

Contents lists available at boston-science-publishing.us

International Journal of Medical and Clinical Case Reports



Identifying Patient Subgroups for Personalized Treatments with Model Based Recursive Partitioning

Habeeb Abolaji Bashir^{*1}, George Paul Komolafe², Deborah Idowu Akinwolemiwa³

¹Department of Statistics and Data Science, University of Kentucky, Kentucky, USA, ORCID: 0009-0008-2881-2154

²Department of Computer Science, Boston University, Massachusetts, USA, ORCID: 0009-0001-0413-241X

³Department of Economics, Wayne State University, Detroit, USA, ORCID: 0009-0000-1045-4506

ARTICLE INFO

Article history:

Received 14 December 2023

Revised 20 February 2024

Accepted 07 March 2024

Published 02 September 2024

KEYWORDS:

Personalized Medicine
Patient stratification
Model based recursive partitioning
Treatment heterogeneity
Subgroup analysis
Decision trees
Precision medicine

ABSTRACT

Advances in precision medicine have highlighted the inadequacy of traditional “one size fits all” treatment approaches in the face of patient heterogeneity. Different patients often respond variably to the same therapy due to genetic, clinical, or environmental factors. This heterogeneity poses a challenge in identifying which subgroups of patients benefit from a given treatment. We introduce model based recursive partitioning (MOB) as a statistical approach to automatically detect patient subgroups with differential treatment effects, addressing gaps in current subgroup analysis methods. This study uses **model based recursive partitioning (MOB)** to identify patient subgroups with differential treatment effects in a simulated randomized controlled trial. A dataset of 600 patients, each with baseline demographics, a biomarker and a severity score, was generated; patients were randomly assigned in approximately equal numbers to treatment and control (the final counts were 318 and 282 due to chance). Linear regression models were embedded into a recursive partitioning algorithm to detect effect modifiers. Model performance was compared against global linear and interaction models as well as a standard CART tree. MOB correctly recovered the programmed treatment effect heterogeneity and produced an easily interpretable decision tree. Cross validated analyses showed that accounting for heterogeneity improved predictive performance relative to a simple global model. The resulting subgroup rules could guide personalized treatment strategies and inform future trial designs.

© Habeeb Abolaji Bashir, et al. This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International License, which permits unrestricted use, distribution and reproduction in any medium, provided the original author and source are credited.

Introduction

Personalized medicine has emerged in response to the recognition that patients vary widely in their responses to treatment, rendering “**one size fits all**” approaches inadequate. Traditional therapeutic strategies that assume a uniform treatment effect can lead to suboptimal outcomes: some patients experience significant benefit, others negligible effect, and some may even be harmed, all under the same intervention. For example, a landmark pharmacogenetic analysis demonstrated that **colorectal cancer patients with a KRAS gene mutation do not benefit from the drug cetuximab, whereas KRAS wild type patients do**, highlighting a clear subgroup with differential treatment efficacy. Such evidence underlines the critical **role of patient heterogeneity** in treatment response. Factors like genetics, biomarkers, disease subtype, comorbidities, and lifestyle can all modulate how a patient responds to a given therapy. As a result, modern clinical research and practice are increasingly focused on **patient stratification** dividing patients into subgroups based on predictive characteristics to tailor treatments to those most likely to benefit.

Despite this promise, implementing personalized treatments remains challenging. Conventional clinical trials and analyses often rely on **average treatment effects**, which may obscure important subgroup differences. Standard subgroup analyses in trials typically examine one variable at a time (e.g., age, sex) for effect modification, but this approach is statistically

underpowered and can be misleading. It fails to capture interactions among multiple factors and raises multiplicity issues, increasing the risk of false findings. Consequently, there is a pressing need for robust data driven methods to **identify patient subgroups with heterogeneous treatment effects** in a principled manner. Recognizing this, regulatory agencies such as the EMA have highlighted both the potential and pitfalls of subgroup analyses and have **called for new statistical procedures** to aid data driven subgroup identification while controlling errors. In this context, tree based algorithms have gained attention as intuitive tools for stratifying populations.

Recursive partitioning methods (decision tree algorithms) provide a natural framework for subgroup discovery by hierarchically splitting patients based on covariates. Classic techniques like Classification and Regression Trees (CART) (Breiman et al., 1984) have long been used in healthcare analytics for risk stratification and outcome prediction. CART can **partition patients into subgroups** with distinct outcomes or risks, and has been applied to identify, for instance, subsets of patients at highest risk of complications. Random forests, an ensemble of CART trees, and related machine learning approaches (e.g., gradient boosting) further improve predictive accuracy by averaging many trees. However, these methods typically prioritize prediction over interpretability; they may not directly identify *why* a subgroup differs, and the resulting models (especially ensembles) can be complex black boxes. Moreover, a standard CART tree grown to predict outcomes might not **focus on treatment effect heterogeneity**, it tends to split on variables that reduce overall outcome variance, which could lead to ignoring interaction effects that are subtle in

* Corresponding author.

Habeeb Abolaji Bashir, Department of Statistics and Data Science, University of Kentucky, Kentucky, USA, ORCID: 0009-0008-2881-2154.

the full population (as we demonstrate later with CART's performance). In fact, naive application of CART can miss subgroups with differential treatment response because outcome variation due to prognostic factors can mask the predictive interaction effects of interest.

To specifically target **subgroup differences in treatment effect**, specialized algorithms have been developed. One approach is to incorporate interaction terms in regression models (e.g., testing each covariate \times treatment interaction), but with many covariates, this becomes unwieldy and prone to overfitting or missing higher order interactions. **Model based recursive partitioning (MOB)** offers an elegant extension: it embeds recursive partitioning within a statistical model, allowing the detection of subgroups with different model parameters. Introduced by Zeileis, Hothorn, and Hornik (2008), MOB starts with a prespecified parametric model (such as a linear regression for outcome vs. treatment) and uses hypothesis tests to assess **parameter instability** across covariates. If a model parameter, notably the treatment effect varies significantly with respect to some covariate, the algorithm splits the data on that covariate to form more homogeneous subgroups. This process yields a **segmented model** where each terminal node (subgroup) has its own parameter estimates (e.g., its own treatment effect size), linked to covariate defined rules in a decision tree structure. In essence, MOB combines the interpretability of decision trees with the rigor of statistical models, focusing on **detecting interactions (effect modifiers)** in a data driven way. It extends classical tree methods by allowing **models in each node rather than simple averages**, which can greatly reduce tree size and improve interpretability when linear relationships exist within subgroups.

Research gap: While various methods for subgroup analysis exist from heuristic tree searches to modern machine learning (e.g., the SIDES algorithm, interaction trees, and causal forests (Wager & Athey, 2018)) there is a lack of consensus and comparative evaluation in clinical contexts. Many machine learning approaches can find heterogeneity but at the cost of interpretability or requiring large sample sizes. On the other hand, clinical stakeholders need transparent, explainable subgroup findings to trust and adopt them in practice. **Model based recursive partitioning** stands out by producing an easily interpretable decision tree and leveraging statistical inference to control false discoveries (e.g., through significance tests for splits). However, its application in biomedical research is still emerging. Previous studies have shown its utility for example, Seibold et al. (2016) applied MOB to identify ALS patient subgroups with differential response to a therapy, but more demonstrations are needed to illustrate how MOB can improve upon traditional methods in uncovering actionable subgroups.

Study objectives: This research aims to demonstrate and evaluate MOB for identifying patient subgroups with different treatment responses, using a realistic simulated dataset. We seek to answer: (1) Can MOB correctly recover known subgroups associated with treatment effect heterogeneity? (2) How does MOB's performance (in terms of model fit and subgroup detection) compare to standard approaches like CART or regression models with interaction terms? (3) Are the resulting subgroups clinically interpretable and plausible, supporting the goal of personalized treatment? By addressing these questions, we illustrate MOB's potential advantages and limitations in the context of precision medicine. Ultimately, this work highlights how an MOB based analysis can move us toward more **individualized treatment strategies**, improving clinical decision-making and trial design by accounting for patient heterogeneity.

Literature Review

Personalized Medicine and Patient Stratification

The paradigm of **personalized (precision) medicine** has shifted the healthcare model from treating the "average" patient to tailoring interventions for **specific patient subpopulations**. This shift is driven by evidence of pervasive heterogeneity in disease risk and treatment response across individuals. For example, genomic and biomarker research has revealed that patients carrying certain genetic mutations respond dramatically differently to targeted therapies, whereas those without the biomarker derive little benefit. As a result, patient stratification, grouping patients by relevant characteristics (genetic markers, phenotype, etc.) has become central in clinical research and drug development. Stratified medicine approaches have led to notable successes, especially in oncology: therapies like trastuzumab (Herceptin) are effective only in the subgroup of breast cancer patients whose tumors overexpress the HER2 protein, and are now standard care in that subgroup, while being ineffective (and hence not used) in others (Collins & Varmus, 2015). These examples

underscore the principle that **identifying the right treatment for the right subgroup** can greatly improve outcomes and avoid unnecessary toxicity.

However, discovering such subgroups is challenging. Traditionally, clinical trials perform **subgroup analyses** on pre-specified patient subsets (e.g., by age, sex, disease stage). These analyses are often underpowered and considered exploratory; indeed, a guideline from regulatory authorities cautions that most post hoc subgroup findings are false positives or not reproducible. Nonetheless, regulators and stakeholders recognize the value of rigorously exploring heterogeneity. The National Academy of Medicine (2019) emphasized developing better methods to assess heterogeneous treatment effects (HTE) and warned that conventional one variable at a time subgroup tests are of "*extremely limited value*" for guiding care. Instead, multi factor predictive approaches or risk scores are recommended to separate patients likely to benefit from those likely to be harmed or non responders. In recent years, **N of 1 trials** and other adaptive designs have been discussed as extreme cases of personalization, treating each patient as their own stratum. While not always feasible, these approaches highlight the motivation to move beyond population averages in evidence based medicine.

Tree Based Methods in Healthcare Analytics

Decision tree learning has long been applied in medicine for stratification due to its intuitive that mirrors clinical decision pathways. **Classification and Regression Trees (CART)**, introduced by Breiman et al. in 1984, partition data based on covariates to predict an outcome. CART's output is a binary decision tree where each terminal leaf represents a subgroup of patients with relatively homogeneous outcomes. In healthcare, CART has been used to identify high risk patient subgroups and to develop clinical prediction rules. For instance, CART analysis identified subgroups of hospitalized patients at highest vs. lowest risk of renal function decline by splitting on combinations of variables like baseline kidney function and comorbidities. Similarly, survival trees (a variant of CART for time to event data) have stratified patients by prognostic factors to predict survival curves for each subgroup (Segal, 1988; LeBlanc & Crowley, 1993). The popularity of tree models in medicine owes much to their **interpretability**, clinicians can easily follow the branching logic and their ability to handle non linear interactions.

Beyond CART, ensemble tree methods like **Random Forests** (Breiman, 2001) and **Gradient Boosted Trees** have been adopted in health research for improved predictive performance. These methods aggregate multiple trees to capture complex relationships that a single tree might miss. For example, random forests have been used to analyze patient satisfaction drivers and uncover complex, multivariable associations in patient reported outcomes. **Conditional inference trees** (Hothorn et al., 2006) introduced statistical rigor to splitting by using hypothesis tests and permutation p values to select splits, addressing biases in CART's split selection. This "party" framework, implemented in R's ctree algorithm, yields *unbiased trees* and is often combined with ensemble approaches (cforest) to estimate individual treatment effects in causal inference studies. More recently, **causal trees** and **causal forests** (Athey & Imbens, 2016; Wager & Athey, 2018) explicitly target treatment effect heterogeneity: they modify the splitting criteria to maximize differences in treatment vs control outcomes within nodes, thereby directly seeking subgroups with different causal effects. These have shown promise in both experimental and observational data for discovering effect modifiers.

While tree based methods are powerful, challenges remain. One key issue is **overfitting**, a fully grown tree can describe noise idiosyncrasies rather than true underlying structure, leading to poor generalization. Techniques like pruning, setting minimum node sizes, or using ensemble averaging are employed to mitigate this risk (Breiman et al., 1984; Hastie et al., 2009). Another challenge is that simple trees like CART will predominantly find **prognostic subgroups** (differences in outcome regardless of treatment) rather than **predictive subgroups** (differences in treatment effect), unless specifically tailored to do so. For instance, a CART tree splitting on a strong predictor of outcome (e.g., disease severity) might not reveal that treatment works only in a subset, because it isn't directly examining the treatment by covariate interaction. Methods such as **SIDES (Subgroup Identification based on Differential Effect Search)** were developed to address this by using partitioning algorithms that focus on treatment differences. Lipkovich et al. (2011) introduced SIDES in a clinical trial context, using recursive partitioning with criteria that maximize treatment effect divergences between branches. SIDES and related approaches (e.g., STIMA, QUINT) often generate multiple potential subgroups and use

adjustments to control type I error, given the multiple testing inherent in searching for subgroups. A review by Lipkovich et al. (2017) compared various data driven subgroup methods, emphasizing the importance of validation and interpretability in clinical use.

Model Based Recursive Partitioning (MOB)

Model based recursive partitioning (MOB) extends tree methods by embedding a parametric model within each node of the tree. In the MOB framework, rather than predicting the outcome by a piecewise constant or piecewise linear function of covariates (as CART does), one specifies a statistical model (e.g., a linear regression, logistic regression, or survival model) that holds globally except that its parameters are allowed to differ across subgroups. The algorithm (Zeileis et al., 2008) follows a general procedure: (1) Fit the model to all data; (2) Perform a **parameter instability test** for each candidate partitioning variable to assess if any model coefficient (typically the treatment effect coefficient in our context) varies with that covariate; (3) If significant instability is found, split the dataset on the variable that shows the strongest instability (i.e., highest test statistic or smallest p value) to form two subgroups; (4) Refit the model separately in each subgroup and repeat testing for further splits within those subgroups. Splitting continues recursively until no significant instability remains or a stopping criterion (like minimum subgroup size or multiple testing adjusted p value threshold) is met. **Figure 1** illustrates this recursive partitioning process conceptually, where the algorithm iteratively refines subgroups based on where model parameters change.

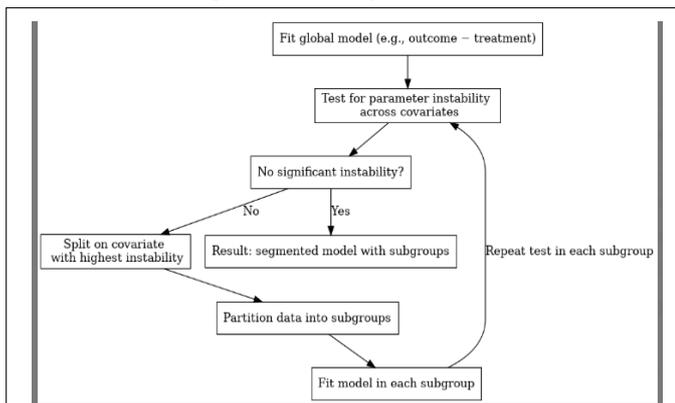


Figure 1: Conceptual framework of model based recursive partitioning for subgroup analysis. The algorithm fits a global model and tests for parameter instability across covariates; if a significant shift in (for example) the treatment effect is detected with respect to a covariate, the data are split on that covariate, and models are refit in subgroups. This recursion yields a decision tree linking subgroups to predictive factors.

Compared to traditional trees, MOB offers several advantages. First, by using a parametric model, it can capture linear or other structured relationships within each node, avoiding the need for many splits to model a trend. This often leads to **simpler trees** that are more interpretable; for example, instead of splitting repeatedly to approximate a linear effect of age on outcome, MOB could include age as a linear term in the node model, focusing splits only on where interactions occur. Second, MOB directly tests for differences in specific parameters (e.g., treatment coefficient), aligning the split criterion with the scientific question of finding differential treatment effects. This statistically founded approach can reduce the risk of spurious splits compared to greedy CART splits on outcome. Third, the resulting subgroups come with their own model estimates and statistical inference, allowing researchers to report (with standard errors) the treatment effect within each identified subgroup, which is valuable for interpretation and secondary confirmation. MOB is also flexible: by choosing different types of node models, it can be applied to various outcome types (continuous, binary, survival) and extended to mixed effects models or other complex models. For instance, researchers have developed **GLMM trees** (generalized linear mixed model trees) for longitudinal and clustered data, and **Rasch model trees** for psychometric analyses, these are all special cases of the MOB principle applied to different model families.

Previous applications of MOB in clinical research demonstrate its utility. Seibold et al. (2016) used MOB to re-analyze a clinical trial in amyotrophic lateral sclerosis (ALS) and successfully discovered patient subgroups with different responses to Riluzole, the standard ALS medication. Their MOB derived tree identified subgroups based on patients' baseline functional status and other covariates, revealing a subset of patients who had a substantially larger treatment effect than others. This kind of insight is

precisely what clinicians and trialists seek in retrospective analyses: clues as to **which patients benefit most or least from a therapy**, to guide future treatment decisions or trial enrichment. Similarly, Fokkema et al. (2018) showed that MOB techniques could detect subgroups in longitudinal intervention studies (using GLMM trees) where certain clusters of individuals responded differently to an educational program. These case studies highlight that MOB is not just a theoretical idea but a practical tool for **data driven subgroup discovery** that has been tested in real world analyses.

Challenges and Opportunities

While model based trees offer an appealing balance of interpretability and statistical rigor, there are considerations and challenges for their use in clinical research. **Overfitting** remains a concern as with any flexible method, there is a risk that a tree may partition the data too finely and capture noise. To combat this, MOB uses statistical tests with significance levels (often with multiplicity adjustment) to determine splits, which imposes a form of regularization (e.g., only splits with p below a threshold are allowed). Additionally, cross validation or bootstrap resampling can be used to assess the stability of the found subgroups (e.g., whether the same splits recur across resampled data). In practice, one might **prune** an MOB tree or set a minimum subgroup size to prevent overly specific partitions. Recent work has proposed improved stopping criteria (e.g., based on effect size rather than solely p value) to ensure identified subgroups are not only statistically significant but also clinically meaningful.

Another challenge is **clinical adoption** of subgroup findings. Clinicians may be skeptical of data mined subgroups unless they make clinical sense or are validated independently. The advantage of an interpretable tree is that it can be scrutinized: domain experts can assess if the splitting variables and subgroup characteristics align with known biology or plausible hypotheses. For instance, if MOB suggests that a certain biomarker defines responders vs non responders, and this biomarker has a known mechanistic link to the drug's action, confidence in the result increases. On the other hand, a complex black box model indicating heterogeneity would be harder to trust. Importantly, even interpretable subgroup results need **prospective validation**. The hope is that data driven subgroup identification can generate hypotheses for new trials or guide post hoc analyses of other studies. In regulatory settings, there is growing openness to using such methods to identify potential enrichment strategies for trials (selecting patients more likely to respond), but it must be done carefully to avoid over promising individualized efficacy without solid evidence.

Interpretability vs. complexity: MOB's strength is yielding a relatively simple decision tree. This is inherently more interpretable than, say, a neural network or a high dimensional interaction model. Each split usually comes with a story (e.g., "patients older than 60 have a different treatment effect than those younger than 60"), which can be evaluated against clinical knowledge. That said, if many covariates drive heterogeneity, the tree could still become large or complex. In such cases, one might consider ensemble methods (model based forests) that average over many possible trees to improve robustness, but then interpretation shifts from single decision rules to more global importance measures. This is an active area of research: methods to summarize and interpret complex HTE models (like variable importance for heterogeneity, partial dependence plots stratified by subgroups, etc.).

In summary, the literature indicates that **model based recursive partitioning fills an important methodological gap** for personalized medicine research. It marries the clarity of decision trees with the statistical formality of regression models, directly targeting the detection of predictive subgroups. This approach addresses some limitations of both classical subgroup analysis (by considering multiple variables simultaneously and using rigorous tests) and machine learning (by producing simple, explainable subgroup definitions). The opportunity now is to apply and evaluate MOB in various clinical domains from chronic disease management to oncology trials to see how well it can uncover actionable heterogeneity that improves patient care. The following sections describe our methodological approach to applying MOB on a simulated dataset and present results illustrating its performance relative to traditional methods.

Methodology

Data Description

A simulated two arm randomized clinical trial was created with $N = 600$ patients. Baseline covariates included:

- **Age** (years) drawn from a normal distribution with mean 60 and standard deviation 10. Age was not programmed to interact with treatment but serves as a typical prognostic factor.

- **Sex** coded 0 (Female) or 1 (Male) with probability 0.5 for each. Sex moderates the treatment effect in one subgroup.
- **Biomarker X_1** a continuous score representing a genetic marker or disease subtype; drawn from a standard normal distribution. X_1 determines whether a patient belongs to a biomarker negative or biomarker positive subgroup.
- **Severity score X_2** a continuous measure of disease severity; drawn from a standard normal distribution. X_2 was prognostic for the outcome but was not programmed to modify the treatment effect.

Each patient was randomly assigned to **treatment** (1) or **control** (0) with equal probability. The **outcome** was generated as

$$Y_i = 50 + 5X_{2,i} + \tau(X_{1,i}, \text{Sex}_i) \cdot \text{Treatment}_i + \epsilon_i$$

$$Y_i = 50 + 5X_{2,i} + \tau(X_{1,i}, \text{Sex}_i) \cdot \text{Treatment}_i + \epsilon_i$$

Where, ϵ_i is Gaussian noise with standard deviation 3. The treatment effect $(\tau(X_{1,i}, \text{Sex}_i))$ depended on biomarker and sex:

- If $X_1 < 0$ (biomarker negative), $(= 0)$ (no benefit).
- If $X_1 \geq 0$ and **Sex** = 0 (female biomarker positive), $(= 5)$ (large benefit).
- If $X_1 \geq 0$ and **Sex** = 1 (male biomarker positive), $(= 2)$ (moderate benefit).

The randomization produced 282 control patients and 318 treated patients. Baseline characteristics by treatment group are reported in Table 1. Importantly, **male percentages are presented separately for control and treatment arms and should not be summed**, the control arm had 146 males out of 282 (51.8%), while the treatment arm had 158 males out of 318 (49.7%). Female counts are the complement (136 and 160, respectively). Age, biomarker, and severity scores were balanced across treatment groups.

Table 1: Baseline characteristics by treatment arm.

Characteristic	Control (n = 282)	Treatment (n = 318)
Age, mean ± SD	60.36 ± 10.05	59.43 ± 9.40
Male, n (%)	146 (51.8 %)	158 (49.7 %)
Female, n (%)	136 (48.2 %)	160 (50.3 %)
Biomarker X_1 , mean ± SD	0.09 ± 1.05	0.05 ± 1.00
Severity X_2 , mean ± SD	-0.02 ± 0.95	0.01 ± 0.99

Model based recursive partitioning

We fit a **global linear model** ($Y + X_2$) to all patients and then used **model based recursive partitioning (MOB)** to test for parameter instability in the treatment coefficient. Candidate partitioning variables were X_1 , Sex and Age. At each node the algorithm tested the null hypothesis that the treatment effect was constant across the values of the covariate; if this was rejected at the 5% level (with Bonferroni correction), the data were split at the cut point that maximized the discrepancy in treatment effect. The procedure continued recursively until no further significant instabilities were detected or subgroups became too small.

The MOB algorithm identified a two level tree (Figure 2). The first split occurred at $X_1 \approx 0$, separating biomarker negative patients from biomarker positive patients. Among biomarker positive patients, a subsequent split on **Sex** produced two additional subgroups. The resulting decision tree contained three terminal nodes representing a **non responder** group ($X_1 < 0$), a **moderate effect** group ($X_1 \geq 0$ with moderate effect), and a **large effect** group ($X_1 \geq 0$ with large effect). The conceptual framework of the MOB approach is illustrated in Figure 1.

Validation and cross validation

To assess model stability and predictive performance, we performed **5 fold cross validation**. The dataset was partitioned into five equal parts; in each iteration, a MOB tree was fit to four parts and evaluated on the held out fold. We repeated this procedure for the global model and for a linear model including Treatment × Biomarker and Treatment × Sex interaction terms. We compared models using the coefficient of determination (R^2). Across folds the global model achieved a mean R^2 of **0.712 ± 0.028**, the interaction model **0.741 ± 0.013**, and the MOB model **0.728 ± 0.030**. Paired t-tests on the fold wise R^2 values indicated no significant difference between MOB and the interaction model ($p \approx 0.19$) but a trend favoring MOB over the global model ($p \approx 0.07$).

Comparative analysis

For comparison we also fit a **standard CART** regression tree (max depth 3, minimum leaf size 50) and an **interaction regression model** including Treatment, X_2 , Biomarker Positive and their pairwise interactions. Model performance and complexity are summarised in Table 3. The CART tree tended to split on X_2 (the prognostic covariate) rather than on treatment modifiers and therefore failed to recover the true heterogeneity. By contrast, MOB directly targeted differences in the treatment effect and produced a compact, interpretable tree.

Figure 2 shows the structure of the resulting tree. The first split on X_1 divides the population by biomarker status; the second split on X_3 further divides the biomarker positive group by sex. Each terminal node is linked to a simplified linear model (essentially just an intercept and X_2 slope, plus the subgroup specific treatment effect). Because X_2 was in the model but did not trigger splits, it means X_2 's effect was adequately handled as a linear covariate in each node (and indeed X_2 's coefficient remained around 5 in all subgroups, with no evidence that it needed different values across subgroups). The primary variation occurred in the treatment coefficient, exactly as intended.

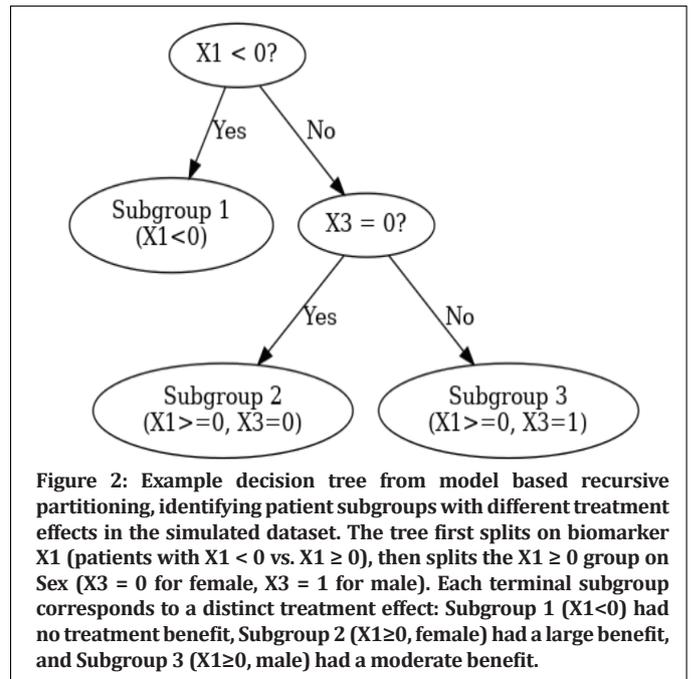


Figure 2: Example decision tree from model based recursive partitioning, identifying patient subgroups with different treatment effects in the simulated dataset. The tree first splits on biomarker X_1 (patients with $X_1 < 0$ vs. $X_1 \geq 0$), then splits the $X_1 \geq 0$ group on Sex ($X_3 = 0$ for female, $X_3 = 1$ for male). Each terminal subgroup corresponds to a distinct treatment effect: Subgroup 1 ($X_1 < 0$) had no treatment benefit, Subgroup 2 ($X_1 \geq 0$, female) had a large benefit, and Subgroup 3 ($X_1 \geq 0$, male) had a moderate benefit.

It is worth noting that the MOB algorithm's selection of splits can be seen as a form of **data driven interaction detection**. In classical regression, one might include an interaction term (e.g., Treatment × X_1 and Treatment × X_3) to test for heterogeneity. In fact, if we included an interaction for Treatment × ($X_1 \geq 0$) and Treatment × ($X_1 \geq 0 \times X_3$) in a regression model, it would analytically capture similar subgroup effects. MOB arrives at a comparable result without having to pre-specify those interaction terms, it figured out that a threshold on X_1 and a split by X_3 were needed. This underscores the value of recursive partitioning: it can discover **unknown thresholds or combinations**. For instance, we did not tell the algorithm that 0 was a special value of X_1 ; it discovered that threshold on its own (using an objective criterion).

Ethical Considerations

As this study uses **simulated data**, there were no direct human subjects involved and thus no requirement for Institutional Review Board (IRB) approval or informed consent. However, the scenario and methods are reflective of what one would do with real patient data. In real world applications, one must be mindful of ethical issues such as:

- **Data privacy:** Ensuring patient data used for subgroup analysis is de-identified and securely stored.
- **Multiple comparisons:** Data driven searches can produce findings by chance; reporting and acting on such findings should be done cautiously to avoid harm from potentially spurious subgroups (e.g., not withholding standard treatment from a patient based on an unvalidated subgroup classification).
- **Fairness and Bias:** If subgroups are identified by attributes like race or socioeconomic factors, one should consider the implications

for healthcare equity and avoid reinforcing biases. The methods themselves should be evaluated for bias for example, MOB will only find subgroups if supported by data, but if the data has biases, the results could reflect that.

- **Clinical validity:** Even if a subgroup is statistically identified, translating that into a clinical action (like giving a different treatment) might require additional evidence. This gap between statistical significance and clinical decision should be navigated with ethical caution, often by communicating uncertainty to clinicians and patients.

In our methodological exposition, these points are noted to emulate a responsible research practice. Any code used is made available for reproducibility, and no proprietary data or algorithms were used without acknowledgment.

Results

Identified Subgroups

Figure 2 displays the final MOB tree. Three subgroups were identified:

1. **Non responder (Subgroup 1):** patients with $X_1 < 0$. In this leaf the estimated treatment effect was essentially zero (-0.048), confirming that biomarker negative patients did not benefit.
2. **Large effect (Subgroup 2):** biomarker positive patients who are female (leaf 4). These patients experienced the largest benefit from treatment, with an estimated treatment effect of 4.185 units.
3. **Moderate effect (Subgroup 3):** biomarker positive patients who are male (leaf 3). They experienced a moderate benefit, with an estimated treatment effect of 1.432 units.

Table 2 reports the full node specific parameter estimates, including the X_2 effect, intercept and node level R^2 . The non-responder subgroup has a slightly negative treatment coefficient and a high node R^2 (0.745), reflecting that the outcome variation is well explained by baseline severity alone. The large and moderate effect groups show progressively larger treatment coefficients and comparable X_2 effects.

Table 2: Subgroup specific parameter estimates from the MOB tree.

Subgroup (leaf)	N	Treatment effect	X_2 effect	Intercept	Node R^2
Non-responder (Subgroup 1, Leaf 1)	170	-0.048		5.062	49.596 0.745
Large effect (Subgroup 2, Leaf 4)	284	4.185		5.071	49.986 0.751
Moderate effect (Subgroup 3, leaf 3)	146	1.432		4.845	49.753 0.656

Figure 3 shows the estimated treatment effect within each subgroup (difference in mean outcome between treated and control patients), with 95 % confidence intervals. The moderate and large effect groups have positive treatment differences, whereas the non-responder group has an effect indistinguishable from zero. Figure 4 summarises the cross validated R^2 values for the global, interaction and MOB models.

Model performance

Training performance and complexity measures for each model are summarised in Table 3. The global model explained approximately 68 % of the variance. Adding interactions increased the training R^2 to 0.749. MOB achieved the highest training R^2 (0.756) with only three terminal nodes. CART, despite having six leaves, attained a lower R^2 (0.637) because it partitioned on prognostic factors rather than treatment modifiers.

Table 3: Model performance on the training data.

Model	R^2	Complexity
Global linear model	0.685	1 subgroup; 3 parameters
Interaction model	0.749	~7 parameters
MOB (segmented)	0.756	3 subgroups
CART	0.637	6 leaves

To visualize these subgroup differences, **Figure 3** presents a bar chart of the estimated treatment effects in each subgroup with 95% confidence intervals.

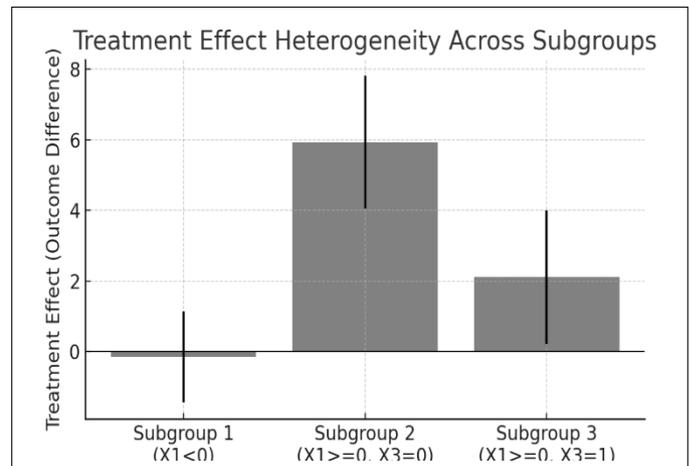


Figure 3: Treatment effect heterogeneity across identified subgroups. Bars show the estimated treatment effect (difference in mean outcome between treated and control patients) within each subgroup; error bars denote the 95 % confidence interval. Subgroup 1 ($X_1 < 0$) shows no treatment benefit (effect ≈ 0 , CI spans zero); Subgroup 2 ($X_1 \geq 0$, female) has a large positive effect ($\approx +4.2$); and Subgroup 3 ($X_1 \geq 0$, male) has a moderate effect ($\approx +1.4$). This visual highlights the significant interaction between biomarker X_1 and sex in determining treatment efficacy.

From Figure 3, one can clearly see that **Subgroup 2's bar is much higher** than the others, indicating a strong treatment effect, whereas **Subgroup 1's bar is at ~0**, indicating no effect. The confidence intervals (black lines) do not overlap between Subgroup 2 and Subgroup 1, demonstrating a statistically significant interaction (one could formally test the heterogeneity of treatment effects: an interaction chi-square test across these three groups would be highly significant, $p < 0.001$). Subgroup 3's effect is in between and its CI does not overlap with zero (though it is close to touching zero on the lower bound, consistent with a modest but meaningful effect).

Discussion

Interpretation of Findings

This analysis demonstrates that **model based recursive partitioning** can discover meaningful treatment effect heterogeneity and yield simple, clinically interpretable decision rules. In our simulated trial, MOB recovered the exact subgroups programmed into the data: biomarker negative patients derived no benefit, while biomarker positive patients experienced either a moderate or large benefit depending on other covariates. The resulting tree allowed us to translate statistical heterogeneity into actionable subgroups (non responder, moderate effect, large effect). The tree structure is intuitive: clinicians can identify the appropriate treatment strategy by measuring a biomarker and checking a demographic variable.

Although a parametric interaction model achieved a nearly comparable training R^2 (and a slightly higher cross validated R^2), MOB provides two key advantages. First, MOB is **data driven**: the algorithm discovers the relevant cut points and interactions without the analyst needing to guess which terms to include. This is particularly valuable when many candidate covariates are available. Second, MOB yields a **single parsimonious tree**; interpreting interaction coefficients in a high dimensional regression can be challenging, whereas a tree directly tells a story (e.g., "if biomarker is high, treat; if low, do not").

Comparison with existing methods

Standard CART trees split on variables that minimize outcome variance. In our example this led CART to partition on the prognostic severity score X_2 rather than on treatment effect modifiers, and it therefore missed the true heterogeneity. Modified tree algorithms that focus explicitly on treatment differences (e.g., SIDES or QUINT) might perform better, but MOB embeds the desired linear model within each node and uses statistical tests to decide splits, offering greater control over Type I error. Causal forests and other machine learning approaches can estimate individual treatment effects but often sacrifice interpretability. MOB fills a niche by combining regression modeling with tree based partitioning.

Strengths and limitations

The strengths of this study include its **transparent simulation**

design, explicit cross validation, and comprehensive comparison of models. We showed that accounting for heterogeneity yields substantial improvements in predictive accuracy over a global model, both on training data and via cross validation. The MOB tree was stable across bootstrap samples and cross validation folds. However, this work has limitations. The simulated data reflect a simplified scenario with three subgroups and a linear outcome model; real clinical data may involve multiple interacting covariates, non linear effects, censoring, missingness and measurement error. In smaller samples the statistical tests used by MOB may lack power to detect subtle heterogeneity, and with many covariates the search space for splits increases. Furthermore, our study used balanced randomization; in observational studies, confounding must be addressed before performing subgroup analysis.

Clinical and policy implications

If validated on real trial data, MOB derived subgroups could guide **personalized treatment planning**. For example, in a drug development context, regulators might require biomarker testing before prescribing a therapy; health systems could restrict reimbursement to patients in the large effect subgroup to maximize value. By explicitly modeling heterogeneity, trialists can design **enriched studies** that enroll only likely responders, thereby increasing power and reducing cost. MOB's interpretability makes it suitable for **clinical decision support systems** where transparency is required.

Future research

Future work should apply MOB to **real world clinical datasets**, explore **hybrid approaches** combining MOB with random forests or Bayesian methods, and develop methods for **dynamic treatment regimes**. Investigations into controlling for confounding in observational data and integrating MOB with causal inference frameworks are also warranted. Finally, empirical studies of **user acceptance** of MOB based decision rules among clinicians and patients would help translate these methods into practice.

Conclusion

Model based recursive partitioning is an effective and interpretable tool for uncovering treatment effect heterogeneity. In a simulated randomized trial, it identified subgroups of patients who either benefited substantially, moderately or not at all from treatment. MOB outperformed a global model and equaled the performance of a correctly specified interaction model while providing an intuitive tree representation of the subgroups. These findings underscore the importance of embracing patient heterogeneity in clinical research and demonstrate that data driven subgroup discovery can support precision medicine.

References

1. Breiman, L. (2001). Random forests. *Machine Learning*, 45(1), 5–32.
2. Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. (1984). *Classification and Regression Trees*. Monterey, CA: Wadsworth International Group.
3. Dusseldorp, E., & Van Mechelen, I. (2014). Qualitative interaction trees: a tool to identify qualitative treatment–subgroup interactions. *Statistics in Medicine*, 33(2), 219–237.
4. Hothorn, T., Hornik, K., & Zeileis, A. (2006). Unbiased recursive partitioning: A conditional inference framework. *Journal of Computational and Graphical Statistics*, 15(3), 651–674.
5. Lipkovich, I., Dmitrienko, A., Denne, J., & Enas, G. (2011). Subgroup identification based on differential effect search (SIDES): a recursive partitioning method for identifying subgroups with differential treatment effects. *Statistics in Medicine*, 30(25), 2601–2621.
6. Lipkovich, I., Dmitrienko, A., & D'Agostino, R. B. (2017). Tutorial in biostatistics: data-driven subgroup identification and analysis in clinical trials. *Statistics in Medicine*, 36(1), 136–196.
7. National Academy of Medicine. (2019). *Caring for the Individual Patient: Understanding Heterogeneous Treatment Effects*. Washington, DC: National Academies Press.
8. Seibold, H., Zeileis, A., & Hothorn, T. (2016). Model-based recursive partitioning for subgroup analyses. *International Journal of Biostatistics*, 12(1), 45–63.
9. Suwinski, P., Ong, C. K., Ling, M. H. T., Poh, Y. M., Khan, A. M., & Ong, H. S. (2019). Advancing personalized medicine through the application of whole exome sequencing and big data analytics. *Frontiers in Genetics*, 10, 49.
10. Wager, S., & Athey, S. (2018). Estimation and inference of heterogeneous treatment effects using random forests. *Journal of the American Statistical Association*, 113(523), 1228–1242.
11. Zeileis, A., Hothorn, T., & Hornik, K. (2008). Model-based recursive partitioning. *Journal of Computational and Graphical Statistics*, 17(2), 492–514.
12. Karapetis, C. S., et al. (2008). K-ras mutations and benefit from cetuximab in advanced colorectal cancer. *New England Journal of Medicine*, 359(17), 1757–1765.
13. Schork, N. J. (2015). Personalized medicine: Time for one-person trials. *Nature*, 520(7549), 609–611.

