

Models for Paired Comparison Data: A Review with Emphasis on Dependent Data

Manuela Cattelan

Abstract. Thurstonian and Bradley–Terry models are the most commonly applied models in the analysis of paired comparison data. Since their introduction, numerous developments have been proposed in different areas. This paper provides an updated overview of these extensions, including how to account for object- and subject-specific covariates and how to deal with ordinal paired comparison data. Special emphasis is given to models for dependent comparisons. Although these models are more realistic, their use is complicated by numerical difficulties. We therefore concentrate on implementation issues. In particular, a pairwise likelihood approach is explored for models for dependent paired comparison data, and a simulation study is carried out to compare the performance of maximum pairwise likelihood with other limited information estimation methods. The methodology is illustrated throughout using a real data set about university paired comparisons performed by students.

Key words and phrases: Bradley–Terry model, limited information estimation, paired comparisons, pairwise likelihood, Thurstonian models.

1. INTRODUCTION

Paired comparison data originate from the comparison of objects in couples. This type of data arises in numerous contexts, especially when the judgment of a person is involved. Indeed, it is easier for people to compare pairs of objects than ranking a list of items. There are other situations that may be regarded as comparisons from which a winner and a loser can be identified without the presence of a judge. Both these instances can be analyzed by the techniques described in this paper.

The objects involved in the paired comparisons can be beverages, carbon typewriter ribbons, lotteries, players, moral values, physical stimuli and many more. Here, the elements that are compared are called objects or sometimes stimuli. The paired comparisons can be performed by a person, an agent, a consumer,

a judge, et cetera, so the terms subject or judge will be employed to denote the person that makes the choice.

The bibliography by Davidson and Farquhar (1976), which includes more than 350 papers related to paired comparison data, testifies to the widespread interest in this type of data. This interest is still present and extensions of models for paired comparison data have been proposed. This paper focuses on recent extensions of the two traditional models, the Thurstone (1927) and the Bradley–Terry (Bradley and Terry, 1952) model, especially those subsequent to the review by Bradley (1976) and the monograph by David (1988), including in particular the work that has been done in the statistical and the psychometric literature.

Section 2 reviews models for independent data. After the introduction of the two classical models for the analysis of paired comparison data and a survey of different areas of application, Sections 2.3 and 2.4 review extensions for ordinal paired comparison data and for inclusion of explanatory variables. Section 3 reviews models that allow for dependence among the observations and outlines the inferential problems related to such an extension. Here, a pairwise likelihood

Manuela Cattelan is Postdoctoral Research Fellow, Department of Statistical Sciences, University of Padua, via C. Battisti 241, 35121 Padova, Italy (e-mail: manuela.cattelan@unipd.it).

approach is proposed to estimate these models, and a simulation study is performed in order to compare the estimates produced by maximum likelihood, a common type of limited information estimation and pairwise likelihood. Section 4 reviews existing R (R Development Core Team, 2011) packages for the statistical analysis of paired comparison data, and Section 5 concludes.

2. INDEPENDENT DATA

2.1 Traditional Models

Let Y_{sij} denote the random variable associated with the result of the paired comparison between objects i and j , $j > i = 1, \dots, n$, made by subject $s = 1, \dots, S$, and let $\mathbf{Y}_s = (Y_{s12}, \dots, Y_{s(n-1)n})$ be the vector of the results of all paired comparisons made by subject s . When $S = 1$ or the difference between judges is not accounted for in the model, then the subscript s will be dropped. If each possible paired comparison is performed, they number $N = n(n - 1)/2$, and $SN = Sn(n - 1)/2$ in a multiple judgment sampling scheme, that is, when all paired comparisons are made by all S subjects. Different sampling schemes are possible. When each paired comparison is performed by a different subject, the outcomes are independent. In other instances, a subject performs more than one paired comparison; in this case, it is conceivable that results of several paired comparisons performed by the same subject will not be independent. In Section 2, independence among observations is assumed while Section 3 addresses the issue of dependent data, assuming that each subject performs all N paired comparisons, except for Section 3.3 which considers the case of dependence not induced by judges.

Let $\mu_i \in \mathbf{R}$, $i = 1, \dots, n$, denote the notional worth of the objects. Traditional models were developed assuming only two possible outcomes of each comparison, so Y_{ij} is a binary random variable, and π_{ij} , the probability that object i is preferred to object j , depends on the difference between the worth of the two objects

$$(2.1) \quad \pi_{ij} = F(\mu_i - \mu_j),$$

where F is the cumulative distribution function of a zero-symmetric random variable. Such models are called linear models by David (1988). When F is the normal cumulative distribution function, formula (2.1) defines the Thurstone (1927) model, while if F is the logistic cumulative distribution function, then the

Bradley–Terry model (Bradley and Terry, 1952) is recovered. Other specifications are possible; for example, Stern (1990) suggests modeling the worth parameters as independent gamma variables with the same shape parameter and different scale parameter. The Thurstone model is also known as the Thurstone–Mosteller model since Mosteller (1951) presented some inferential techniques for the model, while the Bradley–Terry model was independently proposed also by Zermelo (1929) and Ford (1957). Model (2.1) is called unstructured model, and the aim of the analysis is to make inference on the vector $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)'$ of worth parameters which can be used to determine a final ranking of all the objects compared. Note that the specification of model (2.1), through all the pairwise differences $\mu_i - \mu_j$, implies that a constraint is needed in order to identify the parameters. Various constraints can be specified: the most common are the sum constraint, $\sum_{i=1}^n \mu_i = 0$, and the reference object constraint, $\mu_i = 0$ for one object $i \in \{1, \dots, n\}$.

The comparative nature of the data poses inferential and interpretational problems. Consider two different studies, for example, about beverages. If subjects were requested to express an absolute measure of like/dislike for each drink in a categorical scale, then the data obtained from the two studies might be analyzed all together. On the contrary, if the subjects express preferences in paired comparisons, the data can be combined only if at least one object is common to both studies; otherwise the data can be analyzed separately, and no conclusions can be made about relationships between objects in the two different studies. Indeed, the lack of origin implies that no absolute statement can be made about the data and two subjects can provide the same sets of preferences, but one may dislike all items while the other may like all of them. The identification of an origin may be useful for understanding the underlying psychological process, for discriminating between desirable and undesirable objects and for identifying the degree of an option desirability in different conditions. However, it is not possible to recover the origin without further choice experiments and/or further assumptions (Thurstone and Jones, 1957; Böckenholt, 2004). Despite all their limits, paired comparison data are widespread because of their ease of performance and their discriminatory ability since objects that may be judged in the same like/dislike category may be differentiated when compared pairwise.

If the reference object constraint is employed, the identified worth parameters are differences with respect to the reference object. Hence, inference will typ-

ically regard differences between estimated worth parameters with the related statistical problems. For example, for testing $H_0: \mu_i = \mu_j$ by means of the Wald test statistic $(\hat{\mu}_i - \hat{\mu}_j) / \{\widehat{\text{var}}(\hat{\mu}_i - \hat{\mu}_j)\}^{1/2}$, where $\hat{\mu}_i$ is the maximum likelihood estimator of μ_i , the covariance between the estimators of the worth parameters is needed. In general, the whole covariance matrix of the worth parameters should be reported in order to allow the final users to perform the tests they are interested in. However, it is very inconvenient to report that matrix and a useful alternative may be to report quasi-standard errors (Firth and de Menezes, 2004) instead of the usual standard errors since they allow approximate inference on any of the contrasts. Let \mathbf{c} be a vector of zero-sum constants. If the parameters $\boldsymbol{\mu}$ were independent, then the estimated standard error of $\mathbf{c}'\boldsymbol{\mu}$ would be $(\sum_{i=1}^n c_i^2 \hat{v}_i)^{1/2}$, where \hat{v}_i denotes the estimated variance of $\hat{\mu}_i$. Quasi-variances are a vector of constants \mathbf{q} such that

$$\text{var}(\mathbf{c}'\boldsymbol{\mu}) \simeq \sum_{i=1}^n c_i^2 q_i,$$

so they have the property that they add over the components of $\boldsymbol{\mu}$, and hence can be used to approximate variances of contrasts of estimated worth parameters as if they were independent. Let $p(q_i + q_j, \widehat{\text{var}}(\hat{\mu}_i - \hat{\mu}_j))$, be a penalty function which depends on the quasi-variances and the estimated variance of the difference $\hat{\mu}_i - \hat{\mu}_j$, then quasi-variances are computed through minimization of the sum of the penalty function for all contrasts; see Firth and de Menezes (2004, Section 2.1).

Further statistical problems arising from the comparative nature of the data are discussed in Section 3.2.2.

EXAMPLE. A program supported by the European Union offers an international degree in Economics and Management. Twelve universities take part in this program, and in order to receive a degree, a student in the program must spend a semester at another university taking part in the program. Usually, some universities receive more preferences than others, and this may cause organizational problems. A study was carried out among 303 students of the Vienna University of Economics who were asked in which university they would prefer to spend the period abroad, between six universities situated in Barcelona (Escuela Superior de Administracion y Direccion de Empresas), London (London School of Economics and Political Sciences), Milan (Università Luigi Bocconi), Paris (Hautes Études Commerciales), St. Gallen (Hochschule St. Gallen) and

TABLE 1
Universities paired comparison data. 1 and 2 refer to the number of choices in favor of the university in the first and the second column, respectively, while X denotes the number of no preferences expressed

		1	X	2
London	Paris	186	26	91
London	Milan	221	26	56
Paris	Milan	121	32	59
London	St. Gallen	208	22	73
Paris	St. Gallen	165	19	119
Milan	St. Gallen	135	28	140
London	Barcelona	217	19	67
Paris	Barcelona	157	37	109
Milan	Barcelona	104	67	132
St. Gallen	Barcelona	144	25	134
London	Stockholm	250	19	34
Paris	Stockholm	203	30	70
Milan	Stockholm	157	46	100
St. Gallen	Stockholm	155	50	98
Barcelona	Stockholm	172	41	90

Stockholm (Stockholm School of Economics), compared pairwise. This example will be used throughout the paper as an illustration. For an exhaustive analysis of the data refer to Dittrich, Hatzinger and Katzenbeisser (1998, 2001). The data set is available in both the `prefmod` (Hatzinger, 2010) and the `BradleyTerry2` (Turner and Firth, 2010a) R packages; see Section 4. Table 1 reports the aggregated data on the 15 paired comparisons. For example, the first row shows that in the paired comparison between London and Paris, 186 students prefer London, 91 students prefer Paris and 26 students do not have a preference between the two universities. Moreover, 91 students unintentionally overlooked the comparison between Paris and Milan which has only 212 answers. The second column of Table 2 shows the estimate of the worth parameters for the six universities using the Thurstone model and adding half of the number of no preferences to each university in the paired comparison. In Section 2.3 a better way to handle no preference data will be discussed.

The reference object constraint is used, and the worth parameter of Stockholm is set to zero. All estimates are positive, so we can conclude that Stockholm is the least preferred university, while London is the most preferred one, followed by Paris, Barcelona, St. Gallen and Milan. The estimated probability that London is preferred to Paris is $\Phi(0.982 - 0.561) = 0.66$, where Φ denotes the cumulative distribution function

TABLE 2

Estimates (Est.), standard errors (S.E.) and quasi-standard errors (Q.S.E.) of the universities worth parameters employing a two-categorical Thurstone model (Thurstone) and a cumulative extension of the Thurstone model (cumulative Thurstone)

	Thurstone			cumulative Thurstone		
	Est.	S.E.	Q.S.E.	Est.	S.E.	Q.S.E.
Barcelona	0.333	0.043	0.030	0.332	0.041	0.028
London	0.982	0.045	0.033	0.998	0.043	0.031
Milan	0.240	0.044	0.031	0.241	0.041	0.029
Paris	0.561	0.044	0.031	0.566	0.042	0.030
St. Gallen	0.325	0.043	0.030	0.324	0.040	0.028
Stockholm	0	–	0.031	0	–	0.029
τ_2	–	–	–	0.153	0.007	–

of a standard normal random variable. If it is of interest to test whether the worth of St. Gallen is significantly higher than the worth of Milan, the standard error of the difference between these two worth parameters can be approximated by means of the quasi-standard errors as $(0.030^2 + 0.031^2)^{1/2} = 0.043$. Quasi-standard errors are lower than standard errors, thus accounting for the positive covariance between parameter estimates. The value of the test statistic is $(0.325 - 0.240)/0.043 = 1.98$, which yields a p -value of 0.02; hence the hypothesis of equal worth parameters between St. Gallen and Milan is not supported by the data.

2.2 Applications

There are many different areas in which paired comparison data arise. Here, a number of recent applications are described, and further references can be found in Bradley (1976), Davidson and Farquhar (1976) and David (1988).

Despite its simplicity, the basic Bradley–Terry and Thurstone models have found a wide range of applications. Choisel and Wickelmaier (2007) analyze pairwise evaluations of sounds through a standard Bradley–Terry model, while Bäuml (1994) and Kissler and Bäuml (2000) present applications involving facial attractiveness. In Mazzucchi, Linzey and Bruning (2008) the standard Bradley–Terry model is applied to a reliability problem. A panel of wiring experts is asked to state which is the riskier one between different scenarios compared pairwise in order to determine the probability of wire failure as a function of influencing factors in an aircraft environment. Stigler (1994) uses the traditional Bradley–Terry model for ranking scientific journals, and the same model is exploited in genetics by Sham and Curtis (1995).

Maydeu-Olivares and Böckenholt (2008) list 10 reasons to use Thurstone’s model for analyzing subjective health outcomes, including the ease for respondents, the existence of extensions for modeling inconsistent choices and for including covariates and the possibility to investigate which aspects influence the choices of subjects.

In many applications there are more than two possible outcomes of the comparisons. Henery (1992) employs a Thurstone model for ranking chess players and adapts it to three possible results: win, draw and loss. Böckenholt and Dillon (1997a) consider a five-response-categories model for applications to taste testing of beverages and to preferences for brands of cigarettes. Dittrich, Hatzinger and Katzenbeisser (2004) consider motives to start a Ph.D. program using three response categories in the log-linear version of the Bradley–Terry model.

It is often of interest to investigate whether some covariates affect the results of the comparisons. Ellermeier, Mader and Daniel (2004) employ a Bradley–Terry model to analyze pairwise evaluations of sounds and include sound-related covariates, for example, roughness, sharpness, et cetera, to evaluate which of them contribute to the unpleasantness of sounds. Duineveld, Arents and King (2000) use the log-linear formulation of the Bradley–Terry model to investigate consumer preference data on orange soft drinks including an analysis of the factorial design for the drinks compared, while Francis et al. (2002) include subject-specific covariates in the analysis of value orientation of people in different European countries. Applications of the Bradley–Terry model are present also in zoological data in order to investigate aspects of animal behavior considering animal-specific covariates (Stuart-Fox et al., 2006; Whiting et al., 2006; Head et al., 2008). Agresti [(2002), Chapter 10] extends the Bradley–Terry model to account for the home advantage effect in baseball data.

Sometimes it is more realistic to include dependence among observations. Object-specific random effects can be used to introduce correlation between comparisons with common objects, for example, in sports data (Cattelan, 2009). When all judges perform all paired comparisons, random effects can introduce correlation between preferences expressed by the same subject involving a common object as shown in Böckenholt and Tsai (2007) for the university preference data.

When paired comparisons are performed in prolonged time periods, it may be necessary to account

for it. McHale and Morton (2011) estimate a Bradley–Terry model in which tennis matches distant in time are down-weighted since the aim is to predict the results of future matches. Further dynamic extensions for sports data have been proposed by Barry and Hartigan (1993), Fahrmeir and Tutz (1994), Knorr-Held (2000) and Cattelan, Varin and Firth (2012). In tournaments it may happen that a player wins all the comparisons in which he is involved. In this case a standard Bradley–Terry or Thurstone model would estimate an infinity worth parameter for this team. Mease (2003) proposes a penalization of the likelihood which overcomes this problem. The method proposed by Firth (1993) to reduce the bias of the maximum likelihood estimates is an alternative technique to obtain finite estimates in this instance. Finally, the case in which the margin of victory in sport contests is not discrete, but continuous, is analyzed in Stern (2011).

In the context of the log-linear specification of the Bradley–Terry model, Dittrich et al. (2012) account also for missing responses in a study about the qualities of a good teacher.

2.3 Ordinal Paired Comparisons

Sometimes subjects are requested to express a degree of preference. Suppose that objects i and j are compared, and the subject can express strong preference for i over j , mild preference for i , no preference, mild preference for j over i or strong preference for j . If H denotes the number of grades of the scale, then in this example, $H = 5$.

Let $Y_{ij} = 1, \dots, H$, where 1 denotes the least favorable response for i , and H is the most favorable response for i . Agresti (1992) shows how two models for the analysis of ordinal data can be adapted to ordinal paired comparison data. The *cumulative link* models exploit the latent random variable representation. Let Z_{ij} be a continuous latent random variable, and let $\tau_1 < \tau_2 < \dots < \tau_{H-1}$ denote thresholds such that $Y_{ij} = h$ when $\tau_{h-1} < Z_{ij} \leq \tau_h$. Then,

$$(2.2) \quad \text{pr}(Y_{ij} \leq y_{ij}) = F(\tau_{y_{ij}} - \mu_i + \mu_j),$$

where $-\infty = \tau_0 < \tau_1 < \dots < \tau_{H-1} < \tau_H = \infty$, and F is the cumulative distribution function of the latent variable Z_{ij} . F is usually assumed to be either the logistic or the normal distribution function leading to the cumulative logit or the cumulative probit model, respectively. The symmetry of the model imposes that $\tau_h = -\tau_{H-h}$, $h = 1, \dots, H$ and $\tau_{H/2} = 0$ when H is even. When $H = 3$ there are two threshold parameters,

τ_1 and τ_2 , such that $\tau_1 = -\tau_2$ and model (2.2) corresponds to the extension of the Bradley–Terry model introduced by Rao and Kupper (1967) when a logit link is considered, and the extension of the Thurstone model by Glenn and David (1960) when the probit link is employed.

An alternative model proposed by Agresti (1992) is the *adjacent categories* model. In this case the link is applied to adjacent response probabilities, rather than cumulative probabilities and reduces to the Bradley–Terry model when only 2 categories are allowed and to the model proposed by Davidson (1970) when 3 categories are allowed. The adjacent categories model is simpler to interpret than cumulative link models since the odds ratio refers to a given outcome instead of referring to groupings of outcomes (Agresti, 1992). The adjacent categories model, as well as the Bradley–Terry model, has also a log-linear representation (Dittrich, Hatzinger and Katzenbeisser, 2004).

An application of the adjacent categories model to market data is illustrated in Böckenholt and Dillon (1997b). Böckenholt and Dillon (1997a) note that a bias may be caused by the usage of the scale because subjects may use only subsets of all categories. The threshold parameters τ_h can account for the selection bias, for example, in the cumulative probit model the quantity $\Phi(\tau_h) - \Phi(\tau_{h-1})$ gives the category selection bias since it is the probability of selecting category h when the two stimuli are equal. Different latent classes of consumers with different threshold values and worth parameters can be identified. If subjects share the same worth parameters but have different thresholds, it is possible to let thresholds depend on subject-specific covariates and to have a random part (Böckenholt, 2001b). It is also possible to define thresholds that depend on the objects compared, as in Henery (1992).

EXAMPLE. In the paired comparisons of universities, students were allowed to express no preference between two universities. Therefore, the data should be analyzed by means of a model for ordinal data. Columns 5–7 in Table 2 show the estimates of a cumulative probit extension of the Thurstone model for the university data. The estimated threshold parameter $\hat{\tau}_2 = 0.153$ is highly significant. In this particular case, the estimates of the worth parameters and their standard errors are very similar to those of the model with two categories, and the ranking of universities remains the same, but in general, especially when the

number of no preferences is large, results can be different. Moreover, in this case it is possible to estimate the probability of no preference between London and Paris which is $\Phi(0.153 - 0.998 + 0.566) - \Phi(-0.153 - 0.998 + 0.566) = 0.11$, and the estimated probability that London is preferred to Paris reduces to $1 - \Phi(0.153 - 0.998 + 0.566) = 0.61$; hence the estimated probability that Paris is preferred to London is 0.28. There is no much difference from the previous result in the test of equality of worth parameters for universities in St. Gallen and Milan.

2.4 Explanatory Variables

In many instances, it is of interest to investigate whether some explanatory variables affect the results of the comparisons. Explanatory variables can be related to the objects compared, to the subjects performing the comparisons or they can be comparison-specific.

Let $\mathbf{x}_i = (x_{i1}, \dots, x_{iP})'$ be a vector of P explanatory variables related to object i and $\boldsymbol{\beta} = (\beta_1, \dots, \beta_P)$ be a P -dimensional parameter vector. Then, in the context of the Bradley–Terry model, Springall (1973) proposes to describe the worth parameters as the linear combination

$$(2.3) \quad \mu_i = x_{i1}\beta_1 + \dots + x_{iP}\beta_P, \quad i = 1, \dots, n.$$

A paired comparison model with explanatory variables is called a structured model. The same extension can be applied to the Thurstone model. Note that since only the differences $\mu_i - \mu_j = (\mathbf{x}_i - \mathbf{x}_j)'\boldsymbol{\beta}$ enter the linear predictor, an intercept cannot be identified. In some instances, both worth parameters of objects and further object-specific covariates are included, hence the linear predictor assumes the form $\mu_i - \mu_j + (\mathbf{x}_i - \mathbf{x}_j)'\boldsymbol{\beta}$; see Stern (2011).

Model (2.3) has been extended to more flexible models, such as additive combinations of spline smoothers (De Soete and Winsberg, 1993); however large data sets may be necessary to estimate nonlinearities reliably, even though there is no investigation about this issue.

In case worth parameters are specified as in (2.3), standard errors for the worth parameters can be computed through the delta method, while when both the worth parameters and covariates are included in the linear predictor, quasi-standard errors can be computed for the worth parameters.

The results of the comparisons can be influenced also by characteristics of the subject that performs the paired comparisons. In the log-linear representation

of the Bradley–Terry model, Dittrich, Hatzinger and Katzenbeisser (1998) show how to include categorical subject specific covariates, while Francis et al. (2002) tackle the problem of continuous subject-specific covariates and consider also the case in which some of these covariates have a smooth nonlinear relationship.

Dillon, Kumar and De Borrero (1993) consider a marketing application and divide subjects in latent classes to which they belong with a probability that depends on their explanatory variables.

Covariates can be added in the linear predictor (2.3) if they are subject-object interaction effects. For example, the knowledge of a foreign language may influence the preference for a university. An interaction effect can account for whether the student knows, for example, Spanish and one object in the comparison is the university in Barcelona. Unfortunately, subject covariates that do not interact with objects, such as age of respondents, cannot be included.

A semiparametric approach which accounts for subject-specific covariates is proposed by Strobl, Wickelmaier and Zeileis (2011) who suggest a methodology to partition recursively the subjects that perform the paired comparisons on the basis of their covariates. The procedure tests whether structural changes in the parameters occur for subjects with different values of the covariates. Subjects are split according to the test and a different unstructured Bradley–Terry model is fitted for each subgroup. The method allows us to identify which covariates influence the worth parameters without the need to assume a model for them and finds the best cut point in case of continuous covariates. Moreover, it is possible to include subject-specific covariates, not only interaction effects. Attention is needed in setting the minimum number of subjects per class and in setting the significance level of the test in order to avoid overfitting for large data sets. Differently from the usual latent class models, the method allows to divide subjects on the basis of their covariates; however, if some important subject-specific covariates are not available, it may be expected that the usual latent class model will perform better. In Strobl, Wickelmaier and Zeileis (2011) an unstructured Bradley–Terry model is estimated for each subgroup, but it seems possible to extend the method also to structured models.

Finally, there may be also comparison-specific covariates which are related to the objects, but change from comparison to comparison. An example of a comparison-specific covariate is the home advantage effect in sport tournaments since it depends on whether one of the players competes in the home field. This

effect may be accounted for by adding a further term in the linear predictor (2.3). Another example is the experience effect in contests between animals which, in [Stuart-Fox et al. \(2006\)](#), is accounted for through a covariate that counts the number of previous contests fought by animals.

EXAMPLE. In the universities paired comparisons, it may be of interest to assess whether some object-specific covariates influence the results of the comparisons. The universities in London and Milan specialize in economics, the universities in Paris and Barcellona specialize in management science and the remaining two in finance. This aspect may influence the decisions of students. Another element that may affect the comparisons is the location of the universities, in this respect they can be divided in universities in Latin countries (Italy, France and Spain) and universities in other countries.

Some features of the students that performed the universities paired comparisons were collected, too. In particular, it is known whether students have good knowledge of English, Italian, Spanish and French and which is the main topic of their studies. It is conceivable that, for example, students with a good knowledge of French are more inclined to prefer the university in Paris. [Table 3](#) shows the estimates of a model with a linear predictor that includes object specific covariates and subject-object interaction effects. Universities in non-Latin countries are preferred to those in Latin countries, and universities that specialize in finance seem less appealing to students. The good knowledge of a foreign language induces students to choose the university situated in the country where that foreign language is spoken. Consider a student with a good knowledge of both English and French and whose main

discipline of study is management, then the estimated probability that this student prefers London to Paris is $1 - \Phi\{0.160 - (0.141 + 0.757 - 0.652 - 0.789 + 0.835 - 0.238)\} = 0.46$, while the estimated probabilities of no preference and preference for Paris are 0.13 and 0.41, respectively. If this student's main discipline of study was not management, which is the subject in which Paris specializes, then the above estimated probabilities of preferring London, no preference and preferring Paris would become 0.55, 0.12 and 0.33, respectively.

3. MODELS FOR DEPENDENT DATA

3.1 Intransitive Preferences

The models presented so far are estimated assuming independence among all observations. The inclusion of a dependence structure is not only more realistic, but also has an impact on the transitivity properties of the model. Intransitive choices occur when object i is preferred to j , and object j is preferred to k , but in the paired comparison between i and k , the latter is preferred. These are also called circular triads. Paired comparison models can present different transitivity properties. Assume that $\pi_{ij} \geq 0.5$ and $\pi_{jk} \geq 0.5$, then a model satisfies:

- weak stochastic transitivity if $\pi_{ik} \geq 0.5$;
- moderate stochastic transitivity if $\pi_{ik} \geq \min(\pi_{ij}, \pi_{jk})$;
- strong stochastic transitivity if $\pi_{ik} \geq \max(\pi_{ij}, \pi_{jk})$.

The Bradley–Terry and Thurstone models as presented so far satisfy strong stochastic transitivity. This property may be desirable sometimes, for example, when asking wiring experts which is the riskier situation between different scenarios in an aircraft environment. In this case it is desirable that choices are consistent, so [Mazzucchi, Linzey and Bruning \(2008\)](#) use transitivity to check the level of reliability of experts. However, in some situations choices can be systematically intransitive, for example, when the same objects have more than one aspect of interest, and different aspects prevail in different comparisons.

[Causeur and Husson \(2005\)](#) propose a two-dimensional Bradley–Terry model in which the worth parameter of each object is bidimensional and can thus be represented on a plane. A further multidimensional extension is proposed by [Usami \(2010\)](#). However, this methodology does not provide a final ranking of all objects.

TABLE 3

Estimates (Est.) and standard errors (S.E.) of universities data with subject- and object-specific covariates

	Est.	S.E.
Economics	0.757	0.066
Management	0.789	0.080
Latin country	−0.835	0.071
Discipline:Management	0.238	0.054
English:London	0.141	0.075
French:Paris	0.652	0.049
Italian:Milan	1.004	0.094
Spanish:Barcelona	0.831	0.095
τ_2	0.160	0.007

A different method that allows the inclusion in the model even of systematic intransitive comparisons while yielding a ranking of all the objects consists of modeling the dependence structure among comparisons. The development of inferential techniques for dependent data has recently allowed an investigation of models for dependent observations.

3.2 Multiple Judgment Sampling

The assumption of independence is questioned in the case of the multiple judgment sampling, that is, when S people make all the N paired comparisons. It seems more realistic to assume that the comparisons made by the same person are dependent. This aspect has received much attention in the literature during the last decade.

3.2.1 *Thurstonian models.* The original model proposed by [Thurstone \(1927\)](#) includes correlation among the observations. The model was developed for analyzing sensorial discrimination and assumes that the stimuli $\mathbf{T} = (T_1, \dots, T_n)'$ compared in a paired comparison experiment follow a normal distribution, $\mathbf{T} \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma}_T)$, with mean $\boldsymbol{\mu} = (\mu_1, \dots, \mu_n)'$ and variance $\boldsymbol{\Sigma}_T$. [Thurstone \(1927\)](#) proposes different models with different covariance matrices of the stimuli, so the set of models which assume a normal distribution of the stimuli are called Thurstonian models. The single realization t_i of the stimulus T_i can vary, and the result of the paired comparison between the same two stimuli can be different in different occasions. Assume that only either a preference for i or a preference for j can be expressed, so then in a paired comparison when $T_i > T_j$ object i is preferred, or alternatively, when the latent random variable $Z_{ij} = T_i - T_j$ is positive, a win for i is observed; otherwise a win for j occurs. In the context of multiple judgment sampling, [Takane \(1989\)](#) proposes to include a vector of pair specific errors. Let $\mathbf{Z}_s = (Z_{s12}, \dots, Z_{sn-1n})'$ be the vector of all latent continuous random variables pertaining to subject s , then

$$(3.1) \quad \mathbf{Z}_s = \mathbf{A}\mathbf{T} + \mathbf{e}_s,$$

where $\mathbf{e}_s = (e_{s12}, e_{s13}, \dots, e_{sn-1n})'$ is the vector of pair-specific errors which has zero mean, covariance $\boldsymbol{\Omega}$ and is independent of \mathbf{T} and of $\mathbf{e}_{s'}$ for any other subject $s' \neq s$, and \mathbf{A} is the design matrix of paired comparisons whose rows identify the paired comparisons and columns correspond to the objects. For example, if $n = 4$, and the paired comparisons are

(1, 2), (1, 3), (1, 4), (2, 3), (2, 4) and (3, 4), then

$$\mathbf{A} = \begin{pmatrix} 1 & -1 & 0 & 0 \\ 1 & 0 & -1 & 0 \\ 1 & 0 & 0 & -1 \\ 0 & 1 & -1 & 0 \\ 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & -1 \end{pmatrix}.$$

A similar model is employed by [Böckenholt and Tsai \(2001\)](#), who assume that $\mathbf{e}_s \sim N(\mathbf{0}, \omega^2 \mathbf{I}_N)$. The more general analysis of covariance structure proposed by [Takane \(1989\)](#) can accommodate both the *wandering vector* and the *wandering ideal point* models ([Carroll and De Soete, 1991](#)), which are models with different assumptions about the mechanism originating the data. The wandering vector and wandering ideal point models do not impose the number of dimensions which is determined from the data alone, so they are powerful models to analyze human choice behavior and inferring perceptual dimensions.

The model thus specified is over-parametrized. To reduce the number of parameters, [Thurstone \(1927\)](#) proposes different restrictions on the covariance matrix $\boldsymbol{\Sigma}_T$, while [Takane \(1989\)](#) proposes a factor model. Nonetheless, these models with a reduced number of parameters need further identification restrictions; see Section 3.2.2.

A further extension of model (3.1) is proposed by [Tsai and Böckenholt \(2008\)](#) who unify [Tsai and Böckenholt \(2006\)](#) with [Takane \(1989\)](#) to obtain a general class of models that can account simultaneously for transitive choice behavior and systematic deviations from it. In this case the latent variable is

$$(3.2) \quad \mathbf{Z}_s = \mathbf{A}\mathbf{T} + \mathbf{B}\mathbf{V}_s,$$

where $\mathbf{V}_s = (V_{s1(2)}, V_{s1(3)}, \dots, V_{s2(1)}, V_{s2(3)}, \dots, V_{sn(n-1)})'$ is a vector of zero mean random effects designed so as to capture the random variation in judging an object when compared to another specific object, and \mathbf{B} is a matrix with rows corresponding to the paired comparisons and columns corresponding to the elements of \mathbf{V}_s , so, for example, if $n = 3$, $\mathbf{V}_s = (V_{s1(2)}, V_{s1(3)}, V_{s2(1)}, V_{s2(3)}, V_{s3(1)}, V_{s3(2)})'$ and

$$\mathbf{B} = \begin{pmatrix} 1 & 0 & -1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & 1 & 0 & -1 \end{pmatrix}.$$

It is assumed that \mathbf{V}_s , the within-judge variability, is normally distributed with mean 0 and covariance $\boldsymbol{\Sigma}_V$ so that $\mathbf{Z}_s \sim N(\mathbf{A}\boldsymbol{\mu}, \mathbf{A}\boldsymbol{\Sigma}_T\mathbf{A}' + \mathbf{B}\boldsymbol{\Sigma}_V\mathbf{B}')$.

In the remaining it will be assumed that there are only two possible outcomes of the comparisons, but it

is easy to extend this model for ordinal data through the introduction of threshold parameters with a specification analogous to (2.2).

3.2.2 Identification. Psychometricians are interested in understanding the relations between stimuli; hence they are primarily interested in the unstructured and unrestricted Thurstonian models. Unfortunately, due to the comparative nature of the data, some identification restrictions on the covariance matrix are needed. The necessary identification restrictions to estimate model (3.1) are discussed in Maydeu-Olivares (2001, 2003), Tsai and Böckenholt (2002) and Tsai (2003). Consider the covariance matrix $\Sigma_Z = \text{Cov}(\mathbf{Z}_s) = \mathbf{A}\Sigma_T\mathbf{A}' + \mathbf{\Omega}$, where Σ_T is an unrestricted covariance matrix. Because of the difference structure of the judgments Σ_T and $\Sigma_T + \mathbf{d}\mathbf{1}' + \mathbf{1}\mathbf{d}'$ where $\mathbf{1}$ is a vector of n ones and \mathbf{d} is an n -dimensional vector of constants such that the matrix remains positive definite, are not distinguishable (Tsai, 2000). Indeed, let $\mathbf{K} = [\mathbf{I}_{n-1} | -\mathbf{1}]$ be an identity matrix of dimension $n - 1$ to which a column of elements equal to -1 is added, then only $\mathbf{K}\boldsymbol{\mu}$ and $\mathbf{K}\Sigma_T\mathbf{K}'$ are identifiable. For example the matrices

$$\Sigma_{T,1} = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix},$$

$$\Sigma_{T,2} = \begin{pmatrix} 0.750 & 0.125 & 0 \\ 0.125 & 1.5 & 0.375 \\ 0 & 0.375 & 1.250 \end{pmatrix}$$

are not distinguishable because $\mathbf{K}\Sigma_{T,1}\mathbf{K}' = \mathbf{K} \cdot \Sigma_{T,2}\mathbf{K}' = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$, where the second matrix is obtained from the first one by setting $\mathbf{d} = (-1/8, 1/4, 1/8)$. This consideration remains valid for any generic matrix of contrasts that may be used instead of \mathbf{K} . The specifications of the covariance matrix Σ_T with a reduced number of parameters proposed by Thurstone (1927) cannot be recovered from the data and only covariance classes can be considered.

Tsai (2003) shows that $n + 2$ constraints are needed in order to identify model (3.1), including the constraint on the worth parameters. As for the mean parameters, many different constraints can be imposed on the covariance matrix. For example, Böckenholt and Tsai (2001), Tsai and Böckenholt (2002) and Maydeu-Olivares (2003) set all the diagonal elements of Σ_T equal to 1 and either one of the diagonal elements of $\mathbf{\Omega}$ to 1 or one of the nondiagonal elements of Σ_T equal to zero. However, if Σ_T is fixed to be a correlation matrix, the set of matrixes that produce the same sets of

probabilities is limited. Maydeu-Olivares and Böckenholt (2005) set all the covariances involving the last latent utility to zero, which corresponds to assuming independence between the last stimulus and the others, and the variance of the first and last item to one. Maydeu-Olivares and Hernández (2007) suggest to set all diagonal elements of Σ_T equal to one and the sum of the correlations between the first and the other latent variables to one. With these constraints positive entries in the correlation matrix imply that strong preference for one stimulus is associated with strong preference for the other stimulus, while negative entries indicate that strong preference for one stimulus is associated with weak preference for the other stimulus. Thus, it is not necessary to fix any element in the matrix $\mathbf{\Omega}$, since the constraint $\omega = 1$ in $\mathbf{\Omega} = \omega^2\mathbf{I}_N$ could lead to a nonpositive definite matrix Σ_T . After estimation it is possible to recover the class of covariance matrixes that produce the same probabilities (Maydeu-Olivares and Hernández, 2007). However, the initial identification constraints pose limits on the set of covariance matrixes that identify the same model.

There is no discussion or results about the identification restrictions necessary to estimate model (3.2). In order not to incur identification problems, Tsai and Böckenholt (2008) assume that the matrix Σ_V depends on very few parameters.

3.2.3 Models with logit link. The dependence between evaluations made by the same judge has been introduced also in models employing logit link functions. Different specifications have been used for this purpose.

A first inclusion of dependence in logit models is proposed by Lancaster and Quade (1983), who consider multiple judgments by the same person and introduce correlation in the Bradley–Terry model assuming that the worth parameters are random variables following a beta distribution with shape parameters a_{ij} and b_{ij} . The Bradley–Terry model is imposed on the means of the beta distributions, that is, $E(\pi_{ij}) = a_{ij}/(a_{ij} + b_{ij}) = \pi_i/(\pi_i + \pi_j)$, but such a model introduces correlation only between comparisons of the same judge on the same pair of objects, while the other comparisons remain independent. The same limit presents the extension by Matthews and Morris (1995) who consider three possible response categories.

Two different methods have been used for introducing dependence among comparisons made by the same person involving one common object in logit models. The first method exploits the usual association measure for binary data: the odds ratio. Böckenholt and

Dillon (1997a) consider the adjacent categories model for preference data with H categories and suggest a parametrisation in terms of log-odds ratios to account for dependence between observations, while Dittrich, Hatzinger and Katzenbeisser (2002) adopt a similar approach in a two-categorical model using the log-linear formulation of the Bradley–Terry model. This specification is convenient because it allows one to estimate the model through standard software developed for log-linear models, but the number of added parameters can be quite large (Dittrich, Hatzinger and Katzenbeisser, 2002).

Another method used for introducing dependence among observations is the inclusion of random effects in the linear predictor. Böckenholt (2001a) describes the worth of object i for subject s as

$$\mu_{si} = \mu_i + \sum_{p=1}^P \beta_{ip}x_{ip} + U_{si},$$

where U_{si} is a random component, and \mathbf{x}_i is a vector of P subject-specific (and possibly item specific) covariates. Böckenholt (2001a) employs a logit link function and assumes that $\mathbf{U}_s = (U_{s1}, \dots, U_{sn})'$ follows a multivariate normal distribution with mean $\mathbf{0}$ and covariance Σ_U .

Francis, Dittrich and Hatzinger (2010) consider the log-linear representation of the Bradley–Terry model and introduce random effects for each respondent in order to account for residual heterogeneity that is not included in subject-specific covariates. The inclusion of random effects in the linear predictor introduces difficulties in the estimation of the model.

3.2.4 *Choice models.* The work by Thurstone has great importance in the development of models for analyzing discrete choices, not only from a psychometric point of view, but also in economic choice theory. When the idea that choices may be random and not fixed started to develop, the use of the model proposed by Thurstone was suggested (Marschak, 1960). As the Nobel laureate McFadden (2001) states, “when the perceived stimuli are interpreted as levels of satisfaction, or utility, this can be interpreted as a model for economic choice.”

According to the economic theory, models for discrete choice are required to satisfy the utility maximization assumption which states that subjects maximize their utility when making decisions. Let Υ_{si} denote the utility of subject s from alternative i which can be decomposed as $\Upsilon_{si} = M_{si} + \varepsilon_{si}$, where M_{si} denotes a function which relates a set of alternative attributes and subject attributes to the utility gain and ε_{si}

denotes factors that affect utility, but are not included in M_{si} . The probability that subject s chooses alternative i is equal to the probability that the utility gained from i is higher than the utility from every other object in the choice set: $\text{pr}(\Upsilon_{si} > \Upsilon_{sj}, \forall i \neq j) = \text{pr}(\varepsilon_{si} - \varepsilon_{sj} < M_{si} - M_{sj}, \forall i \neq j)$. These models are called random utility models. For each person, a choice is described as $n - 1$ paired comparisons between the preferred alternative and all other options. Note that paired comparisons do not really occur, so inconsistent choices cannot be observed.

From the above specification, different models have been developed depending on the assumptions about the distribution of the errors and the formulation of the mean term M_{si} . If the ε_{si} 's are independent and follow a Gumbel distribution the choice model is a logit model and, when $M_{si} = \mathbf{x}'_{si}\boldsymbol{\beta}$, it corresponds to the structured Bradley–Terry model. A particular concern is caused by the independence from irrelevant alternatives (Luce, 1959) property which characterizes the Bradley–Terry model. Indeed, in the Bradley–Terry model the ratio between probabilities of choosing one option over another is independent from the other available alternatives. Often, this property is not satisfied in real data. This limit is somehow overcome by assuming a type of generalized extreme value distribution for the errors. In the resulting nested logistic model, independence from irrelevant alternatives holds for sets of alternatives within a same subset and not for alternatives in different subsets (Train, 2009). The advantage of these specifications is that models can be estimated easily, but they cannot account for random taste variation or unobserved factors correlated over time.

A further proposal is to assume a multivariate normal distribution for the errors ε_{si} . This model is very flexible since it allows for random taste variation and, when necessary, for temporally correlated errors, but its estimation is not straightforward. The resulting model is a multivariate probit model, like the Thurstone model. In economic choice models it is of interest to consider the influence on decisions of covariates that are included in the mean term M_{si} . Explanatory variables can be considered also in psychometric models (Tsai and Böckenholt, 2002), even though interest is focused on the parameters $\boldsymbol{\mu}$ which are always included in the linear predictor.

Other extensions include further random elements in the mean term M_{si} , so as to allow flexible disturbances or to account for different attitudes and perceptions of different people. All these elements add difficulties in the estimation of the model.

An important aspect in choice theory is the distinction between stated and revealed preferences. This problem has not received much attention in the psychometric literature, but there may be differences between what people say they would choose in a questionnaire survey and what they really choose. The former are called stated preferences and the latter revealed preferences. If both types of preferences are available, it may be useful to analyze them all together. Walker and Ben-Akiva (2002) propose a model that incorporates many of the above extensions; however, care is needed when specifying the model because it may be difficult to understand which parameters can be identified. Moreover, the inclusion of additional disturbances and unobserved covariates requires the approximation of integrals whose dimension can be high.

Random utility models are very useful and widely spread; however, some doubts have been raised about their basic assumption that people act as to maximize their utility since sometimes consumers do not make rational choices (Böckenholt, 2006).

3.3 Object-Related Dependencies

In the multiple judgment sampling the dependence among observations derives from repeated comparisons made by the same person, usually involving a common object. In case paired comparisons are not performed by a judge, the correlation may arise from the fact that the same object is involved in multiple paired comparisons. For example, when contests among animals are analyzed, it is realistic to assume that comparisons involving the same animal are correlated. In this perspective, Firth (2005) suggests to set

$$(3.3) \quad \mu_i = \mathbf{x}'_i \boldsymbol{\beta} + U_i,$$

where U_i is a zero mean object-specific random effect. This approach is investigated in Cattelan (2009). The results of comparisons are related to observed characteristics of the animal and to unobserved quantities that are captured by the random effect U_i .

In this case, the latent random variable can be written as

$$\mathbf{Z} = \mathbf{A}\mathbf{X}\boldsymbol{\beta} + \mathbf{A}\mathbf{U} + \boldsymbol{\eta},$$

where $\mathbf{U} = (U_1, \dots, U_n)$ is the vector of all object-specific random effects, \mathbf{X} is the matrix of covariates with columns \mathbf{x}_i , $\boldsymbol{\eta}$ are independent normally distributed errors with mean 0 and variance 1 while the matrix \mathbf{A} is the design matrix of the paired comparisons with rows that describe which comparisons are observed, not necessarily all possible paired comparisons.

If it is assumed that \mathbf{U} is multivariate normal with mean $\mathbf{0}$ and covariance $\mathbf{I}_n \sigma^2$, then $\mathbf{Z} \sim N(\mathbf{A}\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{A}\mathbf{A}^T + \mathbf{I}_d)$, where d is the number of paired comparisons observed. Again, this model is a multivariate probit model. However, this type of data presents some different features with respect to multiple judgment sampling. While in psychometric applications n is not very large because it is unlikely that a person will make all the paired comparisons when $n > 10$, this will typically happen in sport tournaments or in paired comparison data about animal behavior. Moreover, in the multiple judgment sampling scheme S independent replications of all the comparisons are available, but in other contexts this does not occur, adding further difficulties.

3.4 Inference

3.4.1 *Estimation.* In this section, the multiple judgment sampling scheme is mainly investigated, and only some comments are made about the case of object-related dependencies. There are different methods for estimating models for dependent paired comparison data. A first approach to the computation of the likelihood function requires to integrate out the latent variables \mathbf{T} from the joint distribution of \mathbf{Y} and \mathbf{T} . This integral has dimension n , the number of items, but rewriting it in terms of differences $T_i - T_n, i = 1, \dots, n - 1$, the dimension can be reduced to $n - 1$, which nonetheless may still be quite large when methods such as the Gauss–Hermite quadrature are employed.

Alternatively, it is possible to represent the joint distribution of the observations as a multivariate probit model. Let $\mathbf{Z}_s^* = \mathbf{D}(\mathbf{Z}_s - \mathbf{A}\boldsymbol{\mu})$ be the standardized version of the latent variable \mathbf{Z}_s , where $\mathbf{D} = [\text{diag}(\boldsymbol{\Sigma}_Z)]^{-1/2}$ and $\boldsymbol{\Sigma}_Z$ denotes the covariance matrix of \mathbf{Z}_s expressed as in model (3.1) or in model (3.2). Then, \mathbf{Z}_s^* follows a multivariate normal distribution with mean $\mathbf{0}$ and correlation matrix $\boldsymbol{\Sigma}_{Z^*} = \mathbf{D}\boldsymbol{\Sigma}_Z\mathbf{D}$. Object i is preferred to object j when $z_{sij}^* \geq \tau_{ij}^*$, where the vector of the thresholds is given by $\boldsymbol{\tau}^* = -\mathbf{D}\mathbf{A}\boldsymbol{\mu}$. The likelihood function is the product of the probability of the observations for each judge

$$\mathcal{L}(\boldsymbol{\psi}; \mathbf{Y}) = \prod_{s=1}^S \mathcal{L}_s(\boldsymbol{\psi}; \mathbf{Y}_s),$$

where

$$\mathcal{L}_s(\boldsymbol{\psi}; \mathbf{Y}_s) = \int_{R_{s12}} \dots \int_{R_{sn-1n}} \phi_N(\mathbf{z}_s^*; \boldsymbol{\Sigma}_{Z^*}) d\mathbf{z}_s^*,$$

$\phi_N(\cdot; \boldsymbol{\Sigma}_{Z^*})$ denotes the density function of an N -dimensional normal random variable with mean $\mathbf{0}$ and

correlation matrix Σ_{Z^*} and

$$R_{sij} = \begin{cases} (-\infty, \tau_{ij}^*) & \text{if } Y_{sij} = 1, \\ (\tau_{ij}^*, \infty) & \text{if } Y_{sij} = 2. \end{cases}$$

Note that this approach requires the approximation of S integrals whose dimension is equal to $N = n(n - 1)/2$, the number of paired comparisons, so its growth is quadratic with the increase in the number of objects. However, there is a large literature about methods for approximate inference in multivariate probit models. The algorithm proposed by Genz and Bretz (2002) to approximate multivariate normal probabilities is based on quasi-Monte Carlo methods, and Craig (2008) warns against the randomness of this method for likelihood evaluation. A deterministic approximation is developed by Miwa, Hayter and Kuriki (2003), but it is available only for integrals of dimension up to 20 since even for such a dimension its computation is very slow. Approximations based on Monte Carlo methods can be used (Chib and Greenberg, 1998), but they may be computationally expensive if the dimension of the integral is very large. Böckenholt and Tsai (2001) use an EM algorithm, while in econometric theory a maximum simulated likelihood approach in which multivariate normal probabilities are simulated through the Geweke–Hajivassiliou–Keane algorithm is employed (Train, 2009). A further approach may be based on data cloning (Lele, Nadeem and Schmuland, 2010). When integrals are very large, and the approximation is computationally demanding and time-consuming, it is possible to resort to limited information estimation methods, which are estimation procedures based on low dimensional margins. Here, we compare two different methods. The first one is widely applied in the context of multiple judgment sampling (Maydeu-Olivares, 2001, 2002; Maydeu-Olivares and Böckenholt, 2005) and will be called limited information estimation; the second is proposed in the context of object-specific dependencies in Cattelan (2009) and is called pairwise likelihood.

The limited information estimation procedure considered here consists of three stages. In the first stage the threshold parameters τ^* are estimated exploiting the empirical univariate proportions of wins. In the second stage the elements of Σ_{Z^*} , which are tetrachoric correlations, are estimated employing the bivariate proportions of wins. Finally, in the third stage the model parameters ψ are estimated by minimizing the function

$$(3.4) \quad G = \{\tilde{\kappa} - \kappa(\psi)\}' \hat{\mathbf{W}} \{\tilde{\kappa} - \kappa(\psi)\},$$

where $\tilde{\kappa}$ denotes the thresholds, and tetrachoric correlations, estimated in the first and second stages, $\kappa(\psi)$ denotes the thresholds, and tetrachoric correlations under the restrictions imposed on those parameters by the model parameters ψ and $\hat{\mathbf{W}}$ is a nonnegative definite matrix. Let Ξ denote the asymptotic covariance matrix of $\tilde{\kappa}$. Then it is possible to use $\hat{\mathbf{W}} = \hat{\Xi}^{-1}$ (Muthén, 1978), $\hat{\mathbf{W}} = [\text{diag}(\hat{\Xi})]^{-1}$ (Muthén, Du Toit and Spisic, 1997) or $\hat{\mathbf{W}} = \mathbf{I}$ (Muthén, 1993). The last two options seem more stable in data sets with a small number of objects (Maydeu-Olivares, 2001). This method is very fast, and Maydeu-Olivares (2001) states that it may have an edge over full information methods because it uses only the one and two-dimensional marginals of a large and sparse contingency table.

Pairwise likelihood (Le Cessie and Van Houwelingen, 1994) is a special case of the broader class of composite likelihoods (Lindsay, 1988; Varin, Reid and Firth, 2011). The pairwise likelihood of all the observations is the product of the pairwise likelihoods relative to the single judges $\mathcal{L}_{\text{pair}}(\psi; \mathbf{Y}) = \prod_{s=1}^S \mathcal{L}_{\text{pair}}^s(\psi; \mathbf{Y}_s)$, where

$$\begin{aligned} \mathcal{L}_{\text{pair}}^s(\psi; \mathbf{Y}_s) &= \prod_{i=1}^{n-2} \prod_{j=i+1}^{n-1} \prod_{k=i}^{n-1} \prod_{l=j+1}^n \text{pr}(Y_{sij} = y_{sij}, Y_{skl} = y_{skl}). \end{aligned}$$

Let $\ell_{\text{pair}}^s(\psi; \mathbf{Y}_s) = \log \mathcal{L}_{\text{pair}}^s(\psi; \mathbf{Y}_s)$ denote the logarithm of the pairwise likelihood for subject s and $\ell_{\text{pair}}(\psi; \mathbf{Y}) = \sum_{s=1}^S \ell_{\text{pair}}^s(\psi; \mathbf{Y}_s)$ be the whole pairwise log-likelihood. Under usual regularity conditions on the log-likelihood of univariate and bivariate margins, the maximum pairwise likelihood estimator is consistent and asymptotically normally distributed with mean ψ and covariance matrix $\mathbf{H}(\psi)^{-1} \mathbf{J}(\psi) \mathbf{H}(\psi)^{-1}$, where $\mathbf{J}(\psi) = \text{var}\{\nabla \ell_{\text{pair}}(\psi; \mathbf{Y})\}$ and $\mathbf{H}(\psi) = E\{-\nabla^2 \ell_{\text{pair}}(\psi; \mathbf{Y})\}$ (Molenberghs and Verbeke, 2005; Varin, Reid and Firth, 2011). Unfortunately, the analogous of the likelihood ratio test based on pairwise likelihood does not follow the usual chi-square distribution (Kent, 1982). In the multiple judgment sampling context, it is natural to consider asymptotic properties of pairwise likelihood estimators computed as the number of subjects increases, that is, as $S \rightarrow \infty$. When the number of paired comparisons per subject is bounded, the above properties are satisfied (Zhao and Joe, 2005). Pairwise likelihood reduces noticeably the computational effort since it requires only the computation of bivariate normal probabilities. The standard errors can be computed straightforwardly by exploiting the independence between the observations

TABLE 4

Average (Mn) and median (Md) simulated estimates, average model-based standard errors (s.e.) and simulation standard deviations (s.d.) of parameters estimated by maximum likelihood (ML), limited information estimation (LI) and pairwise likelihood (PL)

	True value	ML			LI				PL			
		Mn	s.e.	s.d.	Mn	Md	s.e.	s.d.	Mn	Md	s.e.	s.d.
μ_1	0.5	0.51	0.13	0.13	0.51	0.50	0.13	0.13	0.50	0.50	0.13	0.13
μ_2	0	0.01	0.12	0.13	0.01	0.01	0.12	0.13	0.01	0.01	0.12	0.13
μ_3	-0.5	-0.49	0.15	0.15	-0.50	-0.48	0.15	0.15	-0.49	-0.48	0.15	0.15
σ_{12}	0.8	0.80	0.12	0.14	0.78	0.80	0.13	0.14	0.79	0.80	0.13	0.15
σ_{13}	0.7	0.70	0.17	0.17	0.69	0.71	0.17	0.17	0.69	0.71	0.18	0.18
σ_{14}	0.8	0.79	0.13	0.14	0.78	0.79	0.13	0.14	0.78	0.80	0.14	0.15
σ_{23}	0.6	0.58	0.19	0.20	0.57	0.60	0.19	0.20	0.57	0.60	0.19	0.20
σ_{24}	0.7	0.68	0.16	0.16	0.66	0.67	0.16	0.17	0.67	0.68	0.16	0.17
σ_{34}	0.6	0.58	0.21	0.20	0.57	0.60	0.20	0.20	0.57	0.60	0.20	0.20

of different judges. In fact, $\mathbf{H}(\psi)$ can be estimated by the Hessian matrix computed at the maximum pairwise likelihood estimate, while the cross-product $\sum_{s=1}^S \nabla \ell_{\text{pair}}^s(\hat{\psi}; \mathbf{Y}_s) \nabla \ell_{\text{pair}}^s(\hat{\psi}; \mathbf{Y}_s)'$ can be used to estimate $\mathbf{J}(\psi)$.

The case of object-related dependencies is not considered in the following simulation study; however, note that some different difficulties arise. As already pointed out, in this context there is a large n and small S , so the limited information estimation method cannot be applied, but pairwise likelihood can still be employed (Cattelan, 2009). However, it is more problematic to consider the asymptotic behavior of the maximum pairwise likelihood estimator when data are a long sequence of dependent observations; see, for example, Cox and Reid (2004). In the context of paired comparison data, results of simulations for increasing n when all possible paired comparisons are performed are encouraging (Cattelan, 2009); however, theoretical results for this instance are still lacking.

3.4.2 *Simulation studies.* Simulation studies were performed considering models (3.1) and (3.2). It is assumed that $n = 4$; hence also a full likelihood approach based on the algorithm by Miwa, Hayter and Kuriki (2003) can be used since the integral has dimension 6.

The first simulation setting is the same as that proposed in Maydeu-Olivares (2001), where the model $\mathbf{Z}_s = \mathbf{A}\mathbf{T} + \mathbf{e}_s$ is assumed with

$$\boldsymbol{\mu} = \begin{pmatrix} 0.5 \\ 0 \\ -0.5 \\ 0 \end{pmatrix}, \quad \boldsymbol{\Sigma}_T = \begin{pmatrix} 1 & & & \\ 0.8 & 1 & & \\ 0.7 & 0.6 & 1 & \\ 0.8 & 0.7 & 0.6 & 1 \end{pmatrix}$$

and the covariance matrix of \mathbf{e} is $\boldsymbol{\Omega} = \omega^2 \mathbf{I}_6$. For identification purposes the diagonal elements of $\boldsymbol{\Sigma}_T$ are set

equal to 1, $\mu_4 = 0$ and $\omega^2 = 1$. Hence, in this case $\boldsymbol{\Sigma}_T$ is actually a correlation matrix. Table 4 shows the mean and medians of the simulated estimates on 1000 data sets assuming $S = 100$ judges. Moreover, the average of model-based standard errors and the simulation standard deviations are reported. In limited information estimation, the matrix $\hat{\mathbf{W}} = \mathbf{I}$ is employed. In this setting all the methods seem to perform comparably well. Table 5 shows the empirical coverages of confidence intervals based on the normal approximation.

The second simulation setting considers model (3.2) proposed by Tsai and Böckenholt (2008). Here, we consider differences with a reference object, so we compute means and variances of the differences $\tilde{T}_i = T_i - T_n$ for $i = 1, \dots, n - 1$. The assumed worth parameters of these differences are $\tilde{\boldsymbol{\mu}} = (-0.2, 1, -1.5)$

TABLE 5

Empirical coverage of confidence intervals for model parameters of limited information estimator (LI) and pairwise likelihood estimator (PL) at nominal levels 95%, 97.5% and 99%

	0.950		0.975		0.990	
	LI	PL	LI	PL	LI	PL
μ_1	0.947	0.958	0.982	0.978	0.992	0.992
μ_2	0.960	0.964	0.978	0.976	0.988	0.988
μ_3	0.941	0.930	0.969	0.972	0.995	0.991
σ_{12}	0.959	0.985	0.975	0.997	0.989	1.000
σ_{13}	0.934	0.939	0.961	0.967	0.968	0.985
σ_{14}	0.941	0.968	0.967	0.996	0.988	1.000
σ_{23}	0.965	0.970	0.973	0.980	0.987	0.995
σ_{24}	0.943	0.933	0.951	0.959	0.967	0.973
σ_{34}	0.953	0.946	0.969	0.966	0.977	0.989

TABLE 6

Average (Mn) and median (Md) simulated estimates, average model-based standard errors (s.e.) and simulation standard deviations (s.d.) of parameters estimated by maximum likelihood (ML), limited information estimation (LI) and pairwise likelihood (PL)

	True value	ML			LI				PL			
		Mn	s.e.	s.d.	Mn	Md	s.e.	s.d.	Mn	Md	s.e.	s.d.
$\tilde{\mu}_1$	-0.2	-0.21	0.19	0.18	-0.23	-0.21	0.21	0.22	-0.22	-0.20	0.19	0.19
$\tilde{\mu}_2$	1	1.00	0.30	0.31	1.07	1.07	0.42	0.47	1.03	1.00	0.33	0.33
$\tilde{\mu}_3$	-1.5	-1.51	0.31	0.32	-1.59	-1.59	0.49	0.51	-1.54	-1.51	0.36	0.35
$\tilde{\sigma}_1^2$	1.5	1.53	0.83	0.81	2.06	1.58	1.97	1.64	1.70	1.44	1.05	0.95
$\tilde{\sigma}_2^2$	4	3.98	1.73	1.75	5.34	4.37	4.42	4.48	4.45	3.92	2.42	2.15
$\tilde{\sigma}_3^2$	3	3.01	1.41	1.42	3.91	3.19	3.17	3.25	3.32	3.04	1.93	1.73
$\tilde{\sigma}_{12}$	1	0.98	0.70	0.64	1.34	1.06	1.44	1.30	1.12	0.97	0.87	0.77
$\tilde{\sigma}_{13}$	1.3	1.29	0.73	0.71	1.72	1.39	1.48	1.49	1.43	1.27	0.95	0.84
$\tilde{\sigma}_{23}$	2.5	2.49	1.09	1.09	3.35	2.72	2.67	2.77	2.77	2.48	1.53	1.33
b	0.5	0.53	0.41	0.39	0.72	0.58	0.82	0.98	0.58	0.50	0.50	0.51

while the covariance matrix is

$$\begin{pmatrix} 1.5 & 1 & 1.3 \\ 1 & 4 & 2.5 \\ 1.3 & 2.5 & 3 \end{pmatrix},$$

and $\tilde{\sigma}_{ij}$ is used to denote the element in row i and column j of the above reduced matrix. Differently from the previous setting, this specification of the model allows one to estimate also the variance of the differences $T_i - T_n$ and to check whether they are different for the various objects. Tsai and Böckenholt (2008) propose a specification of the matrix \mathbf{B} which depends only on one parameter b whose value is set equal to 0.5.

Table 6 presents the results of the simulations. Maximum likelihood based on numerical integration is the method that performs best; however, maximization of the likelihood was not always straightforward, and sometimes the optimization algorithms employed stopped at a point where the Hessian matrix was not negative definite.

Pairwise likelihood estimation seems to perform quite well, especially if compared to limited information estimation, which seems not satisfactory in this case with $S = 100$, as already noticed in Tsai and Böckenholt (2008). Estimating the parameters of the covariance matrix appears more problematic than the estimation of the worth parameters, and the average of the simulated estimates is particularly influenced by some large values, but the median shows a better performance. In particular, while the average simulated estimates for limited information estimation shows a maximum percentage bias equal to 44.1%, for the median it reduces to 15.4%. The maximum bias for the

mean of the simulated estimates using pairwise likelihood is 16.1%, while for the median it is 4%. In both cases, pairwise likelihood shows lower bias. The standard errors of pairwise likelihood estimates are lower, thus yielding shorter confidence intervals. Table 7 reports the empirical coverage of Wald-type confidence intervals for the estimated limited information estimation and pairwise likelihood. The coverage rates of the two methods are very similar, and in both cases the actual coverage for parameters of the covariance matrix appears systematically lower than the nominal levels. In order to obtain accurate coverage probabilities, we may need to resort to a bootstrap procedure for detecting the distribution of the statistic, while with pairwise

TABLE 7

Empirical coverage of confidence intervals for model parameters of limited information estimator (LI) and pairwise likelihood estimator (PL) at nominal levels 95%, 97.5% and 99%

	0.950		0.975		0.990	
	LI	PL	LI	PL	LI	PL
$\tilde{\mu}_1$	0.955	0.935	0.981	0.965	0.994	0.983
$\tilde{\mu}_2$	0.962	0.960	0.973	0.974	0.986	0.986
$\tilde{\mu}_3$	0.920	0.938	0.941	0.960	0.961	0.977
$\tilde{\sigma}_1^2$	0.932	0.922	0.947	0.936	0.959	0.966
$\tilde{\sigma}_2^2$	0.932	0.924	0.949	0.945	0.964	0.961
$\tilde{\sigma}_3^2$	0.936	0.937	0.949	0.953	0.963	0.964
$\tilde{\sigma}_{12}$	0.932	0.937	0.951	0.956	0.966	0.970
$\tilde{\sigma}_{13}$	0.915	0.912	0.929	0.933	0.939	0.945
$\tilde{\sigma}_{23}$	0.922	0.920	0.941	0.937	0.953	0.951
b	0.936	0.936	0.946	0.953	0.963	0.963

TABLE 8

Estimates and standard errors (in brackets) of mean and correlation parameters of model (3.1) for universities data using constraints proposed by Maydeu-Olivares and Hernández (2007). In italics the estimates and standard errors of a model with fixed correlation between Paris and St. Gallen

	Barcelona	London	Milan	Paris	St. Gallen	Stockholm	μ
Barcelona	1 (fixed)	<i>-0.064</i> (0.183)	<i>0.688</i> (0.085)	<i>0.063</i> (0.158)	<i>-0.472</i> (0.146)	<i>0.265</i> (0.145)	0.405 (0.073)
London	0.058 (0.084)	1 (fixed)	<i>0.079</i> (0.185)	<i>-0.069</i> (0.224)	<i>-0.287</i> (0.147)	<i>0.227</i> (0.154)	1.346 (0.087)
Milan	0.724 (0.062)	0.185 (0.097)	1 (fixed)	<i>0.244</i> (0.174)	<i>-0.466</i> (0.137)	<i>0.253</i> (0.160)	0.308 (0.074)
Paris	0.171 (0.094)	0.054 (0.117)	0.331 (0.113)	1 (fixed)	<i>-0.690</i> (fixed)	<i>0.033</i> (0.267)	0.748 (0.086)
St. Gallen	<i>-0.303</i> (0.113)	<i>-0.139</i> (0.139)	<i>-0.298</i> (0.144)	<i>-0.496</i> (0.157)	1 (fixed)	<i>0.194</i> (0.135)	0.371 (0.081)
Stockholm	0.350 (0.079)	0.316 (0.091)	0.339 (0.097)	0.144 (0.113)	0.287 (0.130)	1 (fixed)	0 (fixed)

likelihood it may be possible to obtain intervals based on the pairwise likelihood function.

EXAMPLE. We fit model (3.1) to universities' data by means of pairwise likelihood. A full likelihood approach based on numeric approximation implies computing 303 integrals of dimension 5, in case a university is used as reference object, both for the mean and covariance structure, but methods such as the Gauss-Hermite quadrature are affected by the curse of dimensionality. A multivariate probit approach would require a very slow computation because the algorithm by Miwa would take very long to approximate 303 integrals of dimension 15. It is assumed that $\Omega = \omega^2 \mathbf{I}_{15}$. Table 8 displays the results of the estimates, employing two different sets of constraints. The lower triangle of the covariance matrix shown in Table 8 reports the estimates obtained using the constraints proposed in Maydeu-Olivares and Hernández (2007); see Section 3.2.2. The estimate of the threshold parameter (with standard error in brackets) is $\hat{\tau}_2 = 0.205$ (0.018) while the variance parameter is $\hat{\omega}^2 = 0.180$ (0.026). A high correlation is estimated between Barcelona and Milan, so strong preference for Barcelona is associated with strong preference for Milan. Even though some correlations do not seem significant, it appears that a strong preference for St. Gallen is associated with a weak preference for all the other universities but Stockholm. The worth parameters denote the same ranking of all universities as the one arising from Table 2. However, note that the estimated worth parameters cannot be considered as absolute measures of worth of

items; indeed, it is possible to obtain alternative solutions that give an equivalent fitting. The mean parameters that can be identified in the model are standardized differences, that is, $(\mu_i - \mu_6) / \sqrt{\sigma_i^2 + \sigma_6^2 - 2\sigma_{i6} + \omega^2}$, $i = 1, \dots, 5$, where μ_6 and σ_6^2 are the mean and variance of the latent variable referring to Stockholm, the reference university. From the identified parameters, different covariance matrixes of the universities can be recovered. For example, in this instance where the matrix Σ_T can be interpreted as a correlation matrix, it is shown that the worth parameters $\sqrt{c}\mu$, the correlation matrix $c\Sigma_T + (1 - c)\mathbf{1}\mathbf{1}'$ and the covariance matrix of the pair-specific errors $c\Omega$ produce the same fitting of the model for a positive constant c such that the correlation matrix remains positive definite (Maydeu-Olivares and Hernández, 2007). It is possible to set one of the parameters of the correlation matrix according to some assumption, for example we may presume that a strong preference for Paris is associated with a weak preference for St. Gallen, and determine the value of c which minimizes the correlation between the two universities while yielding a positive definite correlation matrix. The value is $c = 1.13$ which produces a correlation between Paris and St. Gallen equal to -0.690 . The estimates of the correlation matrix with this fixed value of correlation between Paris and St. Gallen are shown in the upper triangle of the matrix in Table 8. The worth parameters can be computed by multiplying the estimates shown in Table 8 by $\sqrt{1.13}$. The fitting of the two models is equal, but in the second case estimation is based on some previous theory about correlation between a certain couple of universities.

This analysis has only an illustrative purpose, in particular Böckenholt (2001b) finds that a model with thresholds that vary among subjects performs better than a model with a constant threshold parameter.

3.4.3 *Model selection and goodness of fit.* Paired comparison data can be arranged in a contingency table. In case of multiple judgment sampling the data can be arranged in a table of dimension 2^N when there are two possible outcomes and H^N when the outcomes are H -categorical. As a result, the contingency table will typically be very sparse, especially if covariates are included so that paired comparisons are observed conditional on the values of the covariates. In this situation the likelihood ratio statistic and the Pearson statistic do not follow a χ^2 distribution, nevertheless these statistics are often employed to assess the model and for model selection. Differences between observed and expected frequencies for subsets of the data, as the 2×2 subtables or triplets of comparisons, are sometimes considered in order to identify where the fitting of the model is not good. In Dittrich et al. (2007) the deviance is used for selection between nested models, but the test of goodness of fit cannot be based on the asymptotic χ^2 distribution so a Monte Carlo procedure is employed.

Since the goodness of fit of the model cannot be assessed through the usual statistics and Monte Carlo procedures are computationally expensive, some statistics based on lower dimensional marginals of the contingency table have been proposed. In general the statistics proposed are quadratic forms of the residuals

$$(3.5) \quad \{\mathbf{p}_r - \boldsymbol{\pi}_r(\hat{\boldsymbol{\psi}})\}' \mathbf{C} \{\mathbf{p}_r - \boldsymbol{\pi}_r(\hat{\boldsymbol{\psi}})\},$$

where \mathbf{C} is a weight matrix, \mathbf{p}_r denotes the sample marginal proportions and r denotes a set of lower order marginals.

Maydeu-Olivares (2001) considers the statistic G as in (3.4) employed for estimation, which corresponds to setting $\mathbf{C} = \mathbf{W}$ in (3.5) and r denoting univariate and bivariate marginal probabilities. The statistic $S\hat{G}$ is analyzed in order to test $H_0: \boldsymbol{\kappa} = \boldsymbol{\kappa}(\boldsymbol{\psi})$. When $\hat{\mathbf{W}} = \hat{\boldsymbol{\Sigma}}^{-1}$, then $S\hat{G} \xrightarrow{d} \chi_d^2$ where $d = N(N+1)/2 - q$ and q is the number of model parameters. However, when $\hat{\mathbf{W}} = [\text{diag}(\hat{\boldsymbol{\Sigma}})]^{-1}$ or $\hat{\mathbf{W}} = I$, the asymptotic distribution of the statistic is a weighted sum of d chi-square random variables with one degree of freedom. Maydeu-Olivares (2001) proposes to rescale the test statistic in order to match the asymptotic chi-square distribution. The same procedure is followed in the proposal for testing $H_0: \boldsymbol{\pi}_2 = \boldsymbol{\pi}_2(\boldsymbol{\psi})$, where $\boldsymbol{\pi}_2$ is the

vector of all univariate and bivariate marginal probabilities. Maydeu-Olivares (2006) considers the testing of further hypotheses but the issue of the asymptotic distribution being a weighted sum of chi-square distributions remains.

Maydeu-Olivares and Joe (2005) consider testing the hypothesis $H_0: \tilde{\boldsymbol{\pi}} = \tilde{\boldsymbol{\pi}}(\boldsymbol{\psi})$ in a multidimensional contingency table, where $\tilde{\boldsymbol{\pi}}$ is the 2^N -dimensional vector of joint probabilities. Again, the use of marginal residuals up to order r is considered. Let $\boldsymbol{\pi}$ denote a vector which stacks all the marginal probabilities: univariate, bivariate, trivariate and so on. There is a one-to-one correspondence between $\tilde{\boldsymbol{\pi}}$ and $\boldsymbol{\pi}$ so that for a particular matrix $\boldsymbol{\Lambda}$ of 0's and 1's $\boldsymbol{\pi} = \boldsymbol{\Lambda}\tilde{\boldsymbol{\pi}}$. If only marginal probabilities up to order r are considered, then $\boldsymbol{\pi}_r = \boldsymbol{\Lambda}_r\tilde{\boldsymbol{\pi}}$ for a sub-matrix $\boldsymbol{\Lambda}_r$ of $\boldsymbol{\Lambda}$. Let $\boldsymbol{\Delta} = \partial\tilde{\boldsymbol{\pi}}/\partial\boldsymbol{\psi}$ and $\boldsymbol{\Gamma} = \mathbf{E} - \tilde{\boldsymbol{\pi}}\tilde{\boldsymbol{\pi}}'$, where $\mathbf{E} = \text{diag}(\tilde{\boldsymbol{\pi}})$. Maydeu-Olivares and Joe (2005) propose the statistic

$$(3.6) \quad M_r = S\{\mathbf{p}_r - \boldsymbol{\pi}_r(\hat{\boldsymbol{\psi}})\}' \mathbf{C}_r(\hat{\boldsymbol{\psi}}) \{\mathbf{p}_r - \boldsymbol{\pi}_r(\hat{\boldsymbol{\psi}})\},$$

where $\mathbf{C}_r(\boldsymbol{\psi}) = \mathbf{F}_r^{-1} - \mathbf{F}_r^{-1}\boldsymbol{\Delta}_r(\boldsymbol{\Delta}_r'\mathbf{F}_r^{-1}\boldsymbol{\Delta}_r)^{-1}\boldsymbol{\Delta}_r'\mathbf{F}_r^{-1}$, $\mathbf{F}_r = \boldsymbol{\Lambda}_r\boldsymbol{\Gamma}\boldsymbol{\Lambda}_r'$ and $\boldsymbol{\Delta}_r = \boldsymbol{\Lambda}_r\boldsymbol{\Delta}$. M_r is asymptotically distributed as a χ_{l-q}^2 random variable where l is the length of \mathbf{p}_r . The M_r statistic asymptotically follows a chi-square distribution not only when $\hat{\boldsymbol{\psi}}$ is the maximum likelihood estimator, but also when it is a \sqrt{S} -consistent estimate, such as the limited information estimator and the pairwise likelihood estimator presented in Section 3.4.1. Since the marginals should not be sparse, Maydeu-Olivares and Joe (2005) suggest to use M_2 when the model is identified using only univariate and bivariate information, also because only up to bivariate sample moments and four-way model probabilities are involved in the computation of M_2 . As the number of cells gets larger, the dimension of the matrices involved in (3.6) increases noticeably, and tricks may be necessary to do the computations. Analysis and extensions of this type of test are considered in Maydeu-Olivares and Joe (2006), Reiser (2008) and Joe and Maydeu-Olivares (2010). All applications considered regard item response theory, so an investigation of their performance in paired comparison data is necessary to understand the sample size needed for obtaining accurate Type I errors using M_2 .

4. SOFTWARE

Fitting models to paired comparison data is facilitated by some R packages which allow fitting of the classical models and, in some cases, also fitting of more complicated models.

The `eba` package (Wickelmaier and Schmid, 2004) fits elimination by aspects models (Tversky, 1972) to paired comparison data. The elimination by aspects model assumes that different objects present various aspects. The worth of each object is the sum of the worth associated with each aspect possessed by the object. When all objects possess only one relevant aspect, then the elimination by aspects model reduces to the Bradley–Terry model. Therefore, in case only one aspect per object is specified, the function `eba` can be used to fit model (2.1) with logit link, while when the link is probit the function `thurstone` can be used. The function `strans` checks how many violations of weak, moderate and strong stochastic transitivity are present in the data.

The `prefmod` package (Hatzinger, 2010) fits Bradley–Terry models exploiting their log-linear representation. Ordinal paired comparisons are allowed, but the software reduces the total number of categories to three or two, depending on whether there is a no preference category or not.

There are three different functions for estimating models for paired comparison data: the `llbt.fit` function which estimates the log-linear version of the Bradley–Terry model through the estimation algorithm described in Hatzinger and Francis (2004), the `llbtPC.fit` function that estimates the log-linear model exploiting the `gnm` (Turner and Firth, 2010b) function for fitting generalized nonlinear models and the `pattPC.fit` function, which fits paired comparison data using a pattern design, that is, all possible patterns of paired comparisons. The latter function handles also some cases in which the responses are missing not at random; see Section 5. A difficulty of this approach is that the response table grows dramatically with the number of objects since, in case of only two possible outcomes, the number of patterns is 2^N , so no more than six objects can be included with two response categories, and not more than five with three response categories. Finally, the function `pattnpml.fit` fits a mixture model to overdispersed paired comparison data using nonparametric maximum likelihood.

The `BradleyTerry2` package (Turner and Firth, 2010a) expands the previous `BradleyTerry` (Firth, 2008) package and allows one to fit the unstructured model (2.1) and extension (2.3) with logit, probit and `cauchit` link functions, including also comparison-specific covariates. Model fitting is either by maximum likelihood, penalized quasi-likelihood or bias-reduced maximum likelihood (Firth, 1993). In case of object

specific random effects, as in model (3.3), penalized quasi-likelihood (Breslow and Clayton, 1993) is used, while when an object wins or loses all the paired comparisons in which it is involved and its estimate worth parameter is infinite, then the bias-reduced maximum likelihood produces finite estimates. If there are missing explanatory variables, an additional worth parameter for the object with missing covariates is estimated. Order effects and more general comparison-specific covariates can be included, but only win-loss responses are allowed.

The package `psychotree` (Strobl, Wickelmaier and Zeileis, 2011) implements the method for recursive partitioning of the subjects on the basis of their explanatory variables and estimates an unstructured Bradley–Terry model for each of the final subgroups of subjects; see Section 2.4.

Although the available packages have many useful features, a combination of those provided by the different packages and also some additional features could be of practical help. The `prefmod` and `BradleyTerry2` packages were built with the aim of analyzing multiple judgment data and tournament-like data, respectively. This is reflected in the different characteristics of the packages. A function that can handle data with at least three-categorical results, thus allowing for the “no preference” category, include different link functions, and an easy implementation of object-, subject- and comparison-specific covariates in a linear model framework would be useful. The available methods for including dependencies between observations are only in a log-linear framework through the introduction of further parameters in the predictor or including object-related random effects, which are estimated by means of penalized quasi likelihood, a method that does not perform well with binary data. At present, there are no available packages for the analysis of paired comparison data that allow the fitting of models as those presented in Section 3.2.1. However, implementation of pairwise likelihood estimation for those models is straightforward since it implies only the computation of bivariate normal probabilities.

5. CONCLUSIONS

This paper reviews some of the extensions proposed in the literature to the two most commonly applied models for paired comparison data, namely the Bradley–Terry and the Thurstone models. However, not every aspect could be considered here, and among

issues that have not been treated, there are the development of models for multi-dimensional data when objects are evaluated with respect to multiple aspects (Böckenholt, 1988; Dittrich et al., 2006), the temporal extension for comparisons repeated in time (Fahrmeir and Tutz, 1994, Glickman, 2001, Böckenholt, 2002, Dittrich, Francis and Katzenbeisser, 2008), the estimation of abilities of individuals belonging to a team that performs the paired comparisons (Huang, Weng and Lin, 2006; Menke and Martinez, 2008) and many more. Another important issue concerns the optimal design of the experiment. Graßhoff et al. (2004) show that the minimum sample size required for maximizing the determinant of the information matrix in an unstructured Bradley–Terry model requires that every comparison is performed once. When objects are specified using factors with a certain number of levels, the required sample size grows exponentially, while the number of parameters grows linearly as the number of factors increases. Some designs, in order to reduce the number of required comparisons, are investigated in Graßhoff et al. (2004). In Graßhoff and Schwabe (2008) a characterization of the locally optimal design in case of two factors design in a Bradley–Terry model is given, but for more complex situations it seems difficult to give general results. Goos and Grossmann (2011) consider also the problem when within-pair order effects are present. It seems that investigation of these issues in other models are not present in the literature.

The methods for independent data are well established, and a lot of literature has been published about them. The problem of the asymptotic behavior of the maximum likelihood estimator has been tackled. The case of a fixed number of objects and increasing number of comparisons per couple does not seem to pose particular difficulties for standard arguments, while more problematic appears the instance of a fixed number of comparisons per couple and increasing number of items. In the context of the unstructured Bradley–Terry model, Simons and Yao (1999) find a condition on the growth rate of the largest ratio between item worth parameters which assures that the maximum likelihood estimator is consistent and asymptotically normally distributed. Yan, Yang and Xu (2012) investigate the case in which the number of comparisons per couple is not fixed, and some comparisons may also be missing, and find a condition that assures normality of the maximum likelihood estimator. We are not acquainted with any other investigation of asymptotic behavior of estimators in models different from the unstructured, independent Bradley–Terry model.

Particular attention has been focused on models for dependent data. Thurstonian models appear particularly suitable to account for dependence between observations. However, the problems posed by the identification restrictions are noticeable. The estimated model has to be interpreted with reference to a class of covariance matrices, and different identification restrictions may lead to different class of matrices. It is possible to rotate the matrix according to a predefined hypothesis about the covariance between certain items (Maydeu-Olivares and Hernández, 2007), but the estimated standard errors vary depending on the fixed parameters and the significance of the other estimated parameters changes.

In the multiple judgment sampling scheme it is often stated that if a judge does not perform all paired comparisons, then it suffices to define subject-specific matrices \mathbf{A}_s (see Section 3.2.1) with rows corresponding only to the comparisons performed by judge s . However, it is expected that this may be problematic for estimation by means of limited information estimation, and there are no studies about the consequences of missing data in this estimation method.

Missing observations cause problems also for testing the goodness of fit since quadratic statistics as (3.5) assume that all comparisons are performed by all subjects.

Missing data may derive from the design of the experiment, for example when n is very large, and only a subset of all comparisons is presented to each subject. Otherwise, if many comparisons are performed by the same subject it may be necessary to account for the fatigue of subjects and/or for the passing of time when comparisons take long in order to be accomplished.

Dittrich et al. (2012) consider the problem of missing data in the context of the log-linear representation of the Bradley–Terry model since the study of the missing mechanism may shed light on the psychological process. It is assumed that the probability that a comparison is missing follows a logistic distribution since this facilitates the fitting of the model. However, the likelihood for such models is not easy to compute, and the function in the `prefmod` package allows one to compute it only for data with up to six objects. It is not easy to discriminate between different types of missing mechanisms, and a very large number of observations may be needed in order to discriminate between a missing completely at random and missing not at random situation.

The economic theory points out some problems in choice data that have not been considered yet. The

main aspects which may need to be incorporated in models include the influence that subjects can have on each other, the influence of one particular subject, that may be some sort of leader, over all the other judges and the dependence on choices caused by the social and cultural context. Inclusions of these aspects will inevitably lead to even more complicated models for paired comparison data.

Finally, methods for object-related dependencies present many open problems. Most of the issues are connected to the dependence among all comparisons which is typically present in this context. Moreover, the scheme of paired comparisons is often much less balanced than in psychometric experiments. Asymptotic theory in models for dependent data when the number of items compared increases has not been developed yet. Maximum pairwise likelihood estimation provided encouraging results, but more extensive studies seem necessary. In this case, computation of standard errors is problematic since there are no independent replications of the data, so a viable alternative lies in parametric bootstrap. Methods for model selection and goodness of fit described in Section 3.4.3 require independent replication of all comparisons; hence they cannot be employed in this setting.

ACKNOWLEDGMENTS

The author would like to thank Cristiano Varin and Alessandra Salvan for helpful comments and the anonymous referees and an Associate Editor for comments and suggestions that led to a substantial improvement of the manuscript.

REFERENCES

- AGRESTI, A. (1992). Analysis of ordinal paired comparison data. *J. R. Stat. Soc. Ser. C Appl. Stat.* **41** 287–297.
- AGRESTI, A. (2002). *Categorical Data Analysis*, 2nd ed. Wiley, New York. MR1914507
- BARRY, D. and HARTIGAN, J. A. (1993). Choice models for predicting divisional winners in major league baseball. *J. Amer. Statist. Assoc.* **88** 766–774.
- BÄUML, K. H. (1994). Upright versus upside-down faces: How interface attractiveness varies with orientation. *Percept. Psychophys.* **56** 163–172.
- BÖCKENHOLT, U. (1988). A logistic representation of multivariate paired-comparison models. *J. Math. Psych.* **32** 44–63. MR0935673
- BÖCKENHOLT, U. (2001a). Hierarchical modeling of paired comparison data. *Psychol. 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. (2002). A Thurstonian analysis of preference change. *J. Math. Psych.* **46** 300–314. MR1920807
- BÖCKENHOLT, U. (2004). Comparative judgments as an alternative to ratings: Identifying the scale origin. *Psychol. Methods* **9** 453–465.
- BÖCKENHOLT, U. (2006). Thurstonian-based analyses: Past, present, and future utilities. *Psychometrika* **71** 615–629. MR2312235
- BÖCKENHOLT, U. and DILLON, W. R. (1997a). Modeling within-subject dependencies in ordinal paired comparison data. *Psychometrika* **62** 411–434.
- BÖCKENHOLT, U. and DILLON, W. R. (1997b). Some new methods for an old problem: Modeling preference changes and competitive market structures in pretest market data. *Journal of Marketing Research* **34** 130–142.
- BÖCKENHOLT, U. and TSAI, R. C. (2001). Individual differences in paired comparison data. *Br. J. Math. Stat. Psychol.* **54** 265–277.
- BÖCKENHOLT, U. and TSAI, R. C. (2007). Random-effects models for preference data. In *Handbook of Statistics* (C. R. Rao and S. Sinharay, eds.) **26** 447–468. Elsevier, Amsterdam.
- BRADLEY, R. A. (1976). Science, statistics, and paired comparisons. *Biometrics* **32** 213–232. MR0408132
- BRADLEY, R. A. and TERRY, M. E. (1952). Rank analysis of incomplete block designs. I. The method of paired comparisons. *Biometrika* **39** 324–345. MR0070925
- BRESLOW, N. E. and CLAYTON, D. G. (1993). Approximate inference in generalized linear mixed models. *J. Amer. Statist. Assoc.* **88** 9–25.
- CARROLL, J. D. and DE SOETE, G. (1991). Toward a new paradigm for the study of multiattribute choice behavior. Spatial and discrete modeling of pairwise preferences. *American Psychologist* **46** 342–351.
- CATTELAN, M. (2009). Correlation models for paired comparison data. Ph.D. thesis, Dept. Statistical Sciences, Univ. Padua.
- CATTELAN, M., VARIN, C. and FIRTH, D. (2012). Dynamic Bradley–Terry modelling of sports tournaments. *J. R. Stat. Soc. Ser. C Appl. Stat.* To appear.
- CAUSEUR, D. and HUSSON, F. (2005). A 2-dimensional extension of the Bradley–Terry model for paired comparisons. *J. Statist. Plann. Inference* **135** 245–259. MR2200468
- CHIB, S. and GREENBERG, E. (1998). Analysis of multivariate probit models. *Biometrika* **85** 347–361.
- CHOISEL, S. and WICKELMAIER, F. (2007). Evaluation of multi-channel reproduced sound: Scaling auditory attributes underlying listener preference. *J. Acoust. Soc. Am.* **121** 388–400.
- COX, D. R. and REID, N. (2004). A note on pseudolikelihood constructed from marginal densities. *Biometrika* **91** 729–737. MR2090633
- CRAIG, P. (2008). A new reconstruction of multivariate normal orthant probabilities. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **70** 227–243. MR2412640
- DAVID, H. A. (1988). *The Method of Paired Comparisons*, 2nd ed. Griffin's Statistical Monographs & Courses **41**. Griffin, London. MR0947340
- DAVIDSON, R. R. (1970). On extending the Bradley–Terry model to accommodate ties in paired comparison experiments. *J. Amer. Statist. Assoc.* **65** 317–328.
- DAVIDSON, R. R. and FARQUHAR, P. H. (1976). A bibliography on the method of paired comparisons. *Biometrics* **32** 241–252. MR0408134

- DE SOETE, G. and WINSBERG, S. (1993). A Thurstonian pairwise choice model with univariate and multivariate spline transformations. *Psychometrika* **58** 233–256.
- DILLON, W. R., KUMAR, A. and DE BORRERO, M. S. (1993). Capturing individual differences in paired comparisons: An extended BTL model incorporating descriptor variables. *Journal of Marketing Research* **30** 42–51.
- DITTRICH, R., FRANCIS, B. and KATZENBEISSER, W. (2008). Temporal dependence in longitudinal paired comparisons. Research report, Dept. Statistics and Mathematics, WU Vienna Univ. Economics and Business.
- DITTRICH, R., HATZINGER, R. and KATZENBEISSER, W. (1998). Modelling the effect of subject-specific covariates in paired comparison studies with an application to university rankings. *J. R. Stat. Soc. Ser. C Appl. Stat.* **47** 511–525.
- DITTRICH, R., HATZINGER, R. and KATZENBEISSER, W. (2001). Corrigendum: “Modelling the effect of subject-specific covariates in paired comparison studies with an application to university rankings.” *J. R. Stat. Soc. Ser. C Appl. Stat.* **50** 247–249. [MR1833276](#)
- DITTRICH, R., HATZINGER, R. and KATZENBEISSER, W. (2002). Modelling dependencies in paired comparison data: A log-linear approach. *Comput. Statist. Data Anal.* **40** 39–57. [MR1921121](#)
- DITTRICH, R., HATZINGER, R. and KATZENBEISSER, W. (2004). A log-linear approach for modelling ordinal paired comparison data on motives to start a PhD program. *Stat. Model.* **4** 1–13.
- DITTRICH, R., FRANCIS, B., HATZINGER, R. and KATZENBEISSER, W. (2006). Modelling dependency in multivariate paired comparisons: A log-linear approach. *Math. Social Sci.* **52** 197–209. [MR2257629](#)
- DITTRICH, R., FRANCIS, B., HATZINGER, R. and KATZENBEISSER, W. (2007). A paired comparison approach for the analysis of sets of Likert-scale responses. *Stat. Model.* **7** 3–28. [MR2749821](#)
- DITTRICH, R., FRANCIS, B., HATZINGER, R. and KATZENBEISSER, W. (2012). Missing observations in paired comparison data. *Stat. Model.* **12** 117–143.
- DUINEVELD, C. A. A., ARENTS, P. and KING, B. M. (2000). Log-linear modelling of paired comparison data from consumer tests. *Food Quality and Preference* **11** 63–70.
- ELLERMEIER, W., MADER, M. and DANIEL, P. (2004). Scaling the unpleasantness of sounds according to the BTL model: Ratio-scale representation and psychoacoustical analysis. *Acta Acustica United with Acustica* **90** 101–107.
- FAHRMEIR, L. and TUTZ, G. (1994). Dynamic stochastic models for time-dependent ordered paired comparison systems. *J. Amer. Statist. Assoc.* **89** 1438–1449.
- FIRTH, D. (1993). Bias reduction of maximum likelihood estimates. *Biometrika* **80** 27–38. [MR1225212](#)
- FIRTH, D. (2005). Bradley–Terry models in R. *Journal of Statistical Software* **12** 1–12.
- FIRTH, D. (2008). BradleyTerry: Bradley–Terry models. Available at <http://CRAN.R-project.org/package=BradleyTerry>.
- FIRTH, D. and DE MENEZES, R. X. (2004). Quasi-variances. *Biometrika* **91** 65–80. [MR2050460](#)
- FORD, L. R. JR. (1957). Solution of a ranking problem from binary comparisons. *Amer. Math. Monthly* **64** 28–33. [MR0097876](#)
- FRANCIS, B., DITTRICH, R. and HATZINGER, R. (2010). Modelling heterogeneity in ranked responses by nonparametric maximum likelihood: How do Europeans get their scientific knowledge? *Ann. Appl. Stat.* **4** 2181–2202. [MR2829952](#)
- FRANCIS, B., DITTRICH, R., HATZINGER, R. and PENN, R. (2002). Analysing partial ranks by using smoothed paired comparison methods: An investigation of value orientation in Europe. *J. R. Stat. Soc. Ser. C Appl. Stat.* **51** 319–336. [MR1920800](#)
- GENZ, A. and BRETZ, F. (2002). Comparison of methods for the computation of multivariate *t* probabilities. *J. Comput. Graph. Statist.* **11** 950–971. [MR1944269](#)
- GLENN, W. A. and DAVID, H. A. (1960). Ties in paired-comparison experiments using a modified Thurstone–Mosteller model. *Biometrics* **16** 86–109.
- GLICKMAN, M. E. (2001). Dynamic paired comparison models with stochastic variances. *J. Appl. Stat.* **28** 673–689. [MR1862491](#)
- GOOS, P. and GROSSMANN, H. (2011). Optimal design of factorial paired comparison experiments in the presence of within-pair order effects. *Food Quality and Preference* **22** 198–204.
- GRABHOFF, U. and SCHWABE, R. (2008). Optimal design for the Bradley–Terry paired comparison model. *Stat. Methods Appl.* **17** 275–289. [MR2425186](#)
- GRABHOFF, U., GROBMANN, H., HOLLING, H. and SCHWABE, R. (2004). Optimal designs for main effects in linear paired comparison models. *J. Statist. Plann. Inference* **126** 361–376. [MR2090864](#)
- HATZINGER, R. (2010). `prefmod`: Utilities to fit paired comparison models for preferences. Available at <http://CRAN.R-project.org/package=prefmod>.
- HATZINGER, R. and FRANCIS, B. J. (2004). Fitting paired comparison models in R. Research report, Univ. Wien. Available at <http://epub.wu.ac.at/id/eprint/740>.
- HEAD, M. L., DOUGHTY, P., BLOMBERG, S. P. and KEOGH, S. (2008). Chemical mediation of reciprocal mother–offspring recognition in the Southern Water Skink (*Eulamprus heatwolei*). *Australian Ecology* **33** 20–28.
- HENERY, R. J. (1992). An extension to the Thurstone–Mosteller model for chess. *The Statistician* **41** 559–567.
- HUANG, T.-K., WENG, R. C. and LIN, C.-J. (2006). Generalized Bradley–Terry models and multi-class probability estimates. *J. Mach. Learn. Res.* **7** 85–115. [MR2274363](#)
- JOE, H. and MAYDEU-OLIVARES, A. (2010). A general family of limited information goodness-of-fit statistics for multinomial data. *Psychometrika* **75** 393–419. [MR2719935](#)
- KENT, J. T. (1982). Robust properties of likelihood ratio tests. *Biometrika* **69** 19–27. [MR0655667](#)
- KISSLER, J. and BÄUML, K. H. (2000). Effects of the beholder’s age on the perception of facial attractiveness. *Acta Psychol. (Amst)* **104** 145–166.
- KNORR-HELD, L. (2000). Dynamic rating of sports teams. *The Statistician* **49** 261–276.
- LANCASTER, J. F. and QUADE, D. (1983). Random effects in paired-comparison experiments using the Bradley–Terry model. *Biometrics* **39** 245–249. [MR0712751](#)
- LE CESSIE, S. and VAN HOUWELINGEN, J. C. (1994). Logistic regression for correlated binary data. *J. R. Stat. Soc. Ser. C Appl. Stat.* **43** 95–108.
- LELE, S. R., NADEEM, K. and SCHMULAND, B. (2010). Estimability and likelihood inference for generalized linear mixed

- models using data cloning. *J. Amer. Statist. Assoc.* **105** 1617–1625. [MR2796576](#)
- LINDSAY, B. G. (1988). Composite likelihood methods. In *Statistical Inference from Stochastic Processes (Ithaca, NY, 1987)*. *Contemp. Math.* **80** 221–239. Amer. Math. Soc., Providence, RI. [MR0999014](#)
- LUCE, R. D. (1959). *Individual Choice Behavior: A Theoretical Analysis*. Wiley, New York. [MR0108411](#)
- MARSCHAK, J. (1960). Binary-choice constraints and random utility indicators. In *Mathematical Methods in the Social Sciences, 1959* (Arrow, K. J., Karlin, S. and Suppes, S., eds.) 312–329. Stanford Univ. Press, Stanford, CA. [MR0118556](#)
- MATTHEWS, J. N. S. and MORRIS, K. P. (1995). An application of Bradley–Terry-type models to the measurement of pain. *J. R. Stat. Soc. Ser. C Appl. Stat.* **44** 243–255.
- MAYDEU-OLIVARES, A. (2001). Limited information estimation and testing of Thurstonian models for paired comparison data under multiple judgment sampling. *Psychometrika* **66** 209–227. [MR1836935](#)
- MAYDEU-OLIVARES, A. (2002). Limited information estimation and testing of Thurstonian models for preference data. *Math. Social Sci.* **43** 467–483. [MR2073576](#)
- MAYDEU-OLIVARES, A. (2003). Thurstonian covariance and correlation structures for multiple judgment paired comparison data. Working Papers Economia, Instituto de Empresa, Area of Economic Environment. Available at <http://econpapers.repec.org/RePEc:emp:wpaper:wp03-04>.
- MAYDEU-OLIVARES, A. (2006). Limited information estimation and testing of discretized multivariate normal structural models. *Psychometrika* **71** 57–77. [MR2272520](#)
- MAYDEU-OLIVARES, A. and BÖCKENHOLT, U. (2005). Structural equation modeling of paired-comparison and ranking data. *Psychometrika* **10** 285–304.
- MAYDEU-OLIVARES, A. and BÖCKENHOLT, U. (2008). Modeling subject health outcomes. Top 10 reasons to use Thurstone’s method. *Medical Care* **46** 346–348.
- MAYDEU-OLIVARES, A. and HERNÁNDEZ, A. (2007). Identification and small sample estimation of Thurstone’s unrestricted model for paired comparisons data. *Multivariate Behavioral Research* **42** 323–347.
- MAYDEU-OLIVARES, A. and JOE, H. (2005). Limited- and full-information estimation and goodness-of-fit testing in 2^n contingency tables: A unified framework. *J. Amer. Statist. Assoc.* **100** 1009–1020. [MR2201027](#)
- MAYDEU-OLIVARES, A. and JOE, H. (2006). Limited information goodness-of-fit testing in multidimensional contingency tables. *Psychometrika* **71** 713–732. [MR2312239](#)
- MAZZUCCHI, T. A., LINZEY, W. G. and BRUNING, A. (2008). A paired comparison experiment for gathering expert judgment for an aircraft wiring risk assessment. *Reliability Engineering and System Safety* **93** 722–731.
- McFADDEN, D. (2001). Economic choices. *American Economic Review* **91** 351–378.
- McHALE, I. and MORTON, A. (2011). A Bradley–Terry type model for forecasting tennis match results. *International Journal of Forecasting* **27** 619–630.
- MEASE, D. (2003). A penalized maximum likelihood approach for the ranking of college football teams independent of victory margins. *Amer. Statist.* **57** 241–248. [MR2016258](#)
- MENKE, J. E. and MARTINEZ, T. R. (2008). A Bradley–Terry artificial neural network model for individual ratings in group competitions. *Neural Computing & Applications* **17** 175–186.
- MIWA, T., HAYTER, A. J. and KURIKI, S. (2003). The evaluation of general non-centred orthant probabilities. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **65** 223–234. [MR1959823](#)
- MOLENBERGHS, G. and VERBEKE, G. (2005). *Models for Discrete Longitudinal Data*. Springer, New York. [MR2171048](#)
- MOSTELLER, F. (1951). Remarks on the method of paired comparisons. I. The least squares solution assuming equal standard deviations and equal correlations. II. The effect of an aberrant standard deviation when equal standard deviations and equal correlations are assumed. III. A test of significance for paired comparisons when equal standard deviations and equal correlations are assumed. *Psychometrika* **16** 3–9, 203–218.
- MUTHÉN, B. (1978). Contributions to factor analysis of dichotomous variables. *Psychometrika* **43** 551–560. [MR0521904](#)
- MUTHÉN, B. (1993). Goodness of fit with categorical and other non normal variables. In *Structural Equation Models* (K. A. Bollen, J. S. Long, eds.) 205–234. Sage, Newbury Park, CA.
- MUTHÉN, B., DU TOIT, S. H. C. and SPISIC, D. (1997). Robust inference using weighted least squares and quadratic estimating equations in latent variable modeling with categorical and continuous outcomes. Technical report.
- R DEVELOPMENT CORE TEAM (2011). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria. ISBN 3-900051-07-0. Available at <http://www.R-project.org>.
- RAO, P. V. and KUPPER, L. L. (1967). Ties in paired-comparison experiments: A generalization of the Bradley–Terry model. *J. Amer. Statist. Assoc.* **62** 194–204. [MR0217963](#)
- REISER, M. (2008). Goodness-of-fit testing using components based on marginal frequencies of multinomial data. *British J. Math. Statist. Psych.* **61** 331–360. [MR2649040](#)
- SHAM, P. C. and CURTIS, D. (1995). An extended transmission/disequilibrium test (TDT) for multi-allele marker loci. *Ann. Hum. Genet.* **59** 323–336.
- SIMONS, G. and YAO, Y.-C. (1999). Asymptotics when the number of parameters tends to infinity in the Bradley–Terry model for paired comparisons. *Ann. Statist.* **27** 1041–1060. [MR1724040](#)
- SPRINGALL, A. (1973). Response surface fitting using a generalization of the Bradley–Terry paired comparison model. *J. R. Stat. Soc. Ser. C Appl. Stat.* **22** 59–68.
- STERN, H. (1990). A continuum of paired comparisons models. *Biometrika* **77** 265–273. [MR1064798](#)
- STERN, S. E. (2011). Moderated paired comparisons: A generalized Bradley–Terry model for continuous data using a discontinuous penalized likelihood function. *J. R. Stat. Soc. Ser. C Appl. Stat.* **60** 397–415. [MR2767853](#)
- STIGLER, S. M. (1994). Citation patterns in the journals of statistics and probability. *Statist. Sci.* **9** 94–108.
- STROBL, C., WICKELMAIER, F. and ZEILEIS, A. (2011). Accounting for individual differences in Bradley–Terry models by means of recursive partitioning. *Journal of Educational and Behavioral Statistics* **36** 135–153.
- STUART-FOX, D. M., FIRTH, D., MOUSSALLI, A. and WHITING, M. J. (2006). Multiple signals in chameleon contests: Designing and analysing animal contests as a tournament. *Animal Behavior* **71** 1263–1271.

- TAKANE, Y. (1989). Analysis of covariance structures and probabilistic binary choice data. In *New Developments in Psychological Choice Modeling* (G. De Soete, H. Feger and K. C. Klauer, eds.). North-Holland, Amsterdam.
- THURSTONE, L. L. (1927). A law of comparative judgment. *Psychological Review* **34** 368–389.
- THURSTONE, L. L. and JONES, L. V. (1957). The rational origin for measuring subjective values. *J. Amer. Statist. Assoc.* **52** 458–471.
- TRAIN, K. E. (2009). *Discrete Choice Methods with Simulation*, 2nd ed. Cambridge Univ. Press, Cambridge. [MR2519514](#)
- TSAI, R.-C. (2000). Remarks on the identifiability of Thurstonian ranking models: Case V, Case III, or neither? *Psychometrika* **65** 233–240. [MR1763521](#)
- TSAI, R.-C. (2003). Remarks on the identifiability of Thurstonian paired comparison models under multiple judgment. *Psychometrika* **68** 361–372. [MR2272384](#)
- TSAI, R.-C. and BÖCKENHOLT, U. (2002). Two-level linear paired comparison models: Estimation and identifiability issues. *Math. Social Sci.* **43** 429–449. [MR2072966](#)
- TSAI, R.-C. and BÖCKENHOLT, U. (2006). Modelling intransitive preferences: A random-effects approach. *J. Math. Psych.* **50** 1–14. [MR2208061](#)
- TSAI, R.-C. and BÖCKENHOLT, U. (2008). On the importance of distinguishing between within- and between-subject effects in intransitive intertemporal choice. *J. Math. Psych.* **52** 10–20. [MR2407792](#)
- TURNER, H. and FIRTH, D. (2010a). Bradley–Terry models in R: The `BradleyTerry2` package. Available at <http://CRAN.R-project.org/package=BradleyTerry2>.
- TURNER, H. and FIRTH, D. (2010b). Generalized nonlinear models in R: An overview of the `gnm` package. Available at <http://CRAN.R-project.org/package=gnm>.
- TVERSKY, A. (1972). Elimination by aspects: A theory of choice. *Psychological Review* **79** 281–299.
- USAMI, S. (2010). Individual differences multidimensional Bradley–Terry model using reversible jump Markov chain Monte Carlo algorithm. *Behaviormetrika* **37** 135–155.
- VARIN, C., REID, N. and FIRTH, D. (2011). An overview of composite likelihood methods. *Statist. Sinica* **21** 5–42. [MR2796852](#)
- WALKER, J. and BEN-AKIVA, M. (2002). Generalized random utility model. *Math. Social Sci.* **43** 303–343. [MR2072961](#)
- WHITING, M. J., STUART-FOX, D. M., O’CONNOR, D., FIRTH, D., BENNETT, N. C. and BLOMBERG, S. P. (2006). Ultraviolet signals ultra-aggression in a lizard. *Animal Behaviour* **72** 353–363.
- WICKELMAIER, F. and SCHMID, C. (2004). A Matlab function to estimate choice model parameters from paired-comparison data. *Behavior Research Methods, Instruments, and Computers* **36** 29–40.
- YAN, T., YANG, Y. and XU, J. (2012). Sparse paired comparisons in the Bradley–Terry model. *Statist. Sinica* **22** 1305–1318.
- ZERMELO, E. (1929). Die Berechnung der Turnier-Ergebnisse als ein Maximumproblem der Wahrscheinlichkeitsrechnung. *Math. Z.* **29** 436–460. [MR1545015](#)
- ZHAO, Y. and JOE, H. (2005). Composite likelihood estimation in multivariate data analysis. *Canad. J. Statist.* **33** 335–356. [MR2193979](#)