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

M. Fokkema[1], N. Smits[2], A. Zeileis[3], T. Hothorn[4], H. Kelderman[5]

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

1            Abstract

2    Identification of subgroups of patients for which treatment A is more effective than
3  treatment B, and vice versa, is of key importance to the development of personalized
4  medicine. Tree-based algorithms are helpful tools for the detection of such interactions,
5  but datasets suited for their detection may often have a clustered or multilevel structure.
6  For example, observations in a dataset may be nested within research centers or trials.
7  Therefore, we propose generalized linear mixed-effects model trees (GLMM tree), allowing
8  for the detection of treatment-subgroup interactions, as well as taking into account the
9  nested structure of a dataset. The algorithm uses model-based recursive partitioning
10  (MOB) to detect treatment-subgroup interactions, and a GLMM for the estimation of
11  random-effects parameters. In a simulation study, we evaluate the performance of GLMM
12  tree and compare it with thattrees without random effects, and GLMMs with pre-specified
13  interactions. GLMM tree was found to accurately recover treatment-subgroup interactions
14  in 90% of datasets, and trees without random effects in 61% of datasets. GLMM tree also
15  outperformed trees without random effects in terms of predictive accuracy (.94 and .88 on
16  average, respectively) and Type-I error rates (4 and 33%, respectively). Furthermore, com-
17  pared to GLMMs with pre-specified interaction effects, GLMM tree showed equal accuracy
18  in predicting treatment-effect differences, but better accuracy in detecting interactions. We
19  illustrate the application of GLMM tree on an individual patient-level data meta-analysis
20  on treatments for depression. We conclude that GLMM tree is a promising algorithm for

the detection of treatment-subgroup interactions in clustered datasets.

## 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 preeminently suited to the detection of treatment-subgroup interactions: they can handle many potential predictor variables at once and can automatically detect (higher order) interactions between predictor variables (Strobl, Malley, & Tutz, 2009). Several tree-based algorithms to the detection of treatment-subgroup interactions have been developed (Dusseldorp & Van Mechelen, 2014; Dusseldorp & Meulman, 2004; Su, Tsai, Wang, Nickerson, & Li, 2009; Foster, Taylor, & Ruberg, 2011; Lipkovich, Dmitrienko, Denne, & Enas, 2011; Zeileis, Hothorn, & Hornik, 2008, 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, which allows users to select from a range of statistical methods, including trees. As such, the Zhang, Tsiatis, Davidian, et al. (2012) method can also be referred to as a tree-based method.

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

The GLMM tree algorithm builds on model-based recursive partitioning (MOB, Zeileis et al., 2008), which offers a flexible framework for subgroup detection. For example, GLM-based MOB has been applied to detect treatment-subgroup interactions for the treatment of depression (Driessen et al., 2016) and amyotrophic lateral sclerosis (Seibold, Zeileis, & Hothorn, 2015). GLMM tree allows for taking into account the clustered structure of datasets (which is not done in other purely tree-based treatment-subgroup interaction detection methods, e.g., Zeileis et al., 2008; Su et al., 2009; Dusseldorp & Van Mechelen, 2014), as well as treatment effect estimation with continuous and non-continuous response variables (which is not available in previously suggested regression trees with random effects, e.g., Hajjem, Bellavance, & Larocque, 2011; Sela & Simonoff, 2012).

In what follows, we introduce an artificial motivating dataset, which will be used to illustrate and explain treatment-subgroup interaction detection with GLM-based MOB and GLMM tree. In the Empirical Evaluation, we evaluate the performance of GLMM tree in simulated datasets, and compare it with that of GLM-based MOB and GLMMs with pre-specified interaction effects. In the Application, we apply the algorithm to an existing dataset of a patient-level meta-analysis on the effects of psycho- and pharmacotherapy for depression. Finally, in the Discussion we summarize the results and point out some directions for future research.

*Artificial motivating dataset*

We created generated a dataset representing observations on 150 participants in a randomized clinical trial. Every participant was randomly assigned to Treatment 1 or Treatment 2, and has a value for the response variable, with which the effect of treatment is assessed: the posttreatment total score on a depression inventory. For all participants, three covariate values are available: duration of depressive symptoms prior to treatment in months (duration, range 0–15); age in years (age, range 18–75); anxiety inventory total

1 score (anxiety, range 3–18).

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

11 Participants were part of one of ten clusters, each with a different value for the
12 intercept. Data was generated such that covariates and cluster-specific intercepts were un-
13 correlated. Also, 43% of variance in posttreatment depression scores was due to treatment-
14 subgroup interactions, and 8% of variance was due to cluster-specific variation.

15 *Model-based recursive partitioning*

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

24 For example, let us take a global GLM to estimate the overall treatment effect in
25 the artificial motivating dataset. The expectation $\mu_i$ of outcome $y_i$ given the treatment
26 regressors $x_i$ is modeled through a linear predictor and suitable link function[1]:

$$
\begin{aligned}
E[y_i|x_i] &= \mu_i, \quad\quad (1)\\
g(\mu_i) &= x_i^\top \beta, \quad\quad (2)
\end{aligned}
$$

27 where $x_i^\top \beta$ is the linear predictor for observation $i$ and $g$ is the link function. $\beta$ is a
28 vector of fixed-effects regression coefficients, the first element representing the intercept,
29 corresponding to the mean value of the linear predictor in the first treatment group, and the
30 second element representing the slope, which is the mean difference in the linear predictor
31 between the first and second treatment group. Thus, for simplicity we assume $x_i$ and $\beta$
32 to have length 2 in the current paper; that is, two treatment conditions and no additional

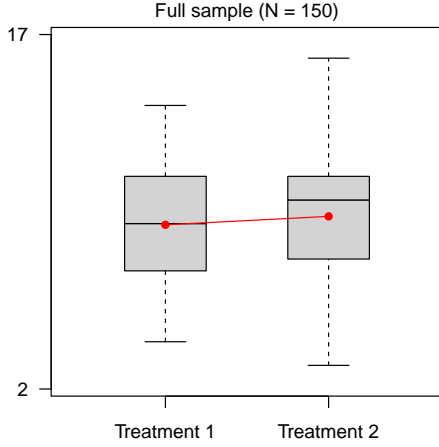[1]An overview of notation used is provided in the appendix.

*Figure 1.* 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 $\beta$.

1  covariates are included in the GLM. Furthermore, we assume a continuous response variable,
2  employing a Gaussian distribution with identity link and denote the error by $\epsilon_i = y_i - \mu_i$ with
3  variance $\sigma_\epsilon^2$. However, the model can easily accomodate additional treatment conditions and
4  covariates, and binary or count outcome variables.

5    The result of fitting a global GLM to the artificial motivating dataset is depicted in
6  Figure 1; the boxplots show the distribution of the posttreatment depression scores in both
7  treatment groups. The global model does not describe the data well: there is substantial
8  residual variance and the slope of the regression line is nearly zero. This does not necessarily
9  mean that posttreatment depression score and treatment type are unrelated, as the effect
10 of treatment may be moderated by variables not yet included in the model.

11    The MOB algorithm detects such moderation, by testing for parameter stability over
12 a set of auxiliary covariates or partitioning variables. When the partitioning is based on a
13 GLM, instabilities are differences in $\hat{\beta}$ across partitions of the dataset, which are defined
14 by one or more auxiliary covariates not included in the linear predictor. To find these
15 partitions, the MOB algorithm cycles iteratively through the following steps (Zeileis et al.,
16 2008): (1) fit the parametric model to the dataset, (2) test for parameter instability over
17 a set of partitioning variables, (3) if there is some overall parameter instability, split the
18 dataset with respect to the variable associated with the highest instability, (4) repeat the
19 procedure in each of the resulting subgroups.

1       More specifically, in step (2), to test for parameter instability, the so-called *scores* are
2 computed, using the score function. By definition, the empirical scores of all observations
3 in a dataset sum to zero. When the model is correctly specified, the expected value of the
4 score for each observation is also zero. Therefore, under the null hypothesis of parameter
5 stability, the scores do not systematically deviate from the expected value of zero, when the
6 observations are ordered by the values of a potential partitioning variable $U_k$ (c.f., Merkle
7 & Zeileis, 2013). To statistically test whether the scores systematically deviate from zero
8 with respect to variable $U_k$, the class of generalized M-fluctuation tests is used (Zeileis,
9 2005; Zeileis & Hornik, 2007).

10       If the null hypothesis of parameter stability in step (2) can be rejected, that is, if at
11 least one of the partitioning variables $U_k$ yields a p-value for the M-fluctuation test below
12 the pre-specified significance level $\alpha$, the dataset is partitioned into two subsets in step (3).
13 The binary partition is created using $U_{k*}$, the variable with the minimal p-value in step
14 (2). The split point for $U_{k*}$ is selected, by taking the value that minimizes the sum of the
15 values of the objective function in both partitions (Zeileis et al., 2008). In step (4), steps (1)
16 through (3) are repeated in each partition, until the null hypothesis of parameter stability
17 can no longer be rejected.

18       Due to the binary recursive nature of MOB, the resulting partition can be represented
19 as a binary tree. If the partitioning is based on a GLM, the result is a GLM tree, which
20 has a local fixed-effects regression model in every $j$-th ($j = 1, \ldots, J$) terminal node. As a
21 result, the value for $\beta$ depends on terminal node $j$ in which observation $i$ 'falls':

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

22 Note that, if the recursive subgroup structure (i.e., the partition) were known, the tree could
23 be estimated as a single GLM where all coefficients interact with the subgroup indicator.
24 Somewhat more formally, the model could then be written: $g(\mu_i) = x_i^{*\top} \beta^*$, where $x_i^*$ are
25 the values of the $2J$ interactions between the subgroups from the tree, and the elements of
26 $x_i$. $\beta^*$ would have length $2J$, and contain the subgroup-specific fixed-effects coefficients.

27       Figure 2 shows the GLM tree grown on the artificial motivating dataset. By using the
28 three auxiliary covariates (anxiety, duration and age), MOB partitioned the observations
29 into four subgroups, each with a different estimate for $\beta_j$. Age was correctly not detected
30 as a partitioning variable, and the left- and rightmost subgroups are in accordance with the
31 treatment-subgroup interactions as described above. However, the two subgroups in the
32 middle do not represent true subgroups, and this may be due to the clustered structure of
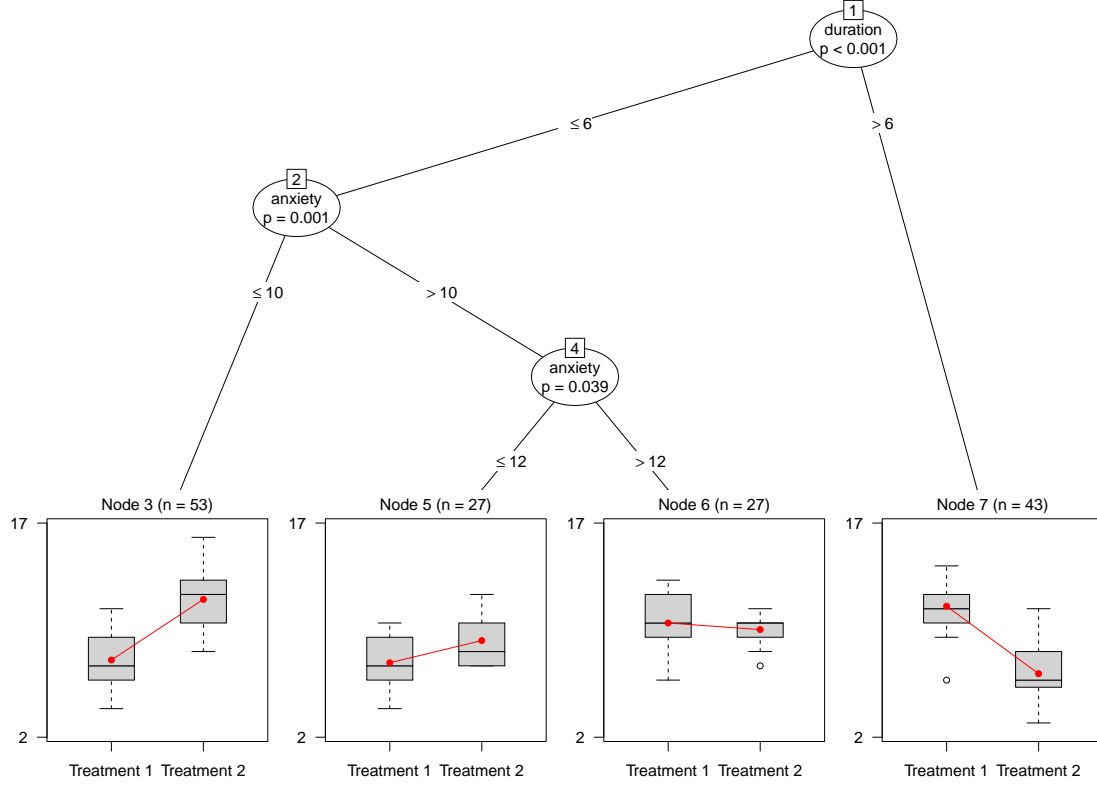33 the dataset not being taken into account.

*Figure 2.* 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 *Generalized linear mixed-effects model trees*

²       For datasets containing observations from multiple clusters (e.g., trials or research
³ centers), application of a GLMM would be more appropriate. The model in Equation 2 is
⁴ then extended to include cluster-specific, or random effects:

$$g(\mu_i) = x_i^\top \beta + z_i^\top b \tag{4}$$

⁵ Where $z_i$ is a unit vector of length $M$, of which the $m$-th element takes a value of 1, and all
⁶ other elements take a value of 0; $m$ $(m = 1, \ldots, M)$ denotes the cluster which observation $i$
⁷ is part of. Further, $b$ is a random vector of length $M$, with every element being the random
⁸ intercept for cluster $m$. Within the GLMM, it is assumed that $b$ is normally distributed,
⁹ with mean zero and variance $\sigma_b^2$. The parameters of the GLMM can be estimated with, for
¹⁰ example, maximum likelihood (ML) and restricted ML (REML).

1      For simplicity, we assume that only cluster-specific intercepts are included in the
2 model in the current paper. However, multiple random-effects coefficients can easily be
3 included.

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

7      Although GLMMs are adequately suited for taking into account the nested struc-
8 ture of a datasets, the global fixed-effects model in Equation 4 may not describe the data
9 well. Therefore, we propose the GLMM tree model, in which the fixed-effects part may be
10 partitioned as in Equation 3 and random effects are incorporated as well:

$$g(\mu_i) = x_i^\top \beta_j + z_i^\top b \tag{5}$$

11 To estimate the parameters of this model, we take an approach similar to that of the mixed-
12 effects regression tree (MERT) approach of Hajjem et al. (2011) and Sela and Simonoff
13 (2012). In the MERT approach, the fixed-effects part of a GLMM is replaced by a CART
14 tree with constant fits in the nodes, and the random-effects part is estimated as usual. To
15 estimate a MERT, an iterative approach is taken, alternating between (1) assuming random
16 effects known, allowing for estimation of the CART tree, and (2) assuming the CART tree
17 known, allowing for estimation of the random effects.

18      For estimating GLMM trees, we take this approach a step further: we take a GLM
19 tree, instead of CART tree with constant fits to estimate the fixed-effects part of the GLMM.
20 This allows not only for detection of differences in main effects, but also for detection of
21 differences in regression effects (e.g., of treatment type) across terminal nodes. In addition,
22 GLMM trees can be estimated for continuous, as well as binary and count variables. The
23 GLMM tree algorithm takes the following steps to estimate the model in Equation 5:

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

25 **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.

26 **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
27      $j(r)$ from the GLM tree.

28 **Step 3:** Repeat Steps 1 and 2 until convergence.

29      The algorithm initializes by setting $b$ to 0, since the random effects are initially un-
30 known. In every iteration, the GLM tree and coefficients $\beta_{j(r)}$ and $b_{(r)}$ are re-estimated.
31 The GLM tree is estimated, given the estimated random effects from the last iteration,
32 and the random effects are estimated, given the estimated GLM tree from the current it-
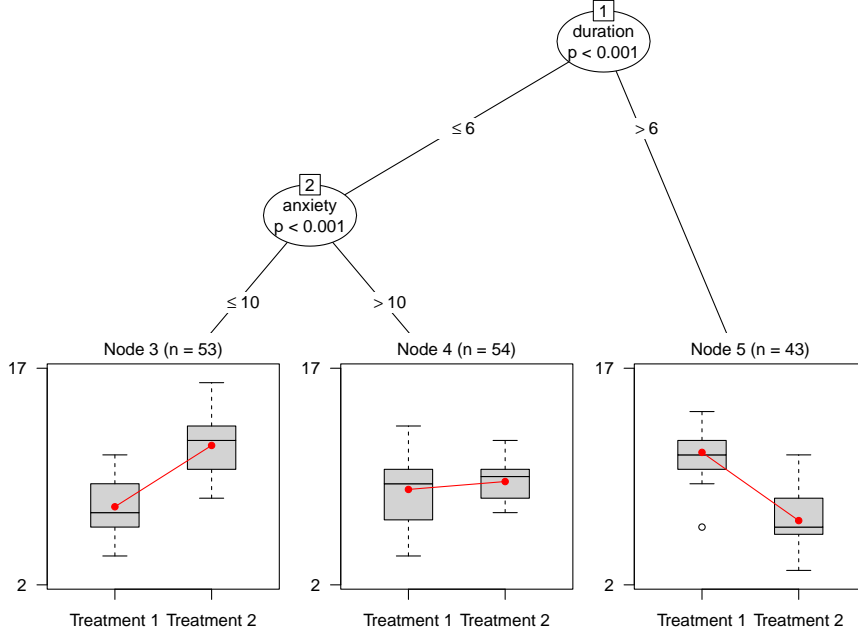
*Figure 3.* GLMM tree of the motivating example dataset. Three covariates (anxiety questionnaire score, duration of depressive symptoms at baseline in months and age) were used as potential splitting variables, and the clustering structure was taken into account by estimating random intercepts.

eration. Iterations are continued until convergence, which is monitored by computing the log-likelihood criterion of the mixed-effects model in Equation 5.

In Figure 3, the result of applying the GLMM tree algorithm to the artificial motivating dataset is presented. As can be seen, by taking into account the clustering of observations, the true treatment subgroups have been recovered, and the spurious split involving the anxiety variable no longer appears in the tree.


## Empirical Evaluation: Method


To assess and compare the performance of GLMM tree, we performed three simulation studies. In Study I, we compared the accuracy of GLMM tree with that of GLM tree in datasets with treatment-subgroup interactions. In Study II, we compared the accuracy of GLMM tree with that of GLM tree in datasets without treatment-subgroup interactions, allowing for assessment of the Type-I error rate. In Study III, we compared the accuracy of GLMM tree with that of a more classical approach of interaction detection: a GLMM with pre-specified interactions.

*Software*

R (R Core Team, 2014) was used for 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. For other response
distributions, the `glmtree` function would be available. For estimation of GLMMs the `lmer`
(or `glmer`, respectively) from the `lme4` package (version 1.1-7; Bates, Maechler, & Bolker,
2014) 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, 2015; available from R-Forge). This
package provides functions `lmertree` and `glmertree` that iterate between estimation of the
`lmtree`/`glmtree` model and the `lmer`/`glmer` model.

In all simulations, the significance level $\alpha$ for the parameter instability tests in the
trees was set to .05, with a Bonferroni correction applied for multiple testing. 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.

*Treatment-subgroup interaction design*

For generating datasets with treatment-subgroup interactions, we used one of the
treatment-subgroup interaction designs from Dusseldorp and Van Mechelen (2014), which
is depicted in Figure 4. Figure 4 shows two subgroups with mean differences in treatment
outcomes, and two subgroups without mean differences in treatment outcomes. The four
subgroups are characterized by their values on the partitioning variables $U_2$, and $U_1$ or $U_5$.
That is, $U_1$, $U_2$ and $U_5$ are true partitioning variables, whereas other potential partitioning
variables ($U_3$, $U_4$, $U_6$ through $U_{15}$) are noise variables.

*Data-generating parameters.* In generating datasets, we varied seven parameters of
the data-generating process:
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 $U_K$:
$K = 5$, $K = 15$ (where only $U_1$, $U_2$ and $U_5$ are true partitioning variables).
3. Two levels of intercorrelations between the 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 are drawn: $\sigma_b = 0$, $\sigma_b = 5$, $\sigma_b = 10$.

---

[2]Note that, when $\sigma_b = 0$, the correlation between $b$ and one of the $U_k$ variables is 0, by definition. However, datasets were created for this condition, to allow for a full factorial design of the simulation study; in reality, $b$ and $U$ are uncorrelated in these instances.
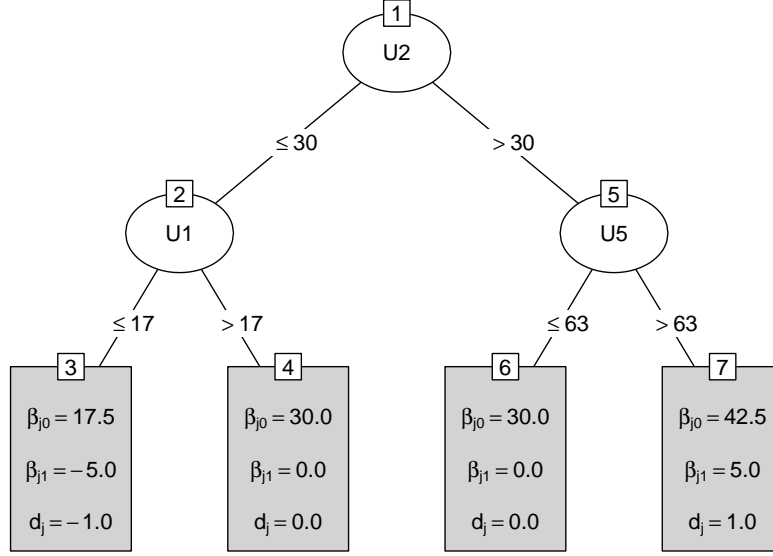
*Figure 4.* Data-generating model for treatment-subgroup interactions. Parameter $d$ denotes the standardized mean difference between the outcomes of Treatment 1 and 2 (i.e., $\beta_{j1}/\sigma_\epsilon$).

6. Three levels for the intercorrelations between $b$ and one of the $U_k$ variables: $b$ and $U_k$ uncorrelated, $b$ correlated with a true partitioning variable (i.e., $U_2$, $U_1$, or $U_5$, introducing a correlation of $\approx 0.42$), $b$ correlated with a non-partitioning covariate (i.e., $U_3$ or $U_4$, introducing a correlation of $\approx 0.42$)$^2$.

7. Two levels for $\beta_{j1}$, the unstandardized mean difference in treatment outcomes. In datasets with treatment-subgroup interactions, the treatment effect difference varied only in nodes 4 and 7, with levels $|\beta_{j1}| = 2.5$ (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 $d = 1.0$).

*Variable distributions*

Following the approach of Dusseldorp and Van Mechelen (2014), all covariates $U_1$ through $U_K$ were drawn from a multivariate normal distribution with means $\mu_{U1}$, $\mu_{U2}$, $\mu_{U4}$, and $\mu_{U5}$ fixed at 10, 30, $-40$, and 70, respectively. The means for all other covariates (i.e., $\mu_{U3}$, and $\mu_{U6}$ through $\mu_{U15}$) were drawn from a discrete uniform distribution on the interval $[-70, 70]$. All covariates $U_1$ through $U_{15}$ have the same standard deviation: $\sigma_{Uk} = 10$. Correlations between the $U_k$ variables vary according to the third facet of the simulation design described above.

To generate the random error term $\epsilon$, for every observation we drew a value from a

normal distribution with $\mu_\epsilon = 0$ and $\sigma_\epsilon = 5$.

To generate the cluster-specific intercepts $b_m$, we partitioned the sample into 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, we drew a single value $b_m$ from a normal distribution with mean 0 and the value of $\sigma_b$ 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.

To generate node-specific fixed effects, we partitioned the sample according to the terminal nodes of the tree in Figure 4.3. In combination with the seventh facet of the simulation design, this determines the values of $\beta_j$. For every observation, we generated a binomial variable (with probability .5) as an indicator for treatment type.

Finally, the response variable was calculated as the sum of the (node-specific) fixed effects, random effects and the error term: $y_i = x_i^\top \beta_j + z_i^\top b_m + \epsilon_i$.

*Assessment of performance*

The total number of nodes in estimated GLM and GLMM tree were calculated to assess Type-I error in datasets without treatment-subgroup interactions. In datasets with treatment-subgroup interactions, the number of nodes allows for assessment of the probability that a dataset is erroneously not partitioned, and the extent to which the algorithms may detect spurious subgroups.

In datasets with treatment-subgroup interactions, we assessed the accuracy of each of the estimated models. For GLM and GLMM trees, an accurately recovered tree was defined as a tree with (1) the true tree size (i.e., tree size = 7 in datasets with treatment-subgroup interactions), (2) the first split in the tree involving variable $U_2$ and a value of $30 \pm 5$, (3) the next split on the left involving variable $U_1$ and a value of $17 \pm 5$, and (4) the next split on the right involving variable $U_5$ and a value of $63 \pm 5$. Note that the allowance of $\pm 5$ equals an allowance of plus or minus half the population standard deviation of the partitioning variable ($\sigma_{U_k}$).

Predictive accuracy of each method was assessed by calculating the correlation between true and predicted treatment-effect differences ($\beta_{j1}$ in Figure 4) in each dataset. Consequently, predictive accuracy was only be assessed in datasets with treatment-subgroup interactions, as the true treatment differences have a constant value in datasets without treatment-subgroup interactions. 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.
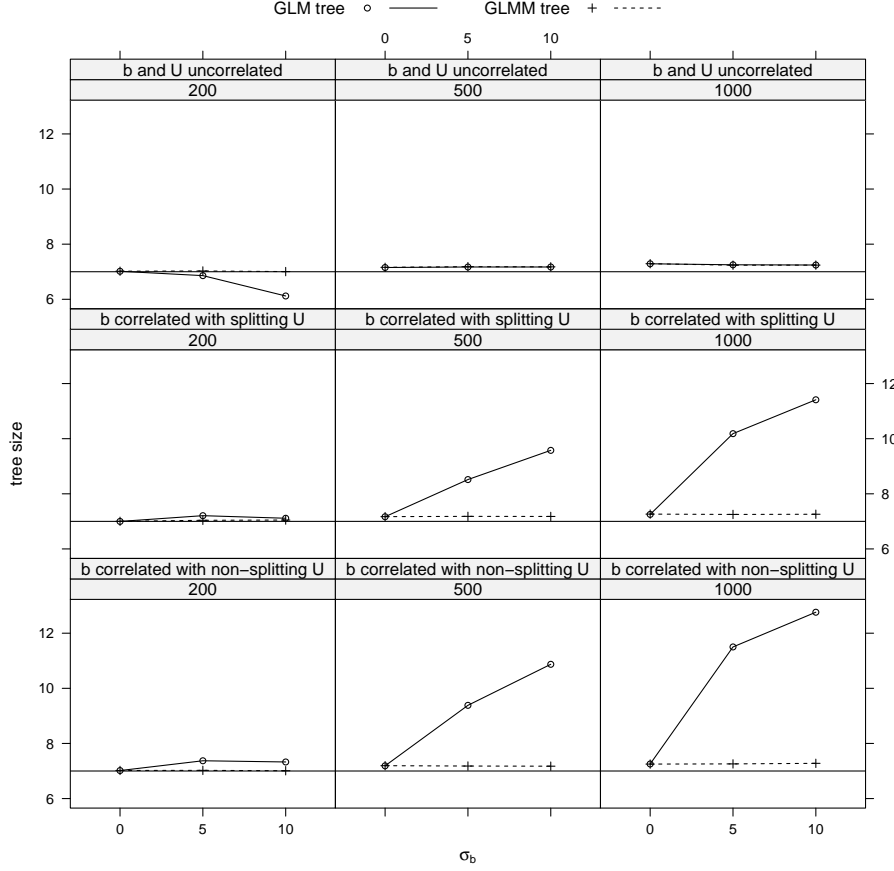
*Figure 5.* Average tree size of GLM and GLMM trees for datasets with treatment-subgroup interactions. Rows represent levels of dependence between random effects ($b$) and one of the partitioning variables $U_k$; columns represent levels of sample size.

# Empirical Evaluation: Results

*Study I: Piecewise interactions*

For each cell of the data-generating design described above, 50 datasets were generated, resulting in $50 \times 3 \times 2 \times 2 \times 3 \times 3 \times 3 \times 2 = 32{,}400$ training datasets with treatment-subgroup interactions. In these datasets, the true tree size was 7 (4 terminal nodes and 3 inner nodes; Figure 4).

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

both algorithms.

To find the most important predictors of tree size, we performed an ANOVA 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 of tree size were algorithm type, sample size, variance of the random intercept, and the correlation between partitioning variables. These effects were further assessed by means of a graphical display (Figure 5).

Figure 5 indicates that average tree size increases slightly with sample size, which may be due to increased power. In the absence of random effects (i.e., $\sigma_b = 0$), both GLM and GLMM trees have about 7 nodes on average. However, clear differences in tree size were observed when $\sigma_b > 0$. When sample size is small and random intercepts are not correlated with a partitioning covariate, GLM tree lacks power and grows trees that are too small, on average. When random intercepts are not correlated with a partitioning covariate and sample size is large (i.e., 500 or 1000), GLM and GLMM trees are about equally sized and have approximately the true tree size, on average. When random intercepts are correlated with a partitioning covariate, GLM starts to create spurious splits, especially when sample size is larger (i.e., 500 or 1000) and when $\sigma_b$ is large (i.e., 10). This effect was somewhat stronger when the random intercept was correlated with a covariate that did not appear in the true tree.

*Accuracy of recovered trees.* To assess the accuracy of trees recovered by both algorithms, we inspected the variables and values that were selected for partitioning. For the first split, GLMM tree selected the true partitioning variable ($U_2$) in all datasets, and GLM tree in all but one datasets. The splitting value for $U_2$ selected for the first split was 29.94 for both GLM and GLMM tree, which is very close to the true splitting value of 30 (Figure 4). GLM tree did show somewhat higher variability in recovering the splitting value for the first split than GLMM tree (SD = 0.154 and SD = 0.127, respectively).

GLMM tree accurately recovered the partition in 90.40% of datasets, whereas GLM tree accurately recovered the partition in only 61.44% of datasets. To find the most important predictors of tree accuracy, we fitted a GLM 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 of tree accuracy were algorithm type, sample size, variance of the random intercept, and the level of dependence between the partitioning variables and the random intercept. These effects were further assessed by means of a graphical display (Figure 6).

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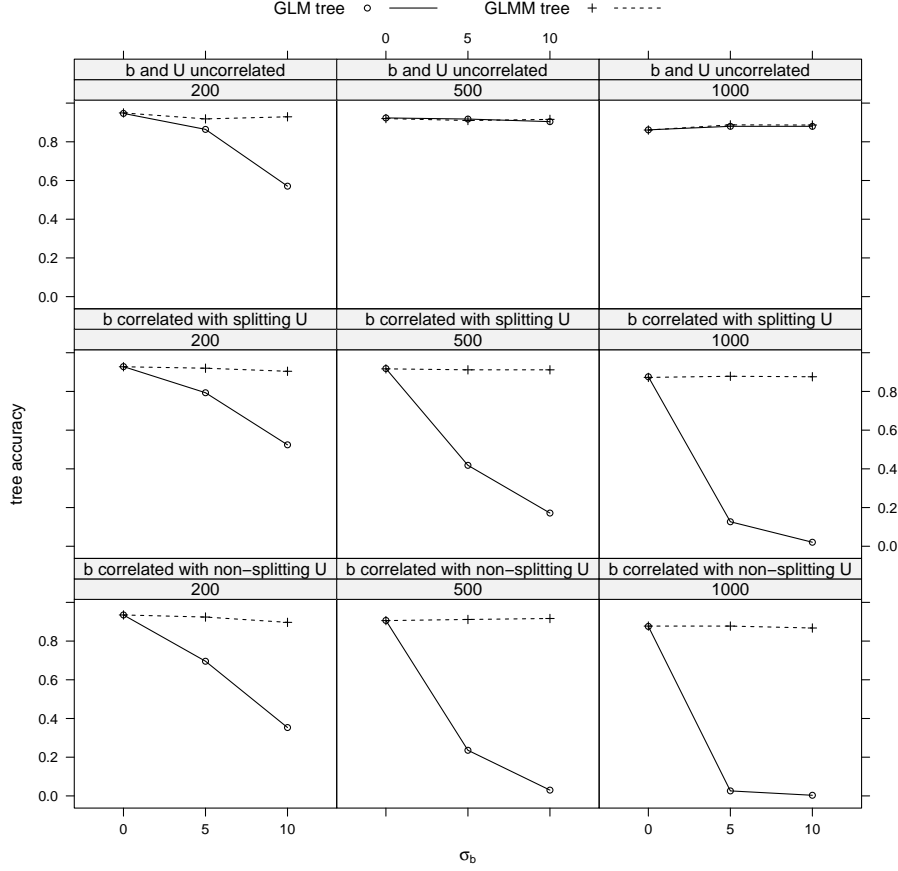
*Figure 6.* Tree accuracy of GLM and GLMM tree. Tree accuracy is defined as the proportion of datasets in which the true tree was accurately recovered. Rows represent levels of dependence between random effects ($b$) and one of the partitioning variables $U_k$, columns represent levels of sample size.

trees were much less accurate than GLMM trees. This was found for all sample sizes, when random effects were correlated with a partitioning covariate. This effect was somewhat stronger when the correlated $U_k$ was not a true partitioning variable. When random intercepts were not correlated with one of the $U_k$ variables, GLMM tree clearly outperformed GLM tree only when sample size was small (i.e., $N = 200$).

*Predictive accuracy of trees.* To assess predictive accuracy, correlations between the true and predicted treatment-effect differences of both algorithms were calculated, using test datasets. 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.
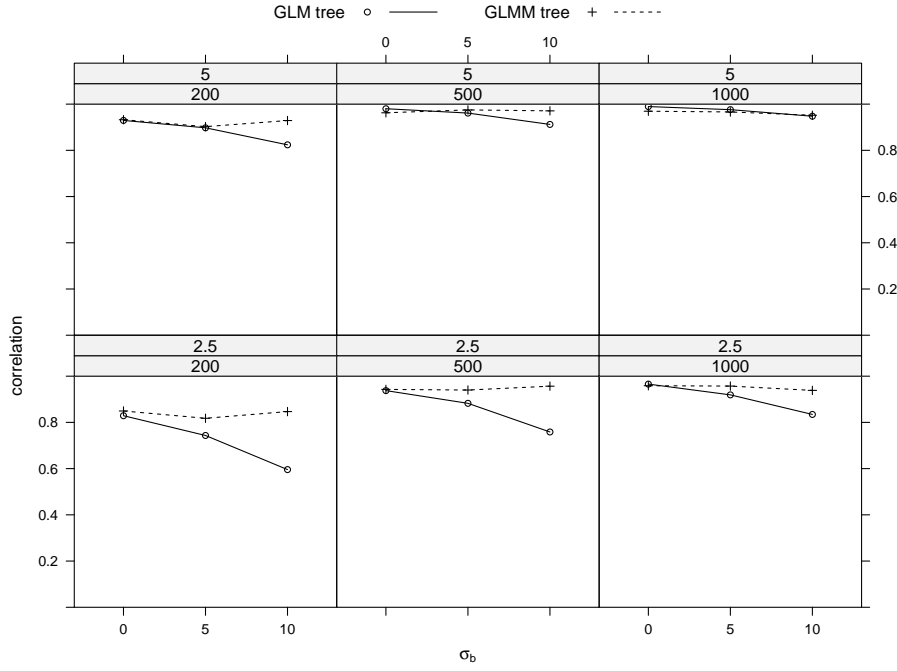
*Figure 7.* Average predictive accuracy of GLM and GLMM trees. Predictive accuracy of trees is defined as the correlation between the true and predicted differences between Treatment 1 and 2. Rows represent the levels of the absolute value of the unstandardized treatment-effect difference in subgroups with treatment-effect differences, columns represent levels of sample size.

An ANOVA indicated that the most important predictors of predictive accuracy were algorithm type, sample size, variance of the random intercept, and the treatment difference effect size. These effects were further assessed by means of a graphical display (Figure 7).

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. However, it should be noted that, in the absence of random effects, GLMM tree tends to outperform GLM tree in smaller sample sizes (i.e., 200), whereas GLM tree tends to outperform GLMM tree in larger sample sizes (i.e., 500 and 1000). In the presence of random effects, GLM tree predictions were much 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 addition, with increasing sample size and treatment-effect differences, the difference in predictive accuracy between GLM and GLMM tree becomes smaller.

*Study II: Type-I error*

*Design.* For assessing Type-I error of GLM and GLMM tree, we generated datasets without treatment-subgroup interactions, in which there is only a main effect of treatment in the population. Put differently, the number of subgroups or terminal nodes in these datasets was $J = 1$, and there was only a single value of $\beta_j = \beta$ in every dataset. For the datasets without treatment-subgroup interactions, we varied the same parameters of the data-generating process as in Study I. However, the sixth facet of the data-generating process had only two levels for these datasets, as there were no 'true' partitioning variables in these datasets. Therefore, $50 \times 3 \times 2 \times 2 \times 3 \times 3 \times 2 \times 2 = 21{,}600$ datasets without treatment-subgroup interactions were generated.

*Tree size and Type-I error.* In datasets without treatment-subgroup interactions, the average tree size was 1.09 (SD = 0.44) for GLMM tree, and 2.02 (SD = 1.68) for GLM tree. The overall Type-I error rate was only .04 for GLMM tree, and .33 for GLM tree. Again, main predictors of tree size were found to be sample size, variance of the random intercepts, and dependence between the random intercept and one of the partitioning variables. These effects were further assessed by means of a graphical display (Figure 8). Figure 8 indicates that in the absence of random effects (i.e., $\sigma_b = 0$), both GLM and GLMM tree tend to create trees of size 1. In the presence of random effects, GLMM tree still tends to create trees of size 1, but GLM tree creates much larger trees, when the random effects are correlated with one of the partitioning covariates. This effect increases with sample size.

*Study III: Linear and piecewise interactions*

*Design.* To compare the performance of GLMM tree with that of a GLMMs with pre-specified interaction effects, we created datasets in which interactions were varied from purely linear to purely piecewise. We reduced the number of cells by limiting the fourth facet of the data-generating design to only one level ($M = 25$ clusters). The fifth facet of the data-generating design was limited to two levels ($\sigma_b = 2.5$ and $\sigma_b = 7.5$). As in Study II, we limited the sixth facet to two levels (the random intercept was either correlated to one of the partitioning variables, or not). To counter unidentifiability due to potentially large numbers of pre-specified interaction terms in the GLMMs, we limited the number of covariates to $k = 5$ and $k = 10$ (instead of $k = 5$ and $k = 15$). Finally, to allow for a smooth transition from linear to piecewise interactions, the seventh facet of the interaction design (effect size of treatment differences) was replaced by a factor with three levels, controlling the type of interactions: purely linear interactions, purely piecewise interactions, and a combination of both. Therefore, $50 \times 3 \times 2 \times 2 \times 1 \times 2 \times 2 \times 3 = 7{,}200$ training datasets were generated for this study.

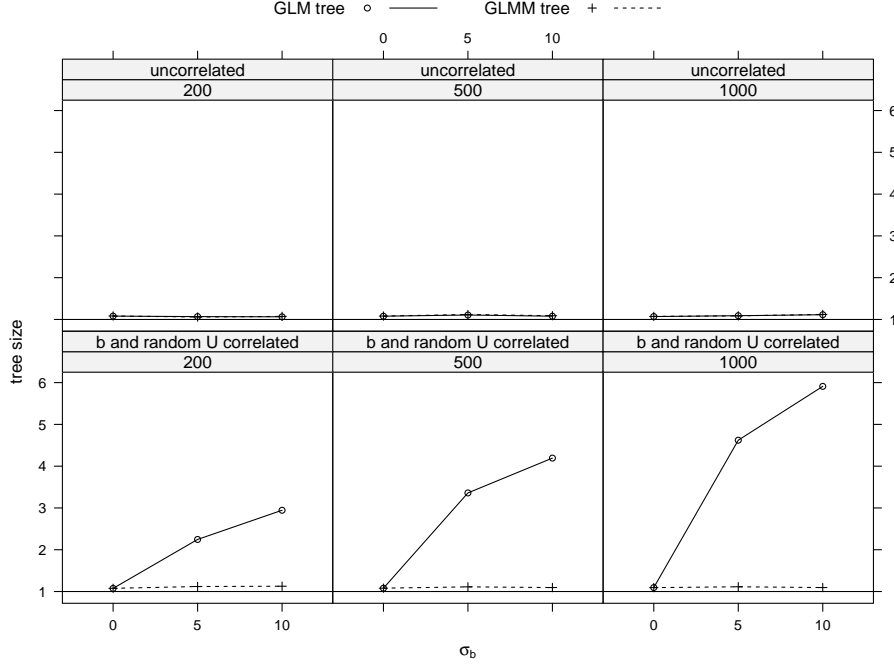The input variables were generated as described above. To generate purely linear

*Figure 8.* Average tree size of GLM and GLMM trees for datasets without treatment-subgroup interactions. Rows represent dependence between random effects ($b$) and one of the partitioning variables $U_k$; columns represent sample size.

interactions, outcome variable $y$ was generated as the sum of:

1. a main effect of $U_2$ (coefficient $= 1.5$)

2. an interaction between $U_1$ and $U_2$ with a negative effect (coefficient $= -.25$), and an interaction between $U_2$ and $U_5$ with a positive effect (coefficient $= .25$)

3. an interaction between $U_1$, $U_2$ and treatment with a negative effect (coefficient $= -.25$), and an interaction between $U_2$, $U_5$ and the treatment indicator with a positive effect (coefficient $= .25$)

Note that the direction and type (main or interaction) of these effects correspond to those in Figure 4. All interactions were created by using centered $U_k$ variables, which were calculated by subtracting the variable means. To generate purely piecewise interactions, we used the partition depicted in Figure 4. We estimated $\bar{y}$ for each of the treatment-subgroups by calculating the treatment-subgroup means in a large dataset ($N = 10^6$) with purely linear interactions. For each observation in a dataset with purely piecewise interactions, the treatment-subgroup specific value for $\bar{y}$ was imputed. To generate combined linear and piecewise interactions, $y$ values according to the purely piecewise and linear models were summed, and divided by 2.

GLMMs were estimated by specifying main effects for all covariates $U_k$ and the treatment indicator, first-order (or two-way) interactions between all pairs of covariates $U_k$, and second-order (or three-way) interactions between all pairs of covariates $U_k$ and treatment. Accurately recovered GLMMs were defined to be models in which significant (i.e., absolute value of the $t$-test statistic $\geq 1.96$) effects were found for $U_2$ and (treatment) interactions of $U_1$ and $U_2$ and $U_2$ and $U_5$. Accurately recovered GLMM trees were defined to be trees in which the first split involved $U_2$, the second split to the right involved $U_1$ and the third split to the right involved $U_5$. Note that this criterion of accuracy is somewhat more lenient than in Study I, as accurate GLMM trees were allowed to include further splits, and accurate GLMMs were allowed to include significant interaction and/or main effects of other variables.

*Accuracy of recovered models.* Overall, GLMM tree recovered the main and interaction effects in 53.58% of datasets, wheareas GLMM recovered the main and interaction effects in 34.81% of datasets. To find the most important predictors of model accuracy, we performed an ANOVA, in which the algorithm type, sample size, the number of covariates and type of interaction (linear, piecewise or both) were found to be most important. These effects were further assessed by means of the graphical display depicted in Figure 9.

Figure 9 indicates a clear effect of sample size: when the number of covariates is large (i.e., 10) GLMMs do not recover all true effects in the dataset, whereas they always seem to recover the true effects when the number of covariates is small (i.e., 5), and the effects are partly linear. However, when the effects are purely piecewise, GLMMs never recover the true effects in the dataset, unless sample size is very large (i.e., $N = 1000$). At the same time, Figure 9 indicates that the accuracy of GLMM trees is unaffected by the number of covariates. Furthermore, GLMM tree performs best in recovering purely piecewise interactions, and poorest in recovering purely linear interactions. Finally, the performance of GLMM tree is affected by sample size: generally, the true effects are more often recovered in datasets with larger sample sizes.

*Predictive accuracy.* Both algorithms showed equal average correlations between the true and predicted treatment-effect differences: average accuracy was .60 (SD = .36) for GLMMs and .60 (SD = .32) for GLMM trees. In terms of bias, GLMM treatment-effect difference predictions deviated more from the true values than the predictions of GLMM tree. The mean true treatment-effect difference averaged over all datasets was .04 (SD = 1.22), whereas the mean predicted treatment-effect differences averaged over all datasets for GLMMs was -.26 (SD = 100.33) and for GLMM trees .03 (SD = 3.88). Furthermore, it should be noted that the full models were used for prediction with GLMMs and GLMM trees. That is, predictions for GLMMs in datasets with 10 covariates may involve 11 (significant or non-significant) main effects, 100 first-order interactions and 100 second-order
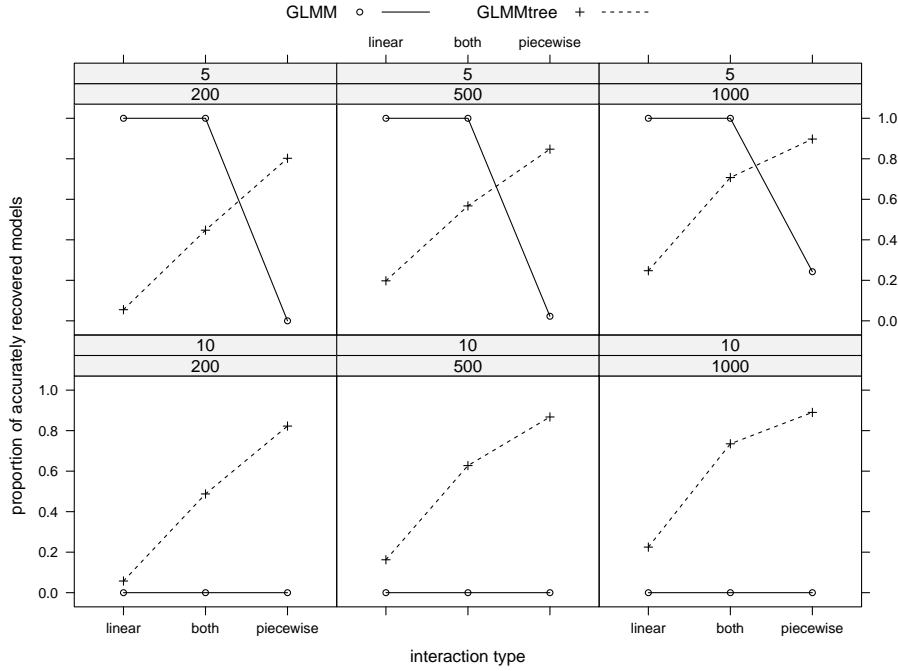
*Figure 9.* Model accuracy for GLMMs and GLMM trees. In an accurate model, the true main and interaction effects were accurately recovered. Columns represent sample size, rows represent the number of covariates in the model.

1 interactions. As the maximum tree depth of GLMM trees was set to 4, GLMM trees used
2 a maximum of only 8 terminal nodes for predictions.

3      The most important predictors of predictive accuracy were algorithm type, sample
4 size, the number of covariates, and type of interaction (linear, piecewise or both). These
5 effects were further assessed by means of the graphical display depicted in Figure 10. As
6 Figure 10 indicates, GLMM trees show highest predictive accuracy in datasets with purely
7 piecewise interactions, whereas GLMMs show highest predictive accuracy in datasets with
8 purely linear interactions. The predictive accuracy of GLMM tree is hardly affected by
9 the number of covariates, whereas the predictive accuracy of GLMMs decreases when the
10 number of covariates increases. Finally, for both algorithms, predictive accuracy improves
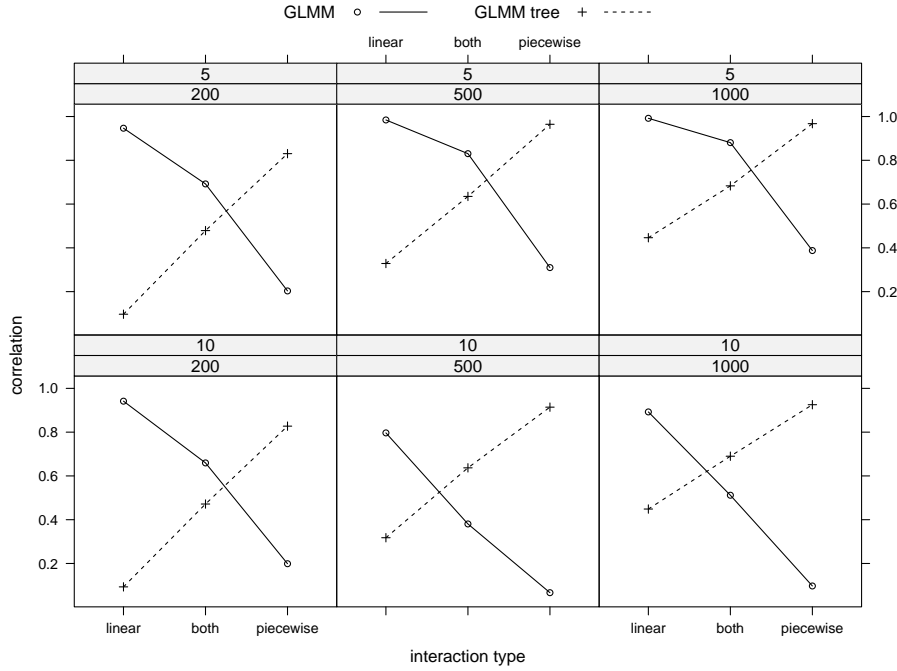11 somewhat when sample size increases.

*Figure 10.* 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 and compare the application of GLM and GLMM tree, we applied both algorithms to a dataset from a 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

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

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

12      To provide a standardized estimate of the treatment effect differences in the final
13 nodes of the trees, we calculated node-specific Cohen's $d$ values. Cohen's $d$ was calculated
14 by dividing the node-specific predicted treatment outcome difference by the node-specific
15 pooled standard deviation. Confidence intervals for Cohen's $d$ could not be calculated,
16 as these can not take into account the exploratory nature of our analyses (i.e., variable
17 and split point selection). The predictive accuracy of GLM and GLMM trees was assessed
18 by calculating average correlations between observed and predicted HAM-D posttreatment
19 scores, based on 50-fold cross validation.

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

*Results*

27      The tree and effects sizes resulting from application of GLM tree are presented in
28 Figures 11 and 12 and those resulting from application of GLMM tree are presented in
29 Figures 13 and 14.

30      The GLM tree (Figure 11) selected level of education as the first partitioning variable,
31 and presence of a comorbid anxiety disorder as a second partitioning variable, for observa-
32 tions with a higher level of education. By taking into account study-specific intercepts, the
33 GLMM tree (Figure 13) indicates that the first split made by GLM tree is a spurious one.
34 The GLMM tree selected presence of a comorbid anxiety disorder as the only partitioning
35 variable. The terminal nodes of Figure 13 show only a single treatment-subgroup interac-
36 tion: for patients without a comorbid anxiety disorder, CBT and PHA provide more or less
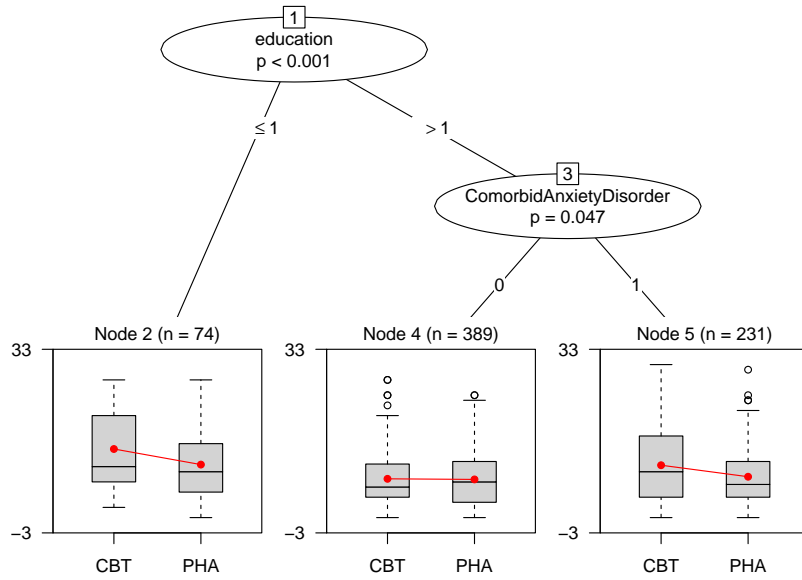37 the same reduction in HAM-D scores (Cohen's $d = 0.05$; Figure 14). For patients with a

*Figure 11.* GLM tree for prediction of posttreatment total scores on the Hamilton Rating Scale for Depression (HAM-D). 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).
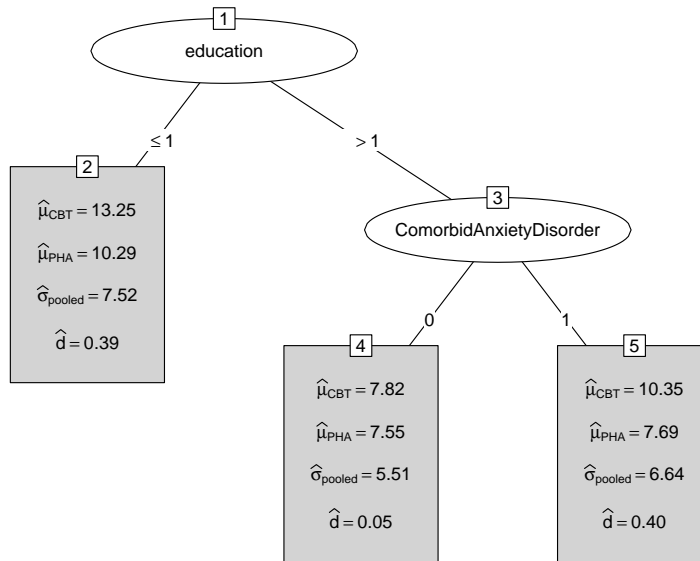


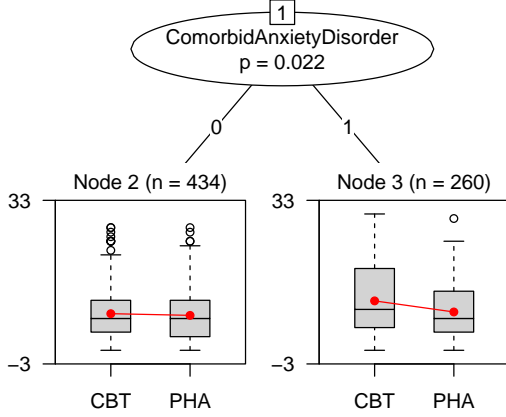*Figure 12.* Statistics and effect sizes for GLM tree in Figure 11.

*Figure 13.* GLMM tree for prediction of posttreatment total scores on the Hamilton Rating Scale for Depression (HAM-D). 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).

comorbid anxiety disorder, PHA provides a greater reduction in HAM-D scores (Cohen's $d = 0.39$; Figure 14). The estimated intraclass correlation coefficient for the GLMM tree was .05.

Assessment of predictive accuracy by means of 50-fold cross validation indicated better predictive accuracy for GLMM tree than GLM tree. The correlation between true and predicted posttreatment HAM-D total scores, averaged over the 50 folds, was .28 (*var =* .067) for GLMM tree, and .19 (*var =* .084) for GLM tree. This indicates that GLMM tree provided higher predictive accuracy, and also somewhat lower variability of predictive accuracy than GLM tree.

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

# 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
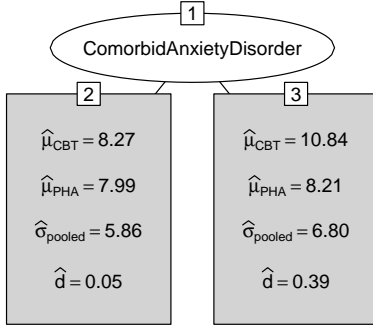
*Figure 14.* Statistics and effect sizes for GLM tree in Figure 13.

Table 1: Variable selection statistics

| | Selection frequency | |
| --- | --- | --- |
| Variable | GLM tree | GLMM tree |
| Education | 0.956 | 0.014 |
| ComorbidAnxietyDisorder | 0.398 | 0.528 |
| HRSDt0 | 0.034 | 0.002 |
| Age | 0.000 | 0.022 |
| Gender | 0.002 | 0.004 |

*Note.* Frequencies are calculated over 500 random subsamples of size $.9 \times N$ of the complete dataset. Frequencies do not add up to 1, as trees may involve multiple, or no splits.

treatment-subgroup interactions in clustered datasets, which was confirmed by our simulation studies.

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

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

teractions.

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

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

Compared to treatment-interaction detection by means of GLMMs with pre-specified interaction effects, GLMM trees provided similar predictive accuracy, on average. When interactions were purely linear, GLMMs outperformed GLMM trees, and when interactions were purely piecewise, GLMM trees outperformed GLMMs. However, GLMM tree may have a clear advantage when there are a large number of potential moderator variables (i.e., $> 5$). With 10 potential moderator variables, $t$-tests of GLMMs were unable to identify the true moderator variables, whereas the number of potential moderator variables did not influence performance of GLMM tree. Furthermore, the number of terms in a GLMM increases quadratically with the number of potential moderator variables, yielding complex predictive models. The trees in our simulations were limited to a maximum number of 8 terminal nodes, and may therefore be much easier to use for prediction in practical decision-making contexts.

In the Application, we found GLMM tree to provide results that are easily interpretable, and also more accurate than a GLM tree without random effects. In addition, node-specific means and variances can be used to calculate effect sizes, to judge clinical relevance of the findings, 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.

These findings are encouraging for the use of GLMM tree in the detection of

treatment-subgroup interactions in datasets with clustered structures. However, it should be noted that the simulations show that GLMM tree performs very well, if the model is correctly specified model. That is, if there are subgroups with respect to the partitioning variables specified, so that there are different parameters of the GLM in each of these subgroups, then GLMM tree will accurately recover those subgroups. However, as with any data-analytic method, misspecification of the model will negatively affect performance. One source of misspecification would be the omission of relevant variables in the GLM, or as potential partitioning variables. Another source of misspecification would be the inclusion of irrelevant variables, either in the linear predictor of the GLM or as partitioning variables, which may reduce the power to detect the actual subgroups. However, the performance of GLMM tree, in contrast to that of GLMMs with pre-specified interactions, was not negatively affected by the number of potential moderator variables.

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

## References

Bates, D., Maechler, M., & Bolker, B. (2014). *lme4: Linear mixed-effects models using S4 classes.* Retrieved from `http://CRAN.R-project.org/package=lme4` (R package version 1.1-7)

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., & 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. (2015). *glmertree: Generalized linear mixed model trees.* Retrieved from `http://R-Forge.R-project.org/R/?group_id=261` (R package version 0.1-0)

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). partykit: A modular toolkit for recursive partitioning in R. *Journal of Machine Learning Research*. (Forthcoming, preprint at http://EconPapers.RePEc.org/RePEc:inn:wpaper:2014-10)

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., & Carolin, S. (2016). A toolkit for stability assessment of tree-based learners. *Working Papers in Economics and Statistics, University of Innsbruck*.

R Core Team. (2014). R: A language and environment for statistical computing [Computer software manual]. Vienna, Austria. Retrieved from `http://www.R-project.org/`

Seibold, H., Zeileis, A., & Hothorn, T. (2015). Model-based recursive partitioning for subgroup analyses. *International Journal of Biostatistics*.

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.