# A mixture model for longitudinal partially ranked data

| | |
|---|---|
| Journal: | *Communications in Statistics – Theory and Methods* |
| Manuscript ID: | LSTA-2013-0027.R1 |
| Manuscript Type: | Special Issue - New boundaries in statistical methods and models |
| Date Submitted by the Author: | n/a |
| Complete List of Authors: | Francis, Brian; Lancaster University, Department of Maths and Statistics Dittrich, Regina; WU Vienna, Institute for Statistics and Mathematics Hatzinger, Reinhold; WU Vienna, Institute for Statistics and Mathematics Humphreys, Les; Lancaster University, Department of Maths and Statistics |
| Keywords: | Nonparametric maximum likelihood, Latent class model, Paired comparisons, Bradley-Terry model, Mixture model, Partially ranked data |
| Abstract: | This paper discusses the use of mixture models in the analysis of longitudinal partially ranked data, where respondents, for example, choose only the preferred and second preferred out of a set of items. To model such data we convert it to a set of paired comparisons. Covariates can be incorporated into the model. We use a nonparametric mixture to account for unmeasured variability in individuals over time. The resulting multi-valued mass points can be interpreted as latent classes of the items. The work is illustrated by two questions on (post)materialism in three sweeps of the British Household Panel Survey. |

**SCHOLARONE™**
Manuscripts

# A mixture model for longitudinal partially ranked data

Brian Francis, Department of Mathematics and Statistics, Lancaster University, UK
Regina Dittrich, Institute for Statistics and Mathematics, WU Vienna, Austria
Reinhold Hatzinger, Institute for Statistics and Mathematics, WU Vienna, Austria
Les Humphreys, Department of Mathematics and Statistics, Lancaster University, UK

**Abstract**

This paper discusses the use of mixture models in the analysis of longitudinal partially ranked data, where respondents, for example, choose only the preferred and second preferred out of a set of items. To model such data we convert it to a set of paired comparisons. Covariates can be incorporated into the model. We use a nonparametric mixture to account for unmeasured variability in individuals over time. The resulting multi-valued mass points can be interpreted as latent classes of the items. The work is illustrated by two questions on (post)materialism in three sweeps of the British Household Panel Survey.

*Keywords*: Nonparametric maximum likelihood, Latent class model, Paired comparisons, Bradley-Terry model, Mixture model, Partially ranked data

Address for correspondence:

Brian Francis
Department of Mathematics and Statistics
Fylde College
Lancaster University, Lancaster, LA1 4YF
UK
E-Mail: `B.Francis@Lancaster.ac.uk`

# 1  Introduction

In this paper we are concerned with the modelling of a partial ranking of a set of items measured repeatedly over time on the same individuals. A partial ranking can be thought of as a ranking where at least two ranks are not collected. Such data might occur in a sample survey, with respondents being asked to choose the most preferred and the next most preferred out of a set of at least four opinions; alternatively this data can arise in elections, and in marketing experiments. The partial ranking might for example identify the top and the second items, with other items unranked – alternatively, it may identify the top and the bottom items. Formally, a partial ranking divides the $J$ items into $Q$ sets with $Q < J$ where the $Q$ sets are ordered, but items within a set are not. Our approach will be to model the partially ranked responses by converting them to a set of paired comparisons, and incorporating individual-level covariates. We will model the individual heterogeneity through a mass-point mixing distribution.

Our work is motivated by two questions in the British Household Panel Survey (Buck et al., 1994) which measure materialistic and postmaterialistic values (Inglehart, 1977) in repeated sweeps of the British Household Panel Survey. As part of the survey, respondents were asked to choose the most preferred and the next most preferred out of a set of four items representing preferred political priorities for the individual (see Figure 1).

In politics it is not always possible to obtain everything one might wish.
On this card several goals are listed.
If you had to choose among them, which would be your first choice?

|  | Highest Priority |
| --- | --- |
| Maintain order in nation | 1 |
| Give people more to say in government decisions | 2 |
| Fight rising prices | 3 |
| Protect freedom of speech | 4 |
| Can't choose | 8 |

And which would be your second choice?

Figure 1: The operationalisation of the Inglehart Index used in the British Household Panel Survey

Table 1 gives the responses for these two questions for the 1991 respondents. It can be seen

2

| Most Important Political Issue | Second Most Important Issue | | | | Missing | Total |
|---|---|---|---|---|---|---|
| | Maintain Order O | People more to say S | Fight rising prices P | Protect freedom of speech F | | |
| Maintain Order - O | 0 | 1027 | <u>1138</u> | 1136 | 14 | 3315 |
| People more to say - S | 1019 | 1 | 1019 | 1049 | 10 | 3098 |
| Fight rising prices- P | <u>688</u> | 779 | 0 | 385 | 8 | 1860 |
| Protect freedom of speech -F | 673 | 594 | 236 | 0 | 6 | 1511 |
| Missing | 1 | 0 | 0 | 0 | - | 1 |
| Total | 2381 | 2401 | 2395 | 2570 | 38 | 9785 |

Table 1: The responses to the two survey questions for the 1991 sweep

that there are 24 types of response which have a non-missing response to at least one of the two questions. Some response patterns are not observed in the survey.

Inglehart identified two of the items "maintain order in nation" (O) and "fight rising prices" (P) as representing materialistic values, with the items "give people more to say" (S) and "protect freedom of speech" (F) representing higher post-materialistic values. The underlying concept is that of a hierarchy of needs, with societies moving from the basic materialistic needs of order and stable prices to higher post-materialistic needs of democracy and rights. A typical analysis would be to calculate the proportion of those choosing the two materialistic items (represented by the two underlined numbers in Table 1) and the proportion of those choosing the two post-materialistic items (represented by the two boxed numbers in Table 1), and to examine changes in these proportions over time. For the 1991 data presented in Table 1 the materialism and postmaterialism proportions would be $\frac{1138+688}{9785} = 0.187$ and $\frac{1049+594}{9785} = 0.168$ respectively. However such an approach has numerous problems - it ignores most of the responses - only four of the 24 cells are used - it ignores the ranked nature of the responses, does not take into account the longitudinal panel design of the survey, and fails to include covariates. In the next section we describe a more appropriate analytic method for longitudinal partial rank data.

3

## 2   Modelling partially ranked data

There are a number of existing methods for modelling partially ranked data. Methods based on metric distances or ranks have been proposed by Critchlow (1980); more recently a mixture of a shifted binomial model combined with a uniform distribution known as the CUB model has been proposed for ranked data (D'Elia and Piccolo, 2005). Both methods however do not focus on the worth or importance of the ranked objects. Among utility-based models, there are two main approaches - the stated choice model (Chapman and Staelin, 1982) and the paired comparison model. This paper uses the paired comparison model, which assumes that the rankings are produced by individuals making internal comparisons of all sets of objects.

The method of paired comparisons, introduced by Thurstone (1927) and popularized by Bradley and Terry (1952), was designed to measure the relative importance or worth of a set of items. Essentially, with $J$ items, each pair of items is taken, and respondents are asked to judge which of the two is most important. In this paper we are concerned with partially ranked items, where respondents are instead asked to determine the partial ordering of a set of items.

It is straightforward to transform fully ranked data into paired comparison (PC) form (Dittrich et al., 2000). Francis et al. (2010) have pointed out the advantages of a PC approach to ranked data compared to the main alternative of choice-based models. The main advantages are that the problem is placed in the standard framework of generalized linear models and respondent covariates can also be incorporated. Additionally, the PC model for ranked data is invariant to the decision process of the respondent whereas a choice-based model assumes that the respondent answers the ranking questions in order.

However, dealing with partially ranked data is more problematic, as the full rank ordering is unknown. Francis et al. (2002) describes the conversion of partial rank data to paired comparisons. For example, in the motivating example above, choosing item O followed by item F will generate six paired comparisons, with O preferred to F,S and P, F preferred to S and P and with no preference between S and P. Similar considerations enable each of the 24 responses to be converted to a set of paired comparison responses including repeated responses and partially missing responses. These latter response patterns will have more "no-preference" paired comparisons as there is only one

4

ranked item.

In this paper, we follow the development of Francis et al. (2010) but consider partial rank responses rather then full rank responses, and a longitudinal rather than cross-sectional design. We refer to the Francis et al. (2010) as the *LARA1* model (LAtent RAnks 1-level) and the model in this paper as the *LARA2* model (LAtent RAnks 2-levels).

We start with a simple model for partial ranks, ignoring for the time being covariates and the longitudinal nature of the survey. Following the *LARA1* model, we assume a multinomial model for the $L$ observed response patterns. Let $N_{\ell i}$ be an $(0, 1)$ indicator as to whether a specific response pattern $\ell$ is observed (1) or not (0) for subject $i$, with $\sum_{\ell=1}^{L} N_{\ell i} = 1$. Then the $N_{\ell i}$ are multinomially distributed, and the likelihood function up to a normalising constant is

$$\text{LIK}_i = \prod_{\ell=1}^{L} P_{\ell i}^{N_{\ell i}} . \tag{1}$$

The probability $P_\ell$ of a specific response pattern $\ell$ is then given by the product over all of the derived PCs. Our model is similar to the Mallows-Bradley-Terry ranking model (Mallows, 1957; Critchlow and Fligner, 1991) and was described in Dittrich et al. (2007). In their model the probability of a response ranking of the items is taken to be proportional to the product of the probabilities of all pairwise comparisons that are consistent with the ranking. For subject $i$, we can then write

$$P_{\ell i} = P_\ell(y_{(12)i}, y_{(13)i}, \ldots, y_{(J-1:J)i}) = \prod_{j<k} P(y_{(jk)i}) . \tag{2}$$

Following Davidson (1970) and Sinclair (1982), the probability for a single PC response $y_{jk}$ between items $j$ and $k$ is defined by

$$P(y_{(jk)}) = a_{jk} \, c_{jk}^{(1-y_{jk}^2)} \left( \frac{\sqrt{\pi_j}}{\sqrt{\pi_k}} \right)^{y_{jk}} ,$$

where $y_{(jk)}$ takes the value of 1, if item $j$ is preferred to $k$, takes the value of $-1$, if item $k$ is preferred to $j$, and takes the value of zero if no preference is stated. In this model, the parameters of specific interest are the $\pi_j$, $j = 1, \ldots, J$ which represent a set of worths or importances of the items. For identifiability, we define $\sum_j \pi_j = 1$ . The $c_{jk}$ are a set of parameters which represent a

5

different probability of no preference for each pair of responses, and $a_{jk} = a_{jk}(\pi_j, \pi_k)$ is a normalising quantity for the comparison $jk$.

This gives

$$P_\ell = \prod_{j<k} a_{jk} \; c_{jk}^{(1-y_{jk}^2)} \; \left(\frac{\sqrt{\pi_j}}{\sqrt{\pi_k}}\right)^{y_{jk}} . \tag{3}$$

The normalising quantity is now $\prod_{j<k} a_{jk}$, which is the same for all patterns. We let $m_{\ell i}$ be the expected value of $N_i P_{\ell i}$. Converting to log-linear form, we have

$$\ln(m_{\ell i}) \;\; = \;\; \phi_i + \sum_{j<k} \left[(1 - y_{(jk)i}^2) \ln(c_{jk}) + y_{(jk)i}(\frac{1}{2} \ln(\pi_j) - \frac{1}{2} ln(\pi_k))\right] \tag{4}$$

$$= \;\; \phi_i + \sum_{j<k} \left[(1 - y_{(jk)i}^2)\psi_{jk} + y_{(jk)i}(\lambda_j - \lambda_k)\right] . \tag{5}$$

For identifability, $\lambda_J$ is set to zero. All of the $\psi_{jk}$ are estimable. The $\phi_i$s are so-called nuisance parameters. The $\lambda$s (location of the preference parameters) are related to the $\pi$s by $\ln \pi = 2\lambda$, and where $\psi_{jk} = \ln(c_{jk})$. The $\pi_k$ can be calculated from the estimated $\lambda$s through the formula

$$\pi_k = \frac{\exp(2\lambda_k)}{\sum_{j=1}^{J} \exp(2\lambda_j)} .$$

The above multinomial model can thus be fitted as a Poisson log-linear model using the standard Poisson-multinomial equivalence, with the constraint that $\sum_\ell m_{\ell i} = 1$ for each $i$. This equivalence requires that the sum of the fitted probabilities over all patterns for each individual is one. As, for each individual, there is only one pattern chosen, and the sum of the observed counts ($N_i$) is one, this is achieved by adding a set of extra parameters to the model which are provided by the $\phi_i$.

Subject covariates $x_{is}$, $s = 1, \ldots, S$ can be included in the model as interaction terms with the items, giving a set of interaction terms $x_{i1s}, x_{i2s}, \ldots, x_{iJs}$ for each covariate $s$, which are added to the linear predictor. The linear model becomes

$$\ln(m_{\ell i}) = \phi_i + \sum_{j<k} \left[(1 - y_{(jk)i}^2)\psi_{jk} + y_{(jk)i}(\lambda_{ji} - \lambda_{ki})\right] . \tag{6}$$

6

where, with $S$ covariates,

$$\lambda_{ji} = \lambda_j + \sum_{s=1}^{S} \beta_{ijs} x_{js} \,.$$

Note that for every covariate $x_{is}$, there are $J$ parameters to estimate - the differential effect of the covariate on each of the items. Typically, for identification, the last of the $J$ items is treated as a reference item, and the parameters associated with this item ($\lambda_J$ and the $\beta_{Js}$) are set to zero.

## 3 Mixture models for longitudinal partially ranked data

We now extend this model to allow for repeated observations of responses over time, which will both allow us to examine change over time in the ranked responses, and to take account of individual heterogeneity. We extend the notation by adding the subscript $t$ ($t = 1, \ldots, T$) for the $T$ time points. We assume that there is a random individual effect which is constant over the sweeps of the survey. The random effect is multivalued as there is a separate random effect for each item (the random effect for the last item is set to zero for identifability). The covariate model is now

$$\lambda_{jit} = \lambda_j + \sum_{s=1}^{S} \beta_{js} x_{ijs} + u_{ij} \,,$$

where $u_{ij}$ is the random effect for individual $i$ on item $j$.

Rather than assuming a multivariate distribution for the random effects, we use a nonparametric $R$-component mass point formulation of the random effects structure, with unknown mass point locations $\Delta_r = (\delta_{r1}, \delta_{r2}, \ldots, \delta_{rJ})$ and masses $q_r$ with $r = 1, \ldots, R$. Each discrete mass point is multi-valued, with a parameter for each item $j$ (Francis et al., 2010). The model is similar to the LARA1 model, but the random effects now represent individual level variability which is constant over time, rather than overdispersed response variability. The resultant model is equivalent to a latent class regression model, where the latent class profiles are provided by the mass point components and the covariates act on the class profiles. This provides an alternative interpretation of the fitted model.

The likelihood for the latent class regression model is

7

$$\mathsf{LIK} = \prod_{\ell it} \Big( \sum_{r=1}^{R} q_r \, P_{\ell itr}(y_{\ell it}|\Delta_r) \Big)^{N_{\ell it}} \quad \text{where} \quad \sum_{\ell} P_{\ell itr} = 1 \quad \forall \, i, r, t \,. \tag{7}$$

We constrain $\sum_r q_r = 1$.

The model for this latent class approach with covariates can be written as

$$\ln(m_{\ell itr}) = \phi_{it} + \sum_{j<k} \left[ (1 - y_{(jk)it}^2)\psi_{jk} + y_{(jk)i}(\lambda_{jitr} - \lambda_{kitr}) \right] \,, \tag{8}$$

where $\lambda_{jitr}$ is now

$$\lambda_{jitr} = \lambda_j + \sum_{s=1}^{S} \beta_{js} x_{ijs} + \delta_{rj} \,.$$

$\lambda_j$ is adjusted by $\delta_{jr}$ for each mass point $r$ and item $j$.

The model deals with attrition over the sweeps of the survey by taking a full information maximum likelihood [FIML] approach which assumes an underlying missing at random [MAR] process (see e.g. Enders, 2001). Thus, all observed data are included in the analysis, whether the individual contributes 1, 2 or 3 responses.

## 3.1   Algorithmic issues

The model is fitted using the EM algorithm (Aitkin, 1999). McLachlan and Peel (2000, p.49) give regularity conditions that need to be satisfied for roots of the likelihood equation to exist for any mixture model. The latent class membership indicators for each individual can be treated as missing data. We can write these as $z_{ir}$, with $z_{ir} = 1$ if individual $i$ belongs to class $r$, and zero otherwise. The expected values of the $z$s are defined to be $w_{ir}$ and are the posterior probabilities of class membership for a respondent $i$. The E-step of the EM algorithm computes the conditional expectation of the complete log-likelihood (involving the calculation of the $w$s), whereas the M-step maximizes the multinomial likelihood with respect to the $\lambda$s and $\delta$s, given the current expected values of the $z$s. The model is implemented through an expanded Poisson log-linear model with weights $w_{ir}$. Fitting the multinomial through a Poisson log-linear model necessitates that a set of nuisance parameters be included in the linear predictor; these constrain the marginal totals over

8

patterns for each individual and sweep to be equal to 1. They are dealt with numerically by using the method of Hatzinger and Francis (2004), which provides an efficient numerical method for fitting large numbers of nuisance parameters. It is usual to start with random values of the $w_{ir}$ with the constraint $\sum_r w_{ir} = 1$.

The problems of fitting latent class models are well known. The first is that of multiple maxima. The EM algorithm may not converge a global solution. To minimize this problem, a number of different starting values are taken for each value of $R$ and for each covariate model, and the model with the lowest value of the deviance( -2 log likelihood) taken (McLachlan and Peel, 2000; Magidson and Vermunt, 2004) .

The model was fitted using the `prefmod` package (Hatzinger and Dittrich, 2012) and the `allvc` function of the `npmlreg` package (Einbeck et al., 2012) of R. The `allvc` function was edited to allow for random start values of the $w_{ir}$.

# 4   An illustrative example

We return to the two Inglehart questions from the British Household Panel Survey. We take repeated responses from three time points (1991,1993 and 1995) to illustrate the method for longitudinal data. Our analysis of similar Inglehart questions from the International Social Science Program 2000 (Francis et al., 2002) identified a number of covariates which were important in explaining changes in response patterns. As well as country of residence, which is not relevant for the current study, these were age, education and gender. To this list we added marital status and year, with the latter allowing us to examine changing worths over year. All covariates apart from year were treated as time stable and measured at the first observed time point.

Table 2 gives the observed number of respondents for each observed response pattern over the three time points. $9,804$ respondents aged 15 or over are surveyed in 1991, and while some drop out, new respondents are surveyed to replace them at later sweeps. In total, there are $11,728$ respondents who contributed between one and three observations, with $27,228$ time-point observations in total. The covariates used, which were all treated as factors, are listed as follows, with sample percentages in square brackets.

Table 2: Response patterns over time: (three time points)

| No. of respondents | 1991 | 1993 | 1995 |
|---:|:---:|:---:|:---:|
| 6810 | √ | √ | √ |
| 921 | √ | √ | |
| 285 | √ | | √ |
| 679 | | √ | √ |
| 1769 | √ | | |
| 411 | | √ | |
| 853 | | | √ |
| Total: 11728 | 9785 | 8821 | 8627 |

- year The year of the observation (1991, 1993, 1995)

- age The age of the respondent in 1991 (under 29 [33.4%], 30-44 [27.0%], 45-64 [24.2%], 65 and over [15.4%])

- edu Highest educational level achieved (no qualification [30.7%], O-level or equivalent [28.7%], A-level or equivalent [16.9%], degree or equivalent [23.8%] )

- sex Gender of respondent (male [46.8%], female [53.2%] )

- mar Marital status of respondent (married or living together [61.5%] never married [25.6%], was married [12.9%] )

# 5  Analysis and results

Our approach to model fitting was to fit a full main-effects model (time+age+sex+mar+edu) interacted with the items (S+P+O+F) with a single latent class ($K = 1$) and then to increase the number of latent classes, examining the BIC value for each value of $K$. Each model was fitted five times with different starting values, and the deviances (minus twice log-likelihood) examined. More complex models may require a larger number of starting values, but five proved adequate for the analysis here.

Table 3 gives the best deviances for $K = 1, 2, 3$ and the equivalent BIC values. For $K = 2$ and $K = 3$ the best deviance was found for four of the five sets of random starting values used. There

10

Table 3: Best deviances and BICs from five different starting values

|  | deviance | no. of parameters | BIC |  |
|---|---|---|---|---|
| 1 latent class model | 239037 | 43 | 239476 |  |
| 2 latent class model | 235240 | 48 | 235731 |  |
| 3 latent class model | 232124 | 53 | 232665 | $\star$ |

Table 4: Change in deviance from main effects model: `time+age+sex+mar+edu`

|  | Deviance | no. of parameters | Deviance change from main effects model | Change in *df* | LRT p-value |
|---|---|---|---|---|---|
| `year+age+sex+mar+edu` | 232123.8 | 53 |  |  |  |
| − `edu` | 232766.1 | 44 | 642.3 | 9 | < 0.001 |
| − `age` | 232435.3 | 44 | 311.5 | 9 | < 0.001 |
| − `mar` | 232202.4 | 47 | 78.6 | 6 | < 0.001 |
| − `sex` | 232235.1 | 50 | 111.3 | 3 | < 0.001 |
| − `year` | 232236.7 | 47 | 112.9 | 6 | < 0.001 |

is strong evidence that three latent classes are needed, with the BIC decreasing dramatically from 1 to 3 classes. The three latent class model with a full main effects model took about ten hours to fit on a high memory 40Gb processor, and in this illustrative example, we do not increase the number of latent classes further.

We then proceeded to test for the effect of each of the five covariates, by examining deletion deviances. All models were again fitted with five different random start values, and the model with the lowest deviance selected. Table 4 contains the results. The removal of the sex covariate from the model removes three interactions, `sex` with S, `sex` with P and `sex` with F and three parameters are lost. The change of deviance of 111.3 on 3 df is highly significant (LRT $p < 0.001$) and the covariate `sex` is retained. Table 4 shows that all covariates are similarly needed in the model. However, judged by the average decrease of deviance per degree of freedom, education is the most important covariate affecting the item worths, followed by age and then by sex.

We now move to interpretation of the three class main effects model. Figure 2 shows the results of fitting the three latent class final model, showing the estimated worths of the items for each latent class at the baseline level of the other covariates (aged under 19, married males with no

11

educational qualifications), and the effect of year on these estimated worths. We can see first of all that the three classes have estimated class sizes $q_r$ (obtained from $q_r = \sum_i w_i r$ ) of $0.316, 0.146$ and $0.539$ respectively. The largest class – class 3 – has the highest worth for "more to say in government" (S) followed by "fight rising prices" (P) – this is neither a materialist nor a materialist group but a mixed group. The next largest class (class 1) is identified as a materialist group, with "maintain order" (O) and "fight rising prices" (P) having the highest worths. Finally, class 2 can be seen to be a postmaterialist group, with "protect freedom of speech" (F) and "more to say in government" (S) having the highest worths.



Figure 2: Item worths for three latent classes

Examining changes by year of survey, we can observe that there is a tendency for "maintain order in the nation" to increase its worth over the three sweeps, and for "protect freedom of speech" to decrease. We emphasise that our model assumes that the covariates act identically on each latent class, and therefore the same decrease in F and increase in O can be seen for each of the three latent classes. The increase of "order" over the sweeps of the survey may be possibly related to the UK withdrawing from the European Exchange Rate Mechanism in 1992 after severe currency speculation against the UK pound and a failure of the government policy of spending £27 billion of reserves in trying to maintain the value of the pound. (This event, known as "black Wednesday"

12

1
2
3
4
5
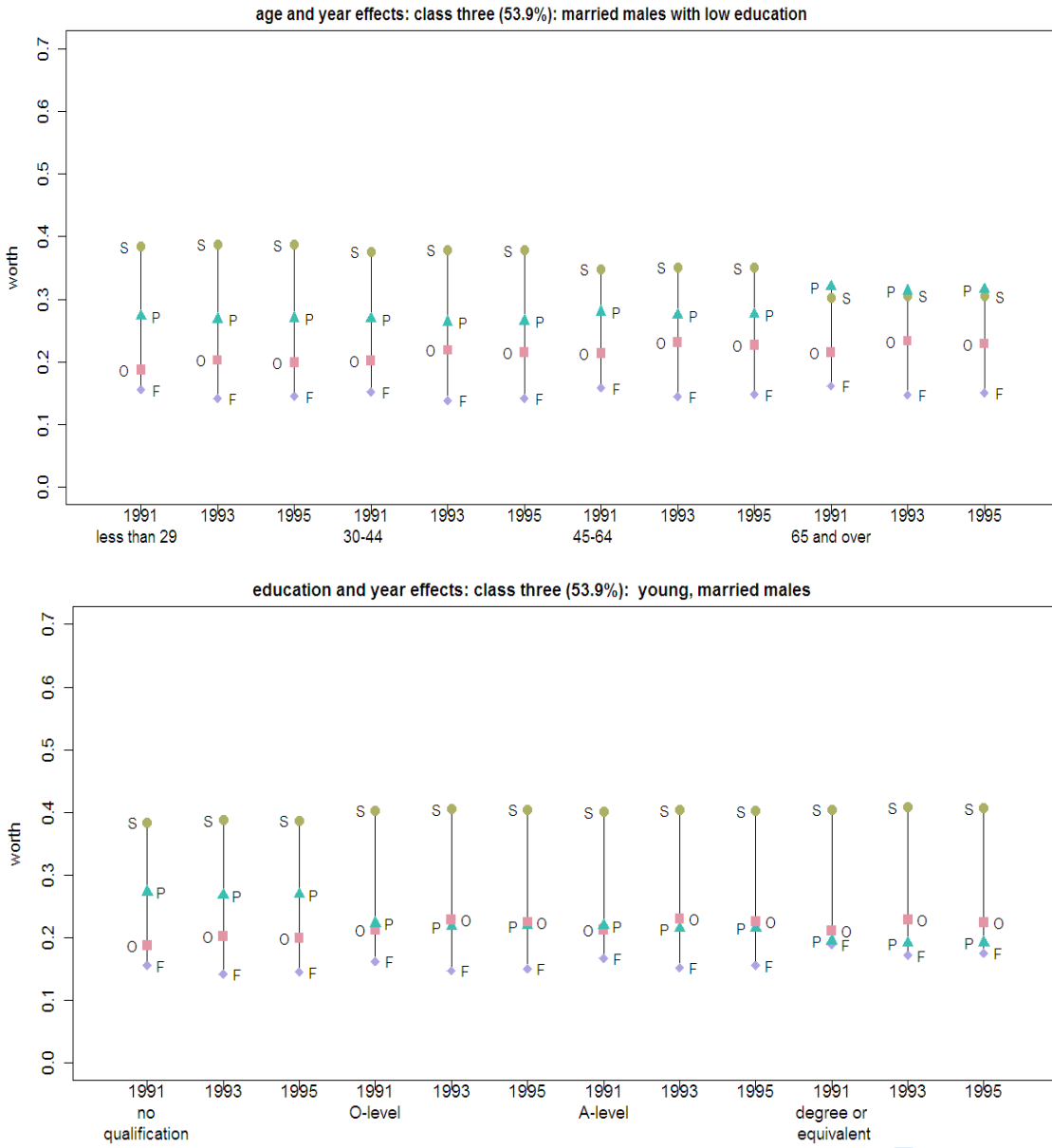6   caused a de facto devaluation of the pound and insecurity in the country).
7
8
9



Figure 3: Item worths by age and education for the biggest latent class

13

Table 4 identified that highest educational qualification and age are the most important covariates affecting the item worths. Figure 3 shows the effects of these covariates on the worths and also showing the effect of yearly sweep. The effects are shown just for the largest latent class, namely class 3.

Examining the age effects first of all, (the top plot in Figure 3) we see that the materialist items ("fight rising prices" and "maintain order") have higher worths for the older age groups, once gender, educational level and marital status are controlled for. We take this to be a generational effect, but it could equally be an aging effect, with respondents becoming more materialist as they age. The bottom plot in Figure 3 shows the effect of highest qualification, controlling for age, gender and marital status. The main observed change is a strong decrease in the worth of "fight rising prices" as educational level increases, with consequent smaller increases for "maintain order" and "protect freedom of speech". It is possible that education is acting as a proxy for income, and thus the effect of rising prices may be felt less for those earning more. However, the increase of worth of "order" and "freedom of speech" is interesting – the first is a materialist item and the second is a post-materialist item. Increasing education seems to lead both to a desire for more security and also to a need for freedom in voicing concerns.

The two remaining covariates show less strong effects and the estimates are not shown. In examining the marital status effect, the main differences occur between the "was married" category and the other two categories. The "was married" category (consisting of divorced, separated and widowed respondents) have a lower worth for "maintain order" and a higher worth for "fight rising prices" compared to the other two marital status groups. The changes however are in the relative importance for two materialistic items rather than a move from post-materialism to materialism. The gender effect is characterised in a similar way – females, similar to the "was married" group – have a lower worth for "maintain order" and a higher worth for "fight rising prices" than their male counterparts. Again the changes seem to be mainly in the relative importance of the two materialistic items.

14

# 6   Discussion

The use of latent class models for analysing repeated partial rank data with an underlying paired comparison model structure has numerous advantages. Firstly, the paired comparison method means that the order in which the questions are asked and answered is ignored, as the method assumes that a rank consensus is made of all items before answering the survey questions. Secondly, the modelling is placed in the general structure of a generalised linear model. Thirdly, the use of a discrete mixing distribution for dealing with a multi-valued individual random effects term simplifies the algorithm considerably as Francis et al. (2010) highlights. Finally, the algorithm in our experience seems relatively stable.

The disadvantage is primarily the speed of the algorithm. Our most complex model (a three latent-class model with 27228 observations, four items and five categorical covariates took around 10 hours of CPU time on a fast processor). It seems feasible that more latent classes are needed but CPU time constraints meant that we could not explore this – the analysis presented in this paper is therefore illustrative rather than the final word.

Some comments should also be made about the attrition process over the sweeps of the survey. Our approach assumes a missing at random process. However, it is possible that the attrition of survey respondents might be informative missing - that the missing partial rank responses of those who have withdrawn from one or more sweeps of the study are more likely to be of a particular form. There are two approaches that could be taken here. The first is to use the response pattern shown in Table 2, (or some summary such as the number of sweeps responded to) as an additional covariate in the model (Hedeker and Gibbons, 2006). This will identify any relationship between item ranking and response pattern. A second approach would be to extend the approach developed by Dittrich et al. (2012), who used composite link models to fit paired comparison models, into similar models for ranked longitudinal panel data. This second approach is under development by the authors.

The flexibility of the algorithm allows for the model to be extended in various ways. It is straightforward to have class-dependent covariates by including an interaction of the class member-ship variable with the desired covariate or covariates. Similarly, we could have included no-preference

15

effects which depended on covariates, by including appropriate interactions in the model formula.

# 7   Acknowledgements

# References

Aitkin, M. (1999). A general maximum likelihood analysis of variance components in generalized linear models. *Biometrics 55*, 117–128.

Bradley, R. and M. Terry (1952). Rank Analysis of Incomplete Block Designs. I. The Method of Paired Comparisons. *Biometrika 39*, 324–345.

Buck, N., J. Gershuny, D. Rose, and J. Scott (1994). *Changing Households: The British Household Panel Survey 1990-1992*. Colchester: ESRC Research Centre on Micro-Social Change, University of Essex.

Chapman, R. and R. Staelin (1982). Exploiting rank ordered choice set data within the stochastic utility model. *Journal of Marketing Research 19*(3), 288–301.

Critchlow, D. (1980). *Metric Methods for Analysing Partially Ranked Data*. Berlin: Springer-Verlag.

Critchlow, D. and M. Fligner (1991). Paired comparison, triple comparison, and ranking experiments as generalized linear models, and their implementation in GLIM. *Psychometrika 56*, 517–533.

Davidson, R. (1970). Extending the Bradley-Terry model to accommodate ties in paired comparison experiments. *JASA 65*, 317–328.

16

D'Elia, A. and D. Piccolo (2005). A mixture model for preferences data analysis. *Computational Statistics and Data Analysis 49*, 917–934.

Dittrich, R., B. Francis, R. Hatzinger, and W. Katzenbeisser (2007). A paired comparison approach for the analysis of sets of Likert scale responses. *Statistical Modelling 7*, 3–28.

Dittrich, R., B. Francis, R. Hatzinger, and W. Katzenbeisser (2012). Missing observations in paired comparison data. *Statistical Modelling 12*(2), 117–143.

Dittrich, R., W. Katzenbeisser, and H. Reisinger (2000). The analysis of rank ordered preference data based on Bradley-Terry type models. *OR Spectrum 22*(1), 117–134.

Einbeck, J., R. Darnell, and J. Hinde (2012). *npmlreg: Nonparametric maximum likelihood estimation for random effect models*. R package version 0.45.

Enders, C. K. (2001). The performance of the full information maximum likelihood estimator in multiple regression models with missing data. *Educational and Psychological Measurement 61*(5), 713–740.

Francis, B., R. Dittrich, and R. Hatzinger (2010). Modeling heterogeneity in ranked responses by nonparametric maximum likelihood: How do europeans get their scientific knowledge? *The Annals of Applied Statistics 4*(4), 2181–2202.

Francis, B., R. Dittrich, R. Hatzinger, and R. Penn (2002). Analysing ranks using paired comparison methods: an investigation of value orientation in Europe. *Applied Statistics 51*(3), 319–336.

Hatzinger, R. and R. Dittrich (2012). prefmod: An R package for modeling preferences based on paired comparisons, rankings, or ratings. *Journal of Statistical Software 48*(10), 1–31.

Hatzinger, R. and B. Francis (2004). Fitting paired comparison models in R. Technical Report 3 - http://epub.wu-wien.ac.at/740/, Institute for Statistics and Mathematics, WU Vienna.

Hedeker, D. and R. Gibbons (2006). *Longitudinal Data Analysis*. New York: John Wiley.

Inglehart, R. (1977). *The silent revolution : changing values and political styles among Western publics*. Princeton: Princeton University Press.

17

Magidson, J. and J. Vermunt (2004). Latent class models. In D. Kaplan (Ed.), *The Sage Handbook of Quantitative Methodology for the Social Sciences*, Thousand Oaks, pp. 175–198. Sage.

Mallows, C. (1957). Non-null ranking models: I. *Biometrika 44*, 114–130.

McLachlan, G. and D. Peel (2000). *Finite Mixture Models*. New York: Wiley.

Sinclair, C. (1982). GLIM for preference. In R. Gilchrist (Ed.), *Proceedings of the International Conference on Generalised Linear Models*, Volume 14, pp. 164–178. Springer Lecture Notes in Statistics.

Thurstone, L. (1927). A law of comparative judgement. *Psychological Review 34*, 273–286.

18