.* GLMMs were estimated by specifying main
 4 effects for all covariates U_k and the treatment indicator, first-order interactions between
 5 all pairs of covariates U_k , and second-order interactions between all pairs of covariates
 6 U_k and treatment. Continuous predictor variables were centered by subtracting the mean
 7 value, before calculating and including the interaction term in the GLMM. Predictions were
 8 obtained by refitting a GLMM on the training data, using only the predictors with p-values
 9 $< .05$ in the original GLMM. Predictions for test observations were calculated using the
 10 refitted GLMMs.

11 *Assessment of performance.* Predictive accuracy was assessed in terms of the cor-
 12 relation between the true and predicted treatment-effect differences in test datasets. To
 13 identify the most important predictors of predictive accuracy, we used an ANOVA with
 14 algorithm type and the parameters of the data-generating process as independent variables.
 15 First-order interactions between algorithm type and each of the data-generating parameters
 16 were also included. The most important predictors were further assessed with a graphical
 17 display.

18 *Study III: Results*

19 GLMM tree showed somewhat higher average correlations between the true and pre-
 20 dicted treatment-effect differences: average accuracy was .54 (SD = .40) for GLMM trees
 21 and .54 (SD = .43) for GLMMs. The most important predictors of predictive accuracy were
 22 algorithm type, sample size, the number of covariates, and type of interaction (continuous,
 23 piecewise or both). These effects were further assessed by means of the graphical display
 24 depicted in Figure 9.

25 As Figure 9 indicates, GLMM trees show highest predictive accuracy in datasets with
 26 purely piecewise interactions, whereas GLMMs show highest predictive accuracy in datasets
 27 with purely continuous interactions. Performance of both methods improves with increasing
 28 sample size. Furthermore, performance of GLMM tree is not affected by the number of
 29 covariates, whereas the predictive accuracy of GLMMs deteriorates when the number of
 30 covariates increases, especially when the true interactions are not purely continuous. This
 31 indicates that GLMM tree may be especially useful for exploratory purposes, where there
 32 are many potential moderator variables. Also, GLMM trees may provide simpler models,
 33 as the average number of fixed-effects predictor variables (main effects and interactions)
 34 evaluated for prediction with GLMMs was 12.30, whereas the average number of inner
 35 nodes that need to be evaluated for prediction with GLMM tree was only 3.38.

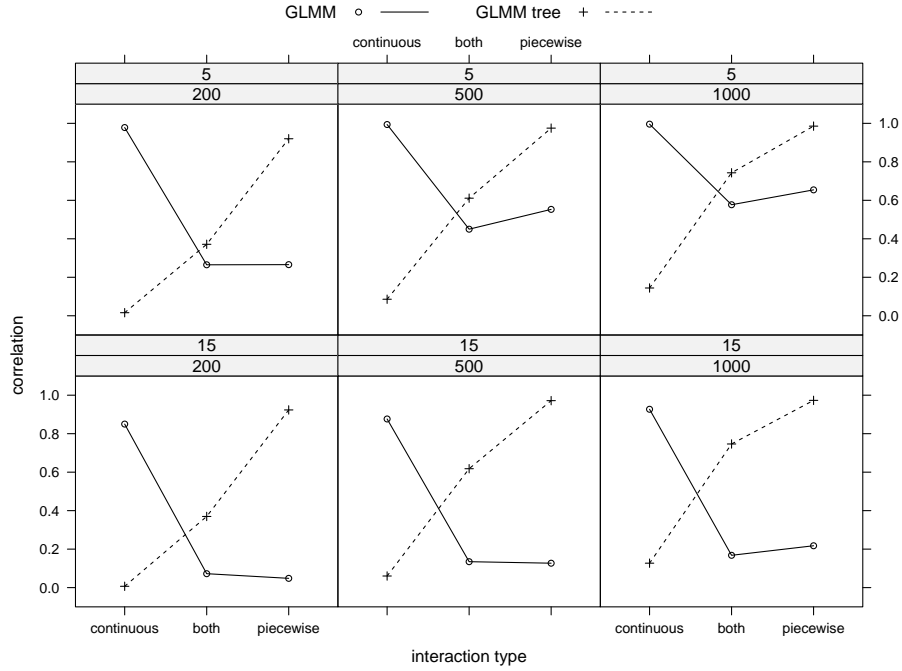


Figure 9. Average predictive accuracy of GLMMs and GLMM trees. Predictive accuracy of trees is defined as the correlation between the true and predicted differences between Treatment 1 and 2. Columns represent sample size, rows represent the number of covariates.

Application: Individual patient-level meta-analysis dataset on
treatments for depression

Method

To further illustrate the use of GLM and GLMM tree, we applied both algorithms to a dataset from an individual-patient data meta-analysis of Cuijpers et al. (2014). This meta-analysis was based on patient-level observations from 14 RCTs, comparing the effects of psychotherapy (cognitive behavioral therapy; CBT) and pharmacotherapy (PHA) in the treatment of depression. The study of Cuijpers et al. (2014) was aimed at establishing whether gender is a predictor or moderator of the outcomes of psychological and pharmacological treatments for depression. Treatment outcomes were assessed by means of the 17-item Hamilton Rating Scale for Depression (HAM-D; Hamilton, 1960). Cuijpers et al. (2014) found no indication that gender predicted or moderated treatment outcomes. Further details on the dataset can be found in Cuijpers et al. (2014).

In our analyses, posttreatment HAM-D score was the outcome variable, and potential partitioning variables were age, gender, level of education, presence of a comorbid anxiety

disorder at baseline, and pretreatment HAM-D score. The predictor variable in the linear model was treatment type (0 = CBT and 1 = PHA). An indicator for study was used as the cluster indicator.

In RCTs, ANCOVAs are often employed, to linearly control posttreatment values on the outcome measure for pretreatment values. Therefore, we estimated the GLM and GLMM trees using posttreatment HAM-D scores, controlled for the linear effects of pretreatment HAM-D scores. The trees were grown using data of patients with complete observations; that is, observations with non-missing values for potential partitioning variables, and pre- and posttreatment HAM-D score. As a result, data from 694 patients from 7 studies were included in the analyses. Results of our analysis may therefore not be representative of the complete dataset of the meta-analysis by Cuijpers et al. (2014).

To provide a standardized estimate of the treatment effect differences in the final nodes of the trees, we calculated node-specific Cohen's d values. Cohen's d was calculated by dividing the node-specific predicted treatment outcome difference by the node-specific pooled standard deviation. Confidence intervals for Cohen's d could not be calculated, as these can not take into account the exploratory nature of our analyses (i.e., variable and split point selection).

To compare the two tree-based approaches with a more classical approach of interaction detection, we also fitted a GLMM with pre-specified interactions on the dataset. In the GLMM, the posttreatment HAM-D scores were controlled for the linear effects of pretreatment HAM-D scores, and then regressed on a random intercept, main effects of treatment and the potential moderators (partitioning variables), and interactions between treatment and the potential moderators. As it is not known in advance how to interact the potential moderators, higher-order interactions were not included.

The predictive accuracy of the GLM and GLMM trees and the GLMM was assessed by calculating average correlations between observed and predicted HAM-D posttreatment scores, based on 50-fold cross validation.

The results of recursive partitioning techniques are known to be potentially unstable, in the sense that small changes in the dataset may substantially alter the variables or values selected for partitioning. Therefore, we used the R package **stablelearner** (Philipp, Zeileis, & Strobl, 2016) to assess the stability of the selected splitting variables and values. Using the **stabletree** function, we calculated the number of times a variable and value were selected for partitioning over 500 subsamples of size $.9 \times N$ of the dataset.

Results

The tree and effects sizes resulting from application of GLM tree are presented in Figure 10. Those resulting from application of GLMM tree are presented in Figure 11.

The GLM tree (Figure 10) selected level of education as the first partitioning variable,

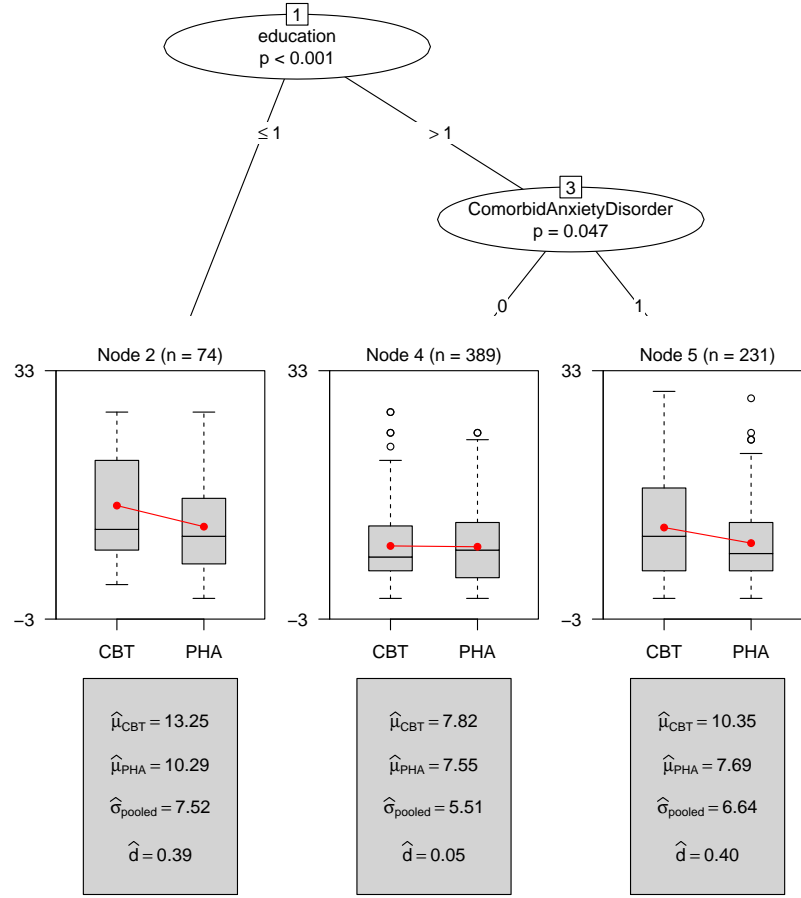


Figure 10. GLM tree for prediction of posttreatment total scores on the Hamilton Rating Scale for Depression (HAM-D). Upper panel: The y-axes of the boxplots represent posttreatment HAM-D scores, and the x-axes represent treatment levels (cognitive behavior therapy, CBT vs. pharmacotherapy, PHA). Lower panel: Subgroup-specific descriptive statistics.

1 and presence of a comorbid anxiety disorder as a second partitioning variable, for observa-
 2 tions with a higher level of education. By taking into account study-specific intercepts, the
 3 GLMM tree (Figure 11) indicates that the first split made by GLM tree may be a spurious
 4 one. The GLMM tree selected presence of a comorbid anxiety disorder as the only parti-
 5 tioning variable. The terminal nodes of Figure 11 show only a single treatment-subgroup
 6 interaction: for patients without a comorbid anxiety disorder, CBT and PHA provide more
 7 or less the same reduction in HAM-D scores (Cohen's $d = 0.05$; Figure 11). For patients
 8 with a comorbid anxiety disorder, PHA provides a greater reduction in HAM-D scores (Co-
 9 hen's $d = 0.39$; Figure 11). The estimated intraclass correlation coefficient for the GLMM

Table 2: Variable selection statistics

Variable	Selection frequency	
	GLM tree	GLMM tree
Education	.956	.014
ComorbidAnxietyDisorder	.398	.528
HRSDt0	.034	.002
Age	.000	.022
Gender	.002	.004

Note. Frequencies are calculated over 500 random subsamples of the complete dataset. Frequencies do not add up to 1, as trees may involve multiple or no splits.

1 tree was .05.

2 The GLMM with pre-specified treatment interactions yielded three significant predic-
 3 tors of treatment outcome: like in the GLMM tree, an effect of the presence of a comorbid
 4 anxiety disorder was found (main effect: $b = 2.29$, $p = .002$; interaction with treatment: b
 5 $= -2.10$, $p = .028$). Also, the GLMM indicated an interaction between treatment and age
 6 ($b = .10$, $p = .018$).

7 Assessment of predictive accuracy by means of 50-fold cross validation indicated better
 8 predictive accuracy for GLMM tree than for GLM tree and the GLMM with pre-specified
 9 interactions. The correlation between true and predicted posttreatment HAM-D total scores
 10 averaged over 50 folds was .272 ($var = .067$) for GLMM tree, .233 ($var = .064$) for the
 11 GLMM with prespecified interactions and .190 ($var = .084$) for GLM tree.

12 Table 2 presents statistics on the variables selected for partitioning in subsamples of
 13 the dataset. Note that the selection frequencies do not add up to 1, as trees may involve
 14 multiple, or no splits. Table 2 indicates that the presence of a comorbid anxiety disorder
 15 was selected for partitioning in the majority of GLMM trees grown on subsamples of the
 16 dataset, and all other variables were selected in less than 3% of the subsamples. As the
 17 comorbid anxiety disorder variable involved only a single splitting value, further assessment
 18 of the stability of splitting values was not necessary.

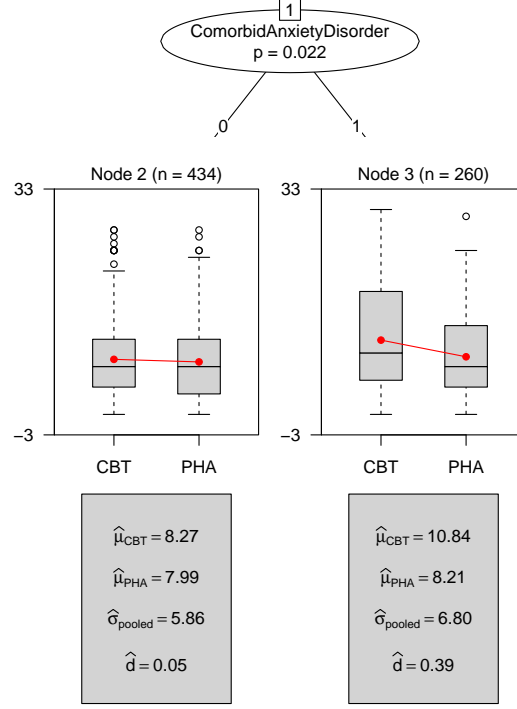


Figure 11. GLMM tree for prediction of posttreatment total scores on the Hamilton Rating Scale for Depression (HAM-D). Left panel: The y-axes of the boxplots represent posttreatment HAM-D scores, and the x-axes represent treatment levels (cognitive behavior therapy, CBT vs. pharmacotherapy, PHA). Right panel: Subgroup-specific descriptive statistics.

Discussion

In the current paper, we presented the GLMM tree algorithm, which allows for estimation of a GLM-based recursive partition, as well as estimation of random-effects parameters. As such, we hypothesized GLMM tree to be well suited for the detection of treatment-subgroup interactions in clustered datasets, which was confirmed by our simulation studies. We performed simulation studies to assess the accuracy of GLMM tree in detecting treatment-subgroup interactions and illustrated the application of GLMM tree on an existing dataset.

GLMM tree performed very well in recovering treatment-subgroup interactions, accurately recovering interactions in 90% of datasets with treatment-subgroup interactions. In contrast, GLM tree accurately recovered interactions in only 61% of these datasets. In the absence of treatment-subgroup interactions, GLMM tree erroneously detected subgroups in 4% of the datasets, whereas GLM tree erroneously detected subgroups in 33% of those

1 datasets. Put differently, the Type-I error rate of GLMM tree very closely approximated
 2 the α level used for testing parameter stability, whereas the Type-I error rate of GLM tree
 3 clearly exceeded this value.

4 The better performance of GLMM tree was mostly observed when random effects in
 5 the datasets were sizable, and random intercepts were correlated with potential partition-
 6 ing variables. In these instances, random effects gave rise to spurious subgroup detection
 7 (spurious splits) by GLM tree, both in datasets with and without treatment-subgroup in-
 8 teractions.

9 GLMM tree also provided better predictive accuracy than GLM tree, with an average
 10 correlation between true and predicted treatment differences of .94 for GLMM tree, and
 11 .88 for GLM tree. GLMM tree clearly outperformed GLM tree when random effects in the
 12 datasets were sizable, differences in treatment effects across terminal nodes were relatively
 13 small (i.e., Cohen's $d = .5$), and/or sample size was small (i.e., $N < 1,000$). Such treatment-
 14 effect differences and sample sizes may be quite common, even in multi-center clinical trials,
 15 and GLMM tree may provide a helpful tool for subgroup detection in those instances.

16 As expected, when random effects were absent from the simulated datasets, GLM and
 17 GLMM tree yielded very similar predictive accuracy. This indicates that GLMM tree can
 18 be applied whenever cluster-specific random effects are expected: In the absence of random
 19 effects, GLM tree and GLMM tree are expected to perform about equally well and in the
 20 presence of random effects GLMM tree is expected to outperform GLM tree.

21 Compared to treatment-interaction detection by means of GLMMs with pre-specified
 22 interaction effects, GLMM trees provided somewhat better accuracy, on average. When
 23 interactions were purely continuous, GLMMs outperformed GLMM trees, and when inter-
 24 actions were partly or purely piecewise, GLMM trees outperformed GLMMs. In addition,
 25 GLMM tree may have a clear advantage when there are a large number of potential mod-
 26 erator variables (i.e., > 5). Therefore, GLMM trees may be better suited than GLMMs for
 27 exploratory analyses, in which moderator variables need to be selected from a larger number
 28 of covariates. Furthermore, the number of terms in a GLMM increases quadratically with
 29 the number of potential moderator variables, yielding complex models. The trees in our
 30 simulations were limited to a maximum number of 7 inner nodes, and may therefore be
 31 much easier to use for prediction in practical decision-making contexts.

32 These findings are encouraging for the use of GLMM tree in the detection of
 33 treatment-subgroup interactions in datasets with clustered structures. However, it should
 34 be noted that the simulations show that GLMM tree performs very well, if the model is
 35 correctly specified. That is, if there are subgroups with respect to the variables specified as
 36 potential partitioning variables and these subgroups have different values for the parameters
 37 of the GLM, then GLMM tree will accurately recover those subgroups. However, misspec-

ification of the model will negatively affect performance. The most important source of misspecification would be the omission of relevant variables, either in the GLM or as partitioning variables. When relevant variables are omitted, GLMM tree can only approximate the true subgroups using the specified variables. Another source of misspecification would be the inclusion of irrelevant variables. Our simulations indicate that the performance of GLMM tree was not negatively affected by the number of potential moderator variables specified, but the power to detect subgroups may still be reduced with larger numbers of potential partitioning variables. Including irrelevant variables in the GLM may also negatively affect performance, although we have not assessed this in our simulations.

In the Application, we found GLMM tree to provide results that were more easily interpretable, and also more accurate than a GLM tree without random effects. In addition, to judge clinical relevance of the findings, we calculated node-specific effect sizes, as would often be done in RCTs or meta-analysis. Although we have limited ourselves to calculating Cohen's d in the current paper, equivalent values of the success rate difference or the number needed to treat can be calculated, but this would involve additional distributional assumptions (Kraemer & Kupfer, 2006). Node-specific effect sizes can also be used to prune trees, when a researcher prefers to have a final tree which is based on statistical as well as clinical significance. A topic for further research would be the development of splitting procedures based on effect sizes, as this would allow for taking into account clinical significant in the tree-growing process.

As discussed in the Introduction, several tree-based methods for treatment-subgroup interaction detection are available. These methods have different objectives, and there is not yet an agreed-upon single best method. In a simulation study, Sies and Van Mechelen (2016) found the method of Zhang, Tsiatis, Davidian, et al. (2012) to perform best, followed by model-based recursive partitioning. However, the method of Zhang et al. performed worst under some conditions of the simulation study in terms of the Type I error rate. Further research comparing tree-based methods for treatment-subgroup interaction detection is needed, especially for clustered datasets, as our simulations were limited to GLM and GLMM-based MOB.

Furthermore, it should be stressed that tree-based methods are exploratory tools. They can be used to detect predictors, interactions and non-linear effects in a data-driven way, but users should take the exploratory nature of such analyses into account. The resulting trees are potentially unstable, and stability of the results should be assessed, preferably in a dataset not used for training, or by multi-fold cross-validation or resampling techniques. In the Application, we have shown how the stability of splitting variable selection can be assessed using resampling techniques.

In conclusion, GLMM tree provided highly accurate recovery of treatment-subgroup

1 interactions and predictions of treatment effect differences, both in the presence and absence
 2 of cluster-specific random effects. Therefore, GLMM tree is a promising algorithm for the
 3 detection of treatment-subgroup interactions in datasets with a clustered structure, like for
 4 example in multi-center trials or individual-level patient data meta-analyses.

References

- Bates, D., Mächler, M., Bolker, B., & Walker, S. (2015). *Fitting linear mixed-effects models using lme4* (Vol. 67) (No. 1). doi: 10.18637/jss.v067.i01
- Breiman, L., Friedman, J., Olshen, R., & Stone, C. (1984). *Classification and regression trees*. New York: Wadsworth.
- Bryk, A. S., & Raudenbush, S. W. (1992). *Hierarchical linear models: Applications and data analysis methods*. Newbury Park, CA: Sage.
- Cohen, J. (1992). A power primer. *Psychological Bulletin*, 112(1), 155–159.
- Cooper, H., & Patall, E. A. (2009). The relative benefits of meta-analysis conducted with individual participant data versus aggregated data. *Psychological Methods*, 14(2), 165.
- Cuijpers, P., Weitz, E., Twisk, J., Kuehner, C., Cristea, I., David, D., ... Hollon, S. D. (2014). Gender as predictor and moderator of outcome in cognitive behavior therapy and pharmacotherapy for adult depression: An “individual-patients data” meta-analysis. *Depression and Anxiety*, 31(11), 941–951.
- Doove, L. L., Dusseldorp, E., Van Deun, K., & Van Mechelen, I. (2014). A comparison of five recursive partitioning methods to find person subgroups involved in meaningful treatment-subgroup interactions. *Advances in Data Analysis and Classification*, 8, 403–425.
- Driessen, E., Smits, N., Dekker, J., Peen, J., Don, F. J., Kool, S., ... Van, H. L. (2016). Differential efficacy of cognitive behavioral therapy and psychodynamic therapy for major depression: A study of prescriptive factors. *Psychological Medicine*, 46(4), 731–744.
- Dusseldorp, E., Doove, L., & Van Mechelen, I. (2015). Quint: An R package for the identification of subgroups of clients who differ in which treatment alternative is best for them. *Behavior Research Methods*, 1–14.
- Dusseldorp, E., & Meulman, J. J. (2004). The regression trunk approach to discover treatment covariate interaction. *Psychometrika*, 69(3), 355–374.
- Dusseldorp, E., & Van Mechelen, I. (2014). Qualitative interaction trees: A tool to identify qualitative treatment-subgroup interactions. *Statistics in Medicine*, 33(2), 219–237.
- Fokkema, M., & Zeileis, A. (2016). *glmertree: Generalized linear mixed model trees*. Retrieved from http://R-Forge.R-project.org/R/?group_id=261 (R package version 0.1-1)
- Foster, J. C., Taylor, J. M. G., & Ruberg, S. J. (2011). Subgroup identification from randomized clinical trial data. *Statistics in Medicine*, 30(24), 2867–2880.
- Hajjem, A., Bellavance, F., & Larocque, D. (2011). Mixed effects regression trees for clustered data. *Statistics & Probability Letters*, 81(4), 451–459.
- Hamilton, M. (1960). A rating scale for depression. *Journal of Neurology, Neurosurgery and Psychiatry*, 23(1), 56.
- Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The elements of statistical learning* (2nd ed.). New York: Springer.

- Higgins, J., Whitehead, A., Turner, R. M., Omar, R. Z., & Thompson, S. G. (2001). Meta-analysis of continuous outcome data from individual patients. *Statistics in Medicine*, 20(15), 2219–2241.
- Hothorn, T., & Zeileis, A. (2015, December). partykit: A modular toolkit for recursive partitioning in R. *Journal of Machine Learning Research*, 16, 3905–3909. Retrieved from <http://www.jmlr.org/papers/v16/hothorn15a.html>
- Koopman, L., Van der Heijden, G. J. M. G., Glasziou, P. P., Grobbee, D. E., & Rovers, M. M. (2007). A systematic review of analytical methods used to study subgroups in (individual patient data) meta-analyses. *Journal of Clinical Epidemiology*, 60(10).
- Kraemer, H. C., Frank, E., & Kupfer, D. J. (2006). Moderators of treatment outcomes: Clinical, research, and policy importance. *Journal of the American Medical Association*, 296(10), 1286–1289.
- Kraemer, H. C., & Kupfer, D. J. (2006). Size of treatment effects and their importance to clinical research and practice. *Biological Psychiatry*, 59(11), 990–996.
- Lipkovich, I., Dmitrienko, A., Denne, J., & Enas, G. (2011). Subgroup identification based on differential effect search – A recursive partitioning method for establishing response to treatment in patient subpopulations. *Statistics in Medicine*, 30(21), 2601–2621.
- Martin, D. (2015). *Efficiently exploring multilevel data with recursive partitioning* (Unpublished doctoral dissertation). University of Virginia.
- Merkle, E. C., & Zeileis, A. (2013). Tests of measurement invariance without subgroups: A generalization of classical methods. *Psychometrika*, 78(1), 59–82.
- Philipp, M., Zeileis, A., & Strobl, C. (2016). A toolkit for stability assessment of tree-based learners. In A. Colubi, A. Blanco, & C. Gatu (Eds.), *Proceedings of COMPSTAT 2016 – 22nd international conference on computational statistics* (pp. 315–325). Oviedo: The International Statistical Institute/International Association for Statistical Computing.
- R Core Team. (2016). R: A language and environment for statistical computing [Computer software manual]. Vienna, Austria. Retrieved from <https://www.R-project.org/>
- Seibold, H., Zeileis, A., & Hothorn, T. (2016). Model-based recursive partitioning for subgroup analyses. *International Journal of Biostatistics*, 12(1), 45–63. doi: 10.1515/ijb-2015-0032
- Sela, R. J., & Simonoff, J. S. (2012). RE-EM trees: A data mining approach for longitudinal and clustered data. *Machine Learning*, 86(2), 169–207.
- Sies, A., & Van Mechelen, I. (2016). *Comparing four methods for estimating tree-based treatment regimes*. (Submitted)
- Strobl, C., Malley, J., & Tutz, G. (2009). An introduction to recursive partitioning: Rationale, application, and characteristics of classification and regression trees, bagging, and random forests. *Psychological Methods*, 14(4), 323.
- Su, X., Tsai, C.-L., Wang, H., Nickerson, D. M., & Li, B. (2009). Subgroup analysis via recursive partitioning. *The Journal of Machine Learning Research*, 10, 141–158.
- Van den Noortgate, W., Opdenakker, M. C., & Onghena, P. (2005). The effects of ignoring a level in multilevel analysis. *School Effectiveness and School Improvement*, 16(3), 281–203.
- Zeileis, A. (2005). A unified approach to structural change tests based on ML scores, F statistics, and OLS residuals. *Econometric Reviews*, 24(4), 445–466.
- Zeileis, A., & Hornik, K. (2007). Generalized M-fluctuation tests for parameter instability. *Statistica*

Neerlandica, 61(4), 488–508.

Zeileis, A., Hothorn, T., & Hornik, K. (2008). Model-based recursive partitioning. *Journal of Computational and Graphical Statistics*, 17(2), 492–514.

Zhang, B., Tsiatis, A. A., Davidian, M., Zhang, M., & Laber, E. (2012). Estimating optimal treatment regimes from a classification perspective. *Stat*, 1(1), 103–114.

Zhang, B., Tsiatis, A. A., Laber, E. B., & Davidian, M. (2012). A robust method for estimating optimal treatment regimes. *Biometrics*, 68(4), 1010–1018.