

Accounting for Individual Differences in Bradley-Terry Models by Means of Recursive Partitioning

Carolin Strobl

Ludwig-Maximilians-Universität München

Florian Wickelmaier

Eberhard Karls Universität Tübingen

Achim Zeileis

Universität Innsbruck

The preference scaling of a group of subjects may not be homogeneous, but different groups of subjects with certain characteristics may show different preference scalings, each of which can be derived from paired comparisons by means of the Bradley-Terry model. Usually, either different models are fit in predefined subsets of the sample or the effects of subject covariates are explicitly specified in a parametric model. In both cases, categorical covariates can be employed directly to distinguish between the different groups, while numeric covariates are typically discretized prior to modeling. Here, a semiparametric approach for recursive partitioning of Bradley-Terry models is introduced as a means for identifying groups of subjects with homogeneous preference scalings in a data-driven way. In this approach, the covariates that—in main effects or interactions—distinguish between groups of subjects with different preference orderings, are detected automatically from the set of candidate covariates. One main advantage of this approach is that sensible partitions in numeric covariates are also detected automatically.

Keywords: *Bradley-Terry-Luce model; subject covariates; recursive partitioning*

1. Introduction

The choice model suggested by Bradley and Terry (1952) is the most widely used means for deriving a latent preference scale from paired comparison data

The authors would like to thank the participants of the “Psychometric Computing 2009” workshop for feedback and fruitful discussions—especially Regina Dittrich and Reinhold Hatzinger for sharing both their code and their expertise for paired comparison models.

when no natural measuring scale is available. In a measurement-theoretic approach, Luce (1959) showed that the model can be derived from a simple axiom for the choice probabilities. Therefore, the model is referred to either as the Bradley-Terry-Luce (BTL) model or the Bradley-Terry (BT) model. (In the following we will use the latter abbreviation.)

The BT model has been applied in a variety of fields in psychology and related disciplines. Early applications and developments are summarized in an extensive bibliography compiled by Davidson and Farquhar (1976) containing more than 350 references. More recent applications include, for example, surveys on health care, educational, and political choice (Dittrich, Francis, Hatzinger, & Katzenbeisser, 2006; Dittrich, Hatzinger, & Katzenbeisser, 1998; McGuire & Davison, 1991) as well as psychophysical studies on the sensory evaluation of pain, sound, and taste (Choisel & Wickelmaier, 2007; Matthews & Morris, 1995; Oberfeld, Hecht, Allendorf, & Wickelmaier, 2009).

In many applications it is reasonable to assume that the preference scaling of a group of subjects not only depends on characteristics of the stimuli to be judged by the subjects but also on characteristics of the subjects themselves. It is common practice to fit separate BT models, for example, for younger and older participants (Kissler & Bäuml, 2000; McGuire & Davison, 1991). In more advanced approaches (such as Böckenholt, 2001a, 2001b; Dittrich et al., 1998), the covariates are explicitly included in the model. An overview over seminal and advanced methods on preference scaling in the field of psychology is given by Böckenholt (2006). An overview of methods applied in the field of market segmentation, including mixture approaches for paired comparisons, is given by Wedel and Kamakura (2000).

Here, we suggest a new approach for incorporating subject covariates in BT models: The approach of model-based recursive partitioning, that is well established in the field of statistics and machine learning, can be applied intuitively to identifying groups of subjects that differ in their preference scalings.

The approach of model-based recursive partitioning in general, as well as the framework for treating the BT model with this approach, is introduced in the following section. Two application examples are presented to illustrate the usage and benefits of this new technique for incorporating subject covariates in BT models. The main differences between model-based recursive partitioning and previous approaches for incorporating subject covariates in choice models are reviewed in the discussion.

2. Method

Model-based partitioning employs the same principle as the more widely known classification and regression trees (CART; Breiman, Friedman, Olshen, & Stone, 1984): The covariate space is recursively partitioned to distinguish between groups of subjects with different characteristics. In the following, we will

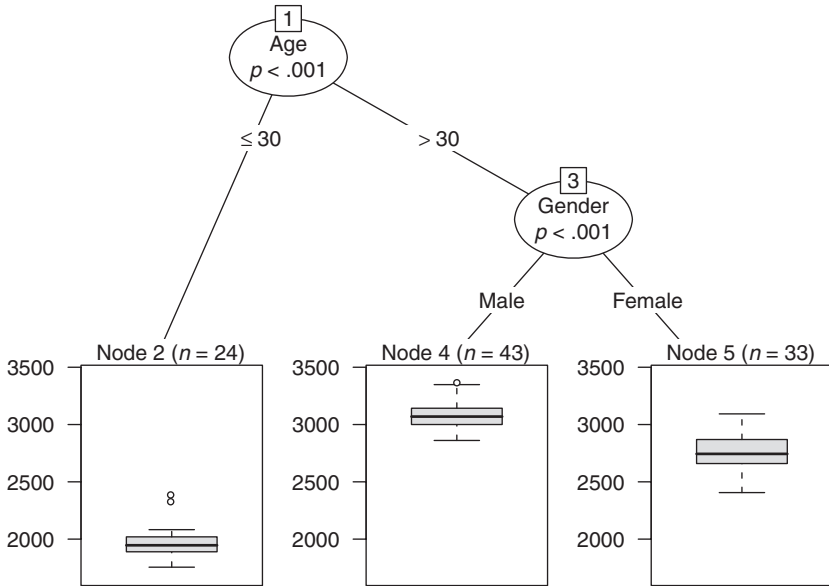


FIGURE 1. Exemplary regression tree indicating that the average income varies in groups defined by a combination of the covariates age and gender (artificial data).

briefly outline the rationale of recursive partitioning in general, before we provide the framework and technical details for model-based partitioning of the BT model.

2.1. A Brief Introduction to Recursive Partitioning

Following the principle of recursive partitioning, CART produce a tree-structured partition of the covariate space, where, starting with the entire sample, subjects are divided into groups according to their values of selected covariates. The splitting rules represented by the tree are chosen such that the subjects within the resulting groups have similar values of the response variable, whereas their response values differ from the subjects in the other groups.

The illustrative example in Figure 1 shows a regression tree as applied to an artificial data set with response income and covariates *gender* and *age*. The tree detects three groups with different income levels: Subjects under the age of 30 (in the leftmost Node 2) have a low average income, male subjects over 30 (in the middle Node 4) have a high average income, and females over 30 (in the rightmost Node 5) have a medium average income. (Note that the node numbers are only labels assigned recursively from left to right starting from the top node.)

In comparison to simple CART, model-based partitioning does not aim at finding groups of subjects with different values of the response variable but with

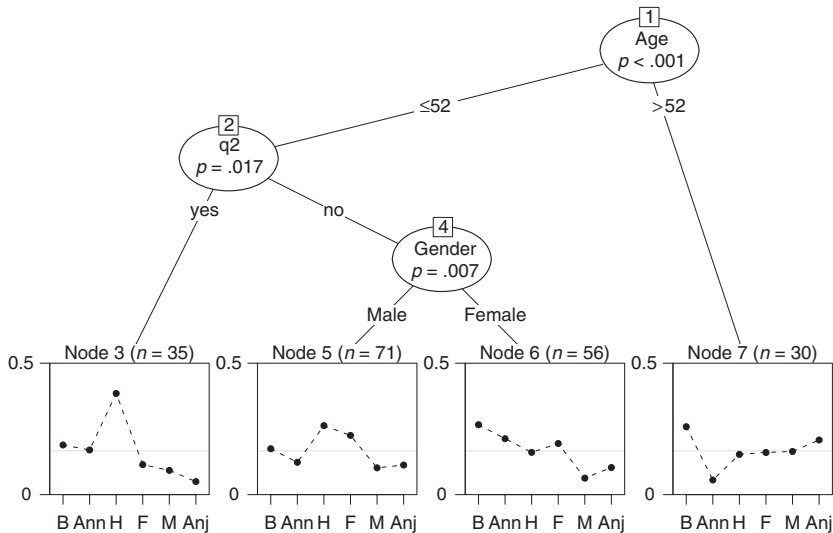


FIGURE 2. Partitioned paired comparison model for the attractiveness data, indicating that judged attractiveness varies in groups defined by combinations of the covariates age, q2, and gender (B: Barbara, Ann: Anni, H: Hana, F: Fiona, M: Mandy, and Anj: Anja).

different values of certain model parameters. Such parameters could be, for example, intercepts and slopes in a linear regression or—as in our case—the worth parameters of the stimuli in a BT model, which may vary between groups of subjects.

An example for a BT-based tree is displayed in Figure 2. Here, the preference scales for the stimuli derived from the BT model vary between the groups of subjects represented by the terminal nodes: Subjects up to the age of 52 who answered yes to Question 2 (in the leftmost Node 3) clearly prefer the third stimulus over all other stimuli, while, for example, subjects over 52 (in the rightmost Node 7) prefer all other stimuli over the second stimulus, and so on.

In the remainder of this section, the construction of the tree is described in detail, serving as an illustration of the general method for estimation of BT tree models.

The data underlying the tree displayed in Figure 2 were collected in a study at the Department of Psychology, Universität Tübingen: $n = 192$ subjects were interviewed and asked to judge the attractiveness of the candidates of the second season of “Germany’s Next Topmodel,” which aired March through May 2007.

“Germany’s Next Topmodel” is a casting show for topmodel hopefuls on German television—an adaptation of the corresponding U.S. show “America’s Next Topmodel.” Based on photos of the contestants taken at the beginning of the season, the participants of the study judged the attractiveness of the $k = 6$ finalists (Barbara Meier, Anni Wendler, Hana Nitsche, Fiona Erdmann, Mandy

Graff, and Anja Platzer, listed here in decreasing order, that is, starting with the winner of the show, Barbara Meier) in a forced choice experiment.

The stimuli were digital portrait photographs (resolution 303×404 pixels) of the contestants. Participants were presented with all 15 pairs of photographs. In each trial, their task was to judge which of the two women on the photos was more attractive.

The sample was stratified with respect to *gender* and *age* (younger vs. older than 30 years) with an equal number of subjects in each group. Overall, the sample contained 96 female and 96 male raters between the ages of 15 and 77.

Additionally, several subject-specific covariates about the raters are available: *gender*, *age*, and the answers (yes/no) to the following three questions:

Question 1 (q1): Do you know the women on the photos? Do you know the TV show Germany's Next Topmodel?

Question 2 (q2): Did you watch the latest season of Germany's Next Topmodel regularly?

Question 3 (q3): Have you seen the final of the latest season of Germany's Next Topmodel? Do you know who won the latest season of Germany's Next Topmodel?

For Questions 1 and 3, a positive answer to at least one of the subquestions resulted in a positive overall answer.

As explained in detail in the remainder of this section, the recursively partitioned BT model displayed in Figure 2 was generated by means of a simple algorithm consisting of the following consecutive steps:

1. Fit a BT model to the paired comparisons of all subjects in the current subsample, starting with the full sample.
2. Assess the stability of the BT model parameters with respect to each available covariate.
3. If there is significant instability, split the sample along the covariate with the strongest instability and use the cutpoint with the highest improvement of the model fit.
4. Repeat Steps 1–3 recursively in the resulting subsamples until there are no more significant instabilities (or the subsample is too small).

We will now go through each of the steps of this algorithm, provide the statistical tools, and explain how they were used to generate the model-based partition of the BT model depicted in Figure 2.

2.2. Fitting Bradley-Terry Models

To fix notation, we consider paired comparison models with possible ties (see e.g., Critchlow & Fligner, 1991; Section 4): Each comparison of two stimuli can result in (1) the first stimulus being preferred, (2) the second stimulus being

preferred, or (3) the subject being undecided between the two stimuli (i.e., a tie). The common forced choice experiments, where ties are prohibited by the experimental design, can be considered as a special case of this more general view.

In a notation similar to Critchlow and Fligner (1991), we consider $i = 1, \dots, n$ subjects who judge all unordered pairs of $j = 1, \dots, k$ stimuli. Thus, each subject performs $k^* = k \cdot (k - 1)/2$ comparisons—each resulting in a choice for an answer c in 1, 2, 3. According to the Davidson (1970) extension of the BT model, the three possible outcomes have probabilities:

$$\begin{aligned}
 p_{jj'1} &= \frac{\pi_j}{\pi_j + \pi_{j'} + \nu\sqrt{\pi_j\pi_{j'}}}, \\
 p_{jj'2} &= \frac{\pi_{j'}}{\pi_j + \pi_{j'} + \nu\sqrt{\pi_j\pi_{j'}}}, \\
 p_{jj'3} &= \frac{\nu\sqrt{\pi_j\pi_{j'}}}{\pi_j + \pi_{j'} + \nu\sqrt{\pi_j\pi_{j'}}},
 \end{aligned}$$

where $\pi_j \geq 0$ ($j = 1, \dots, k$) are stimulus-specific parameters, also called *worth parameters* or *merits*, and $\nu \geq 0$ is a discrimination constant governing the probability of ties.

This formulation of the model is easy to interpret but has two drawbacks when it comes to parameter estimation: It is overidentified (multiplication of all π_j with a constant does not change the probabilities) and the parameters are constrained to be non-negative. Hence, for parameter estimation, one parameter is typically kept fixed and all others are considered on a log-scale yielding the k -dimensional parameter $\theta = (\log(\pi_1), \dots, \log(\pi_{k-1}), \log(\nu))^T$. Without loss of generality $\log(\pi_k)$ is fixed at zero; equivalently, the sum of the worth parameters can be constrained to 1. This latter view will be adopted for reporting the π_j in our empirical results.

Note that the classical BT model for forced choice experiments without ties follows as the simple special case when $\nu = 0$ and thus $p_{jj'3} = 0$. Consequently, the parameter θ is only $k - 1$ -dimensional for the BT model.

Given $i = 1, \dots, n$ observations $y_i \in \{1, 2, 3\}^{k^*}$, that is, each y_i containing the k^* comparisons with outcome $c = 1, 2, 3$, the joint log-likelihood is given by:

$$\begin{aligned}
 \log L(\theta|y_1, \dots, y_n) &= \sum_{i=1}^n \sum_{j < j'}^3 \sum_{c=1}^3 I(y_{i,jj'} = c) \log(p_{jj'c}) \\
 &= \sum_{i=1}^n \Psi(y_i, \theta),
 \end{aligned}$$

where $\Psi(y_i, \theta)$ denotes the likelihood contribution of the i -th observation and $I(\cdot)$ is the indicator function. The parameter estimates $\hat{\theta}$ can then be obtained by maximum likelihood (ML) estimation:

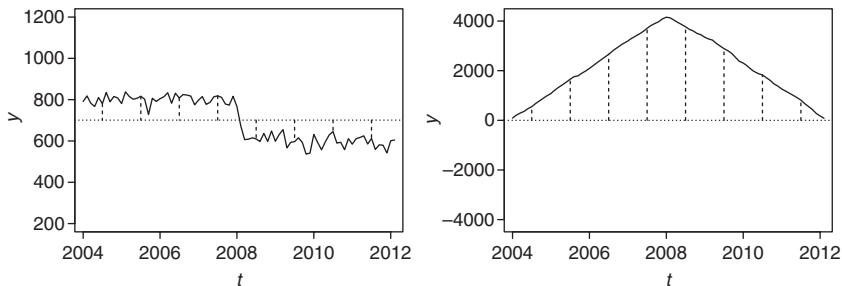


FIGURE 3. *Structural change in the mean over time (artificial data). In the left plot, the dotted line indicates the overall mean, the dashed lines illustrate that deviations from the overall mean are positive before the structural change and negative afterward. In the right plot, the positive and negative deviations are cumulated and the structural change is noticeable from the triangular shape of the path of the cumulative sum process.*

$$\hat{\theta} = \operatorname{argmax}_{\theta} \sum_{i=1}^n \Psi(y_i, \theta).$$

Typically, the ML estimate is not derived by direct maximization of the multinomial likelihood above but instead by fitting a surrogate log-linear Poisson model for the aggregate frequencies $n_{j'c} = \sum_{i=1}^n I(y_{i,jj'} = c)$. This can be easily performed in many statistical software packages, see Critchlow and Fligner (1991) for more details.

2.3. Assessing Parameter Instability in Bradley-Terry Models

After fitting the BT model, the next question is whether the set of parameters θ (i.e., the preference scale) is really the same for all n subjects or there are subsamples with differing sets of parameters. As our aim is to capture potential effects of the available covariates, it should be formally tested whether there are parameter instabilities over one of the covariates. This step will be repeated recursively within the newly created subsamples. However, in order to keep the notation simple, we only use the full-sample notation in the following.

The simplest example of parameter instability (also termed structural change or structural break in the literature) is a change in a single parameter, such as a shift in the mean. The artificial example of a change in the mean illustrated in Figure 3 (left) could depict, for example, a drop in stock returns or consumer spendings after a financial crisis.

Technically speaking, what is depicted in Figure 3 (left) can be considered as the change in the parameter “mean return” over the order implied by the variable “time”—and from this understanding, it is only a small step to describing the change in any kind of parameter over the range of any variable: The values of

the model parameter of interest can be ordered with respect to each candidate variable, and the significance of the structural change over the range of this variable can be tested statistically.

Various approaches are conceivable for this objective. A particularly convenient one—due to its generality and ease of computation—is the usage of fluctuation tests (Zeileis & Hornik, 2007) as adopted in the model-based recursive partitioning framework of Zeileis, Hothorn, and Hornik (2008). The idea of this class of tests is to compute subject-wise model deviations that should fluctuate randomly around zero under the null hypothesis of parameter stability.

In our example in Figure 3 (left), under the null hypothesis of parameter stability, the overall mean (dotted line) should hold over the entire time range. Accordingly, the deviations from the overall mean (dashed lines) should not show any systematic variation under the null hypothesis, while under the alternative of a structural break, we would expect the deviations to differ systematically from zero before and after the cutpoint, as is actually the case in Figure 3.

A general measure of deviation for likelihood-based models is the subject-wise *score function* or *estimating function*: the derivative of the likelihood contributions with respect to the parameter vector. For the BT model, these are given by:

$$\psi(y_i, \theta) = \frac{\partial \Psi(y_i, \theta)}{\partial \theta} = \sum_{j < j'}^3 \sum_{c=1}^3 I(y_{i,jj'} = c) \frac{\partial \log(p_{jj'c})}{\partial \theta}.$$

Thus, for computing the estimating functions for the parameters θ , the gradients of the log-probabilities $\log(p_{jj'c})$, $c = 1, 2, 3$, just need to be aggregated suitably. These can be shown to be:

$$\begin{aligned} \frac{\partial \log(p_{jj'1})}{\partial \theta_h} &= \begin{cases} 1 & -p_{jj'1} & -0.5p_{jj'3} & h = j \\ & -p_{jj'2} & -0.5p_{jj'3} & h = j' \\ & -p_{jj'3} & & h = k \\ 0 & & & \text{otherwise} \end{cases} \\ \frac{\partial \log(p_{jj'2})}{\partial \theta_h} &= \begin{cases} & -p_{jj'1} & -0.5p_{jj'3} & h = j \\ 1 & -p_{jj'2} & -0.5p_{jj'3} & h = j' \\ & -p_{jj'3} & & h = k \\ 0 & & & \text{otherwise} \end{cases} \\ \frac{\partial \log(p_{jj'3})}{\partial \theta_h} &= \begin{cases} 0.5 & -p_{jj'1} & -0.5p_{jj'3} & h = j \\ 0.5 & -p_{jj'2} & -0.5p_{jj'3} & h = j' \\ 1 & -p_{jj'3} & & h = k \\ 0 & & & \text{otherwise} \end{cases} \end{aligned}$$

With this notion of model deviation available, it can be assessed whether systematic deviations occur along one of the m covariates: $x_{i\ell}$ ($i = 1, \dots, n$,

$\ell = 1, \dots, m$). To do so, the deviations are cumulatively aggregated along each of the m covariates:

$$W_\ell(t) = \hat{V}^{-1/2} n^{-1/2} \sum_{i=1}^{\lfloor n \cdot t \rfloor} \psi(y_{(i|\ell)}, \hat{\theta}) \quad (0 \leq t \leq 1),$$

where the index $(i|\ell)$ denotes the i -th ordered observation with respect to the ℓ -th covariate, $\lfloor \cdot \rfloor$ denotes the integer part, and $\hat{V} = \sum_{i=1}^n \psi(y_i, \hat{\theta}) \psi(y_i, \hat{\theta})^T$ is the outer-product-of-gradients estimate of the covariance matrix.

For the artificial example of a structural change in the mean in Figure 3 (left), the cumulative aggregation is depicted in Figure 3 (right): Since the deviations from the overall mean are positive in the first half of Figure 3 (left), the cumulative aggregation in Figure 3 (right) rises up to the point of the structural change and decreases again when the negative deviations from the second half are added. Thus, the sharp kink in the path of the cumulative aggregation in Figure 3 (right) is an indicator of the strong structural change in the mean in Figure 3 (left).

The cumulative aggregation is used here to incorporate the order of the individual deviations with respect to the considered covariate: The $i = 1, \dots, n$ individual deviations are ordered with respect to the ℓ -th covariate and aggregated up to the $\lfloor n \cdot t \rfloor$ -th element in each step. When $W_\ell(t)$ is considered as a function of the fraction t of the sample size, the null model with no structural change corresponds to the path of a random process with constant zero mean.

The advantage of this approach is that the model does not have to be reestimated for all subsamples, because the individual deviations remain the same and only their order (and the corresponding path of $W_\ell(t)$) needs to be adjusted for evaluating the different covariates.

Under the null hypothesis of parameter stability, the cumulative sum process $W_\ell(\cdot)$ can be shown to converge to a k -dimensional Brownian bridge (Zeileis & Hornik, 2007), which can be used as the basis for statistical inference. To capture systematic deviations in $W_\ell(\cdot)$, different test statistics can be used depending on whether the ℓ -th covariate is a numeric or a categorical variable. If it is numeric, Zeileis et al. (2008) point out that a natural test statistic is:

$$S_\ell = \max_{i=\underline{i}, \dots, \bar{i}} \left(\frac{i \cdot n - i}{n \cdot n} \right)^{-1} \left\| W_\ell \left(\frac{i}{n} \right) \right\|_2^2.$$

This can be interpreted as the maximum Lagrange-multiplier statistic (also known as score statistic) for a single shift alternative over all conceivable cut-points in $[\underline{i}, \bar{i}]$. The limiting distribution is the supremum of a tied-down Bessel process from which p values can be computed (see Zeileis et al., 2008, for details).

If, on the other hand, the ℓ -th covariate is categorical (with $q = 1, \dots, Q$ categories, say), it is more natural to use the following test statistic:

TABLE 1
Parameter Instability Test Statistics S_1, \dots, S_5 and Corresponding P Values for the Full-Sample Bradley-Terry (BT) Model for the Attractiveness Data

	<i>Gender</i>	<i>Age</i>	<i>q1</i>	<i>q2</i>	<i>q3</i>
Statistic	17.0880	32.3566	12.6320	19.8392	6.7586
<i>p</i> value	.0217	.0008	.1283	.0067	.7452

$$S_\ell = \sum_{q=1}^Q n \left(\sum_{i=1}^n I(x_{i\ell} = q) \right)^{-1} \left\| \Delta_q W_\ell \left(\frac{i}{n} \right) \right\|_2^2,$$

where Δ_q is the increment within the q -th category. This test statistic is invariant to reordering of the Q categories and the subjects within each category. The test statistic captures the instability over the Q subsamples. Its limiting distribution is χ^2 with $(Q - 1) \cdot k$ degrees of freedom from which p values can be computed.

Although the technical details of this testing procedure are somewhat challenging, the results are easy to interpret: Parameter instability test statistics S_ℓ ($\ell = 1, \dots, m$) with associated p values (corrected for multiple testing) are provided for each candidate variable. The variable with the smallest p value is then used for determining the subsamples in the current step of the recursive partitioning algorithm—unless all p values exceed the significance level (commonly 5%), indicating that there is no (more) significant parameter instability and thus no need for partitioning.

For the attractiveness data example, the parameter instability test statistics and p values of each candidate splitting variable in the full sample are displayed in Table 1. Accordingly, the variable *age* associated with the smallest p value is used for the first split in Figure 2. The choice of the cutpoint within the chosen splitting variable is discussed in the next section. In the subsamples resulting from splitting in the first cutpoint, splitting continues recursively in the same or other splitting variables, until no more significant parameter instability is detected or until the number of observations in a subsample falls below a given threshold. Note that, while a sudden structural change in the data, as in our artificial data example in Figure 3 can be adequately described by a single cutpoint, other patterns, such as multiple, continuous, or nonlinear changes, can be captured by sequences of splits.

2.4. Cutpoint Selection in Bradley-Terry Models

After the ℓ -th covariate was chosen for splitting, the optimal cutpoint within this variable is selected by maximizing the partitioned likelihood (i.e., the sum

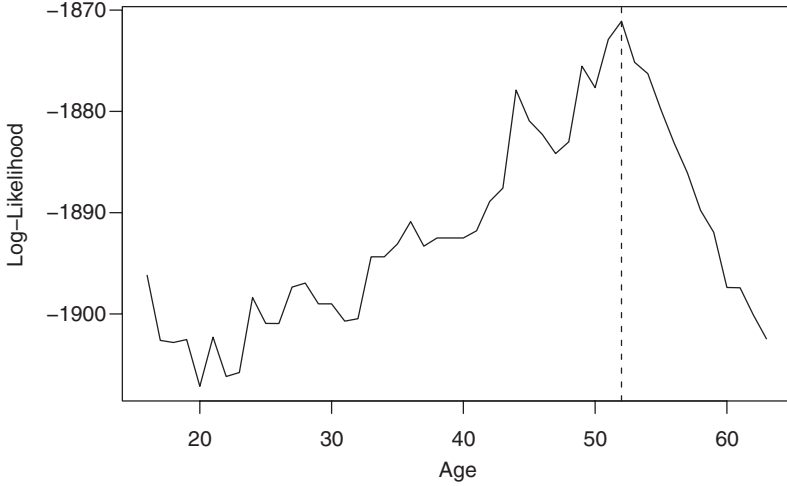


FIGURE 4. Log-likelihood of partitioned BT model for the first split in the covariate age.

of the likelihoods for the observations before and after the cutpoint) over all candidate cutpoints.

More formally, for a numeric splitting variable, we can define the subsamples $L(\xi) = \{i | x_{i\ell} \leq \xi\}$ and $R(\xi) = \{i | x_{i\ell} > \xi\}$ on the left and right, respectively, of some cutpoint ξ . For both subsamples, the parameters $\hat{\theta}^{(L)}$ and $\hat{\theta}^{(R)}$ can be estimated as described above. To determine the optimal cutpoint ξ , the partitioned likelihood

$$\sum_{i \in L(\xi)} \Psi(y_i, \hat{\theta}^{(L)}) + \sum_{i \in R(\xi)} \Psi(y_i, \hat{\theta}^{(R)})$$

is maximized over all candidate cutpoints ξ (typically requiring a certain minimal subsample size) as illustrated in Figure 4.

While this approach can be applied to numeric and ordered covariates, it is inappropriate for unordered categorical covariates. Instead, the Q categories of an unordered categorical covariate can be split into any two groups. From all these candidate binary partitions, again the one with the maximal partitioned likelihood is chosen. (Note that, in principle, the partitioning idea is not limited to binary splits—however, binary splits are typically found convenient in practice. See Zeileis et al., 2008 for strategies to compute multi-way splits.)

For the attractiveness data example, Figure 4 depicts the partitioned log-likelihood for all candidate cutpoints within the range of the numeric covariate *age*, which was selected for the first split. Note that the sharp kink in the log-likelihood in Figure 4—just like the sharp kink in the cumulative sum of scores

in Figure 3 (right)—indicates a strong instability in at least one of the model parameters over the range of the variable *age*.

The maximum of the log-likelihood in Figure 4 is achieved for $\xi = 52$. Accordingly, this value is used as the cutpoint and the sample is split into two subsamples with $age \leq 52$ and $age > 52$, as displayed in Figure 2.

Within the subsamples created by this first split, splitting is again repeated recursively as illustrated in Figure 2. However, for the following splitting variables *q2* and *gender* no cutpoint selection is necessary because there are only two subsamples associated with the two categories of both variables. Thus, the cutpoint is already determined by the selection of these covariates for splitting.

This concludes the discussion of the recursive partitioning procedure for BT models: The four basic steps—(a) BT model estimation, (b) parameter instability tests for splitting variable selection, (c) maximization of the segmented likelihood for cutpoint selection, and (d) sample splitting—are repeated recursively until there are no more significant instabilities or the subsample size is too small to consider further splitting. Note that the significance level and minimal subsample size required for further splitting need to be defined by the researcher. While in most cases the common significance level of 5% will be appropriate, lower values should be chosen when the overall sample size is very large in order to avoid growing too complex trees that may induce overfitting. The minimal subsample-size, on the other hand, should be chosen such as to provide a sufficient basis for parameter estimation in each subsample and should thus be increased when the number of stimuli becomes large.

For the application examples presented so far and in the following, which have moderate sample sizes and numbers of stimuli, a significance level of 5% and a minimal subsample size of 5 subjects were employed.

3. Application Examples

For the illustrative attractiveness data example already presented in the previous section, a more thorough discussion is provided here to highlight the straightforward interpretability of the tree-structured model. Additionally, the BT tree method is applied to a well-known data set from the field of education: Following Dittrich et al. (1998) and Böckenholt (2001b), we investigate which covariates influence business students' choice of a university for studying abroad.

3.1. Attractiveness Data

The model-based partitioning procedure for the attractiveness data was outlined in the previous section, with the resulting tree displayed in Figure 2. This data set is particularly useful for illustrating the BT tree method because, in addition to binary covariates with only a single potential cutpoint, it contains the numeric covariate *age*. As emphasized above, one important advantage of model-based partitioning for including subject-covariate information in BT

TABLE 2

Estimates of Worth Parameters in Terminal Nodes From the Partitioned Paired Comparison Model for the Germany's Next Topmodel 2007 Data

	Barbara	Anni	Hana	Fiona	Mandy	Anja
Node 3	.19	.17	.39	.11	.09	.05
Node 5	.17	.12	.26	.23	.10	.11
Node 6	.27	.21	.16	.19	.06	.10
Node 7	.26	.06	.15	.16	.16	.21

models is that such a numeric covariate does not need to be discretized in advance but can be directly included in the analysis, where an appropriate cut-point is selected in a data-driven way.

In addition to the graphical representation of the partitioned model in Figure 2, the results can be summarized by reporting the worth-parameter estimates (scaled to sum to unity) in each subsample, as in Table 2. These show that the rating of those subjects up to age 52 who watched the show on a regular basis (Node 3) essentially conforms with the assessment of the jury—except for the rating of the candidate Hana, who was judged by viewers of the show to be about twice as attractive as Barbara, the actual winner. This extreme preference for Hana cannot be found in any of the groups who did not watch the show on a regular basis. Of the subjects up to age 52 who did not watch the show on a regular basis, males (Node 5) have preferences for Hana and Fiona, while females (Node 6) rank Barbara highest, followed by Anni and Fiona.

Interestingly, the preferences of older participants (Node 7) are completely different from all other groups: Unlike the other groups, subjects over 52 judged Anja to be almost as attractive as Barbara, while they strongly dislike Anni (her attractiveness scale value is only about 20% of Barbara's). In addition to that, this group shows almost no discrimination between Hana, Fiona, and Mandy.

While the latter finding supports the common perception that the ideal of beauty varies between generations, the fact that those subjects who regularly watched the show have such a strong preference for one candidate may indicate that the candidates' personality, rather than their physical appearance, can be crucial for the audience's appreciation of candidates in casting shows.

3.2. CEMS University Choice Data

Students of the WU Wirtschaftsuniversität Wien can spend part of their study abroad, visiting one of currently 17 CEMS (Community of European Management Schools and International Companies) universities. Dittrich et al. (1998) conduct and analyze a survey of $n = 303$ first-year students to examine the students' preferences for $k = 6$ CEMS universities located in different European

TABLE 3

Observed Frequencies of Comparisons for the Community of European Management Schools and International Companies (CEMS) University Choice Data

	>	=	<	Missing
London:Paris	186	26	91	0
London:Milano	221	26	56	0
Paris:Milano	121	32	59	91
London:St. Gallen	208	22	73	0
Paris:St. Gallen	165	19	119	0
Milano:St. Gallen	135	28	140	0
London:Barcelona	217	19	67	0
Paris:Barcelona	157	37	109	0
Milano:Barcelona	104	67	132	0
St. Gallen:Barcelona	144	25	134	0
London:Stockholm	250	19	34	0
Paris:Stockholm	203	30	70	0
Milano:Stockholm	157	46	100	0
St. Gallen:Stockholm	155	50	98	0
Barcelona:Stockholm	172	41	90	0

cities: London School of Economics, Haut Etudes Commerciales (Paris), Università Commerciale Luigi Bocconi (Milan), Universität St. Gallen, Escueala Superior de Administracion y Direccion de Empresas (Barcelona), and Handelshögskolan i Stockholm. To identify reasons for the students' preferences, several subject-specific covariates have been assessed as well.

The covariates included in the analysis are *gender*, *major field of study*, and indicators of whether the students *work full time*, aim for an *international degree*, and have good skills in *French*, *Spanish*, and *Italian*. The aggregate observed frequencies $n_{j|c}$ for the $k^* = 15$ possible comparisons are listed in Table 3. Note that in this examples ties are possible if a subject is undecided between two stimuli. For 91 subjects, the comparison Paris:Milan has unintentionally been left out (Dittrich et al., 1998).

To assess the influence of the subject-specific covariates, the paired comparison model is recursively partitioned. Figure 5 shows the resulting tree. The corresponding worth-parameter estimates (scaled to sum to unity) in each of the subsamples are displayed in Table 4.

The results show that the preference scaling of the universities highly depends on the subject covariates: While for the majority of students, London is the most appealing option, students with good skills in *Italian* and *Spanish* (Node 3) have the strongest preference for Barcelona (more than twice as strong as for London), students with good skills in *Italian* but not in *Spanish* (Node 4) have a preference for Milan that is almost as pronounced as that for London, and for students with

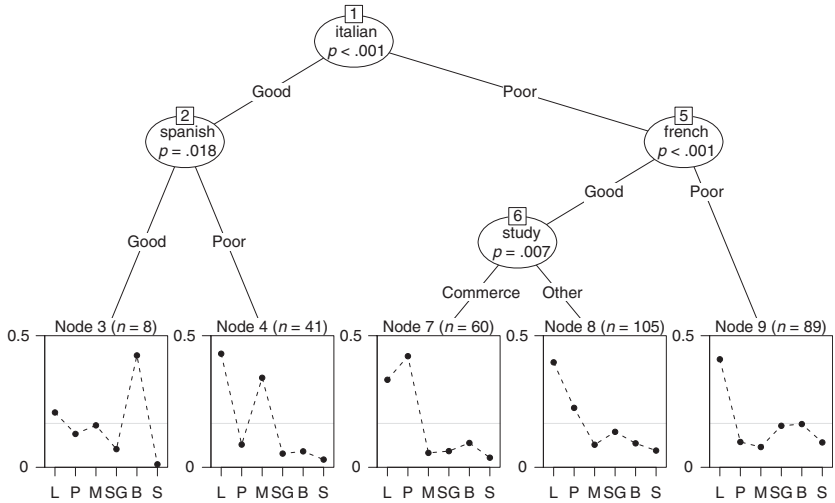


FIGURE 5. Partitioned paired comparison model for the Community of European Management Schools and International Companies (CEMS) university choice data (L: London, P: Paris, M: Milan, SG: St. Gallen, B: Barcelona, and S: Stockholm).

TABLE 4

Estimates of Worth Parameters in Terminal Nodes From the Partitioned Paired Comparison Model for the Community of European Management Schools and International Companies (CEMS) University Choice Data

	London	Paris	Milano	St. Gallen	Barcelona	Stockholm
Node 3	.21	.13	.16	.07	.43	.01
Node 4	.43	.09	.34	.05	.06	.03
Node 7	.33	.42	.05	.06	.09	.04
Node 8	.40	.23	.09	.13	.09	.06
Node 9	.41	.10	.08	.16	.16	.09

poor skills in *Italian* but good skills in *French*, the preference depends on the students' major field of study: Those students with an emphasis on commerce (Node 7) have a preference for Paris, which has a high reputation in this field, while the remaining students share the preference for London, most likely due to the fact that all Austrian university students have been exposed to several years of English language training.

Interestingly, Dittrich et al. (1998) point out that the low preference for Stockholm throughout the entire sample is most likely due to the fact that most students believe that lectures at Handelshögskolan i Stockholm are held in Swedish—while in fact they are held in English, too.

4. Discussion and Comparison to Existing Methods

Our results for the two data examples illustrate that the model-based partitioning approach for incorporating subject covariates in BT models is well suited for identifying groups of subjects with common preference scales. All covariates found relevant for partitioning in the university choice data example were also included in the model of Dittrich et al. (1998) for the same data set. However, the graphical representation as a tree makes the fitted models more accessible and intuitive to interpret compared to the parametric approaches of Dittrich et al. (1998) and Böckenholt (2001b).

Besides the straightforward graphical representation, the main difference between the semiparametric partitioning approach introduced here, and the fully parametric approaches of Dittrich et al. (1998) and Böckenholt (2001a, 2001b), is the specification of the influence of the covariates: The recursive partitioning approach presented here is data driven in the sense that the covariates enter the model in a non-parametric way, which allows for the selection of relevant covariates from a larger set of candidates and even leaves the functional form of the effects of the covariates unspecified.

This flexibility with respect to the functional form is often considered as the major advantage of recursive partitioning approaches (see also Strobl, Malley, & Tutz, 2009), because it allows for the detection of both nonlinear effects and interactions between the covariates. Another advantage of this approach is the natural treatment of both numeric and categorical covariates: While numeric covariates, such as age, are often discretized arbitrarily when separate models are fit to different groups of subjects (as, e.g., in Kissler & Bäuml, 2000), in recursive partitioning the optimal cutpoint for splitting a numeric covariate is automatically selected in a data-driven way.

As opposed to that, fully parametric approaches like those of Dittrich et al. (1998), who employ a loglinear model framework, and Böckenholt (2001a, 2001b), who employs a mixed-effects framework, require not only an active selection of the covariates but also a distinct choice of the functional form in which the covariates are included in the model. Hence these parametric models are particularly well suited for hypothesis-based modeling of psychological processes, while the recursive partitioning approach can also be applied when no or only partial information on the influence of a variety of potential covariates is available.

A different class of methods, which shares the goal of identifying groups of subjects with homogenous model parameters, is latent class or mixture modeling (an overview with applications in market segmentation, including a mixture model for paired comparisons of food preferences, is given by Wedel & Kamakura, 2000). However, in latent class approaches, the covariates responsible for the heterogeneity in the model parameters are not used for partitioning the data. Only in a second step of the analysis—if at all—are the covariates used for

characterizing the latent classes. In contrast to that, the recursive partitioning approach presented here already employs the covariates when identifying the groups, so that the interpretation is straightforward. Thus, when no covariates are available, latent class analysis is a useful tool to identify groups or clusters of subjects with homogenous model parameters. However, when covariates are available, this additional information can be readily incorporated in our recursive partitioning approach.

5. Summary and Outlook

Model-based recursive partitioning is a flexible semiparametric method adopted from machine learning, which is extended to BT models for identifying groups of subjects with different latent preference scales. The method employs splits in different covariates for partitioning the subjects, relying on the well-established statistical inference framework of fluctuation tests for detecting structural change points. Advantages of the resulting BT trees for paired comparison data are that (a) they are easy to interpret by means of visualization, (b) numeric covariates do not need to be discretized in advance, but suitable cut-points are detected in a data-driven way, (c) from a potentially large number of candidate covariates those that correspond to a significant change in the model parameters are automatically detected, and (d) interactions between covariates are also included in the same way.

Future work will aim at expanding applications of model-based partitioning in psychometrics to cover extensions of the BT model including observed stimulus-covariates (Dittrich et al., 1998) and latent characteristics of the stimuli as in the elimination by aspects (EBA) model (Tversky, 1972), as well as the Rasch model (Rasch, 1960) and its extensions.

Computational Details

Our results were obtained using R 2.9.2 (R Development Core Team, 2009) using the package `psychotree` 0.9-0 (Zeileis, Strobl, & Wickelmaier, 2009), which implements BT trees as introduced in this article. The package also contains the data for the attractiveness and the university choice examples. It relies on package `party` 0.9-999 (Hothorn, Hornik, Strobl, & Zeileis, 2009) for recursive partitioning. R itself and all packages used are freely available under the terms of the General Public License from the Comprehensive R Archive Network at <http://CRAN.R-project.org/>. Code for replicating our analysis is available in the `psychotree` package via `example("bttree", package = "psychotree")`.

References

Böckenholt, U. (2001a). Hierarchical modeling of paired comparison data. *Psychological Methods*, 6, 49–66.

- Böckenholt, U. (2001b). Thresholds and intransitivities in pairwise judgments: A multilevel analysis. *Journal of Educational and Behavioral Statistics*, 26, 269–282.
- Böckenholt, U. (2006). Thurstonian-based analyses: Past, present and future utilities. *Psychometrika*, 71, 615–629.
- Bradley, R. A., & Terry, M. E. (1952). Rank analysis of incomplete block designs. I. The method of paired comparisons. *Biometrika*, 39, 324–345.
- Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. (1984). *Classification and regression trees*. New York, NY: Chapman and Hall.
- Choisel, S., & Wickelmaier, F. (2007). Evaluation of multichannel reproduced sound: Scaling auditory attributes underlying listener preference. *Journal of the Acoustical Society of America*, 121, 388–400.
- Critchlow, D. E., & Fligner, M. A. (1991). Paired comparison, triple comparison, and ranking experiments as generalized linear models, and their implementation on GLIM. *Psychometrika*, 56, 517–533.
- Davidson, R. R. (1970). On extending the Bradley-Terry model to accommodate ties in paired comparison experiments. *Journal of the American Statistical Association*, 65, 317–328.
- Davidson, R. R., & Farquhar, P. H. (1976). A bibliography on the method of paired comparisons. *Biometrics*, 32, 241–252.
- Dittrich, R., Francis, B., Hatzinger, R., & Katzenbeisser, W. (2006). Modelling dependency in multivariate paired comparisons: A log-linear approach. *Mathematical Social Sciences*, 52, 197–209.
- Dittrich, R., Hatzinger, R., & Katzenbeisser, W. (1998). Modelling the effect of subject-specific covariates in paired comparison studies with an application to university rankings. *Journal of the Royal Statistical Society C*, 47, 511–525.
- Hothorn, T., Hornik, K., Strobl, C., & Zeileis, A. (2009). *party: A Laboratory for Recursive Partytioning*. R package version 0.9-999. Retrieved from <http://CRAN.R-project.org/package=party>
- Kissler, J., & Bäuml, K. H. (2000). Effects, of the beholder's age on the perception of facial attractiveness. *Acta Psychologica*, 104, 145–166.
- Luce, R. D. (1959). *Individual choice behavior: A theoretical analysis*. New York, NY: John Wiley.
- Matthews, J. N. S., & Morris, K. P. (1995). An application of Bradley-Terry-type models to the measurement of pain. *Journal of the Royal Statistical Society C*, 44, 243–255.
- McGuire, D. P., & Davison, M. L. (1991). Testing group differences in paired comparisons data. *Psychological Bulletin*, 110, 171–182.
- Oberfeld, D., Hecht, H., Allendorf, U., & Wickelmaier, F. (2009). Ambient lighting modifies the flavor of wine. *Journal of Sensory Studies*, 24, 797–832.
- Rasch, G. (1960). *Probabilistic models for some intelligence and attainment tests*. Chicago, IL: The University of Chicago Press. (Reprinted 1980).
- R Development Core Team. (2009). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0. Retrieved from <http://www.R-project.org/>
- 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, 323–348.

- Tversky, A. (1972). Elimination by aspects: A theory of choice. *Psychological Review*, 79, 281–299.
- Wedel, M., & Kamakura, W. (2000). *Market segmentation—Conceptual and methodological foundations* (2nd ed.). Norwell, MA: Kluwer Academic Publishers.
- Zeileis, A., & Hornik, K. (2007). Generalized M-fluctuation tests for parameter instability. *Statistica Neerlandica*, 61, 488–508.
- Zeileis, A., Hothorn, T., & Hornik, K. (2008). Model-based recursive partitioning. *Journal of Computational and Graphical Statistics*, 17, 492–514.
- Zeileis, A., Strobl, C., & Wickelmaier, F. (2009). psychotree: Recursive partitioning based on psychometric models. R package version 0.9-0. Retrieved from <http://CRAN.R-project.org/package=psychotree>

Authors

CAROLIN STROBL is assistant professor at the Department of Statistics, Ludwig-Maximilians-Universität München, Ludwigstr. 33, 80539 Munich, Germany; carolin.strobl@stat.uni-muenchen.de. Her research interests include the development and application of psychometric models and machine learning techniques.

FLORIAN WICKELMAIER is assistant professor at the Department of Psychology, Universität Tübingen, Friedrichstr. 21, 72072 Tübingen, Germany; florian.wickelmaier@uni-tuebingen.de. His research focuses on psychoacoustics, cognitive modeling, and psychophysical scaling.

ACHIM ZEILEIS is professor at the Department of Statistics, Universität Innsbruck, Universitätsstr. 15, 6020 Innsbruck, Austria; Achim.Zeileis@R-project.org. His research interests include statistical computing, applied statistics in economics and the social sciences, and statistical learning.

Manuscript received April 23, 2009

Revision received September 25, 2009

Accepted October 30, 2009