



The heterogeneous thresholds ordered response model: identification and inference

Franco Peracchi

Tor Vergata University and Einaudi Institute for Economics and Finance, Rome, Italy

and Claudio Rossetti

Libera Università Internazionale degli Studi Sociali 'Guido Carli', Rome, Italy

[Received November 2011. Revised April 2012]

Summary. Although surveys routinely ask respondents to evaluate various aspects of their life on an ordered scale, there is concern about interpersonal comparability of these self-assessments. Statistically, the problem is one of identification in ordered response models with heterogeneous thresholds. As a solution to the identification problem, King and his colleagues proposed the use of anchoring vignettes, namely brief descriptions of hypothetical people or situations that survey respondents are asked to evaluate on the same scale as they use to rate their own situation. Although vignettes have been introduced in several social surveys and are increasingly employed in a variety of fields, the reliability of this approach hinges crucially on the validity of the assumptions of response consistency and vignette equivalence. The paper proposes a joint test of these key assumptions based on the fact that the underlying statistical model is overidentified if the two assumptions hold. Monte Carlo results show that the test proposed has good size and power properties in finite samples. We apply our test to self-assessment of pain by using data from the first wave of the Survey of Health, Ageing and Retirement in Europe. We find that, when using only one of the three available vignettes, or when the test is carried out separately by subgroups of respondents, the overidentifying restrictions are less likely to be rejected.

Keywords: Anchoring vignettes; Differential item functioning; Minimum distance methods; Ordered response models; Reporting heterogeneity; Self-assessment of health

1. Introduction

Surveys respondents are often asked to evaluate various aspects of their life on an ordered scale. Examples include questions on life satisfaction and self-rated health in household surveys or questions on customer satisfaction in consumer surveys. Although such questions are widely used, there is a concern that different people may interpret and answer them differently. This is especially true when comparing subjective assessments across groups that are characterized by different culture, nationality, socio-economic status, age or gender. For example, when asked to rate their own health on a given categorical scale, people may answer differently because their true or perceived health differs, but also because they interpret differently the various levels of the scale. As a consequence, differences in self-reports between otherwise similar individuals may depend on differences in style of response, namely the mapping of true or perceived health into reported health (Sen, 2002). Lack of interpersonal comparability of responses to subjective survey questions is often referred to as 'differential item functioning' (DIF), which is a term

Address for correspondence: Franco Peracchi, Department of Economics and Finance, Tor Vergata University, via Columbia 2, Rome 00133, Italy.
E-mail: franco.peracchi@uniroma.it

that originated in the educational testing literature (Holland and Wainer, 1993) where a test question is said to have DIF if equally able individuals have unequal probabilities of answering the question correctly. From the view point of statistical modelling, the DIF problem is essentially one of identification in ordered response models where the observed responses are derived from latent continuous random variables discretized through a set of heterogeneous thresholds or cut-off points.

Following the seminal paper of King *et al.* (2004), anchoring vignettes have been developed as a new component of survey instruments that may be used to solve the DIF problem. They are brief descriptions of hypothetical people or situations that survey respondents are asked to evaluate on the same scale as they use to rate their own situation. Because the people or situations that are described in the vignettes are the same for all respondents, vignettes have the potential to identify individual variation in subjective thresholds. Several social surveys such as the Survey of Health, Ageing and Retirement in Europe (SHARE), the US Health and Retirement Study, the English Longitudinal Study of Ageing and the World Health Organization's World Health Surveys have introduced specific modules with vignette questions. However, introducing anchoring vignettes implies substantial costs in terms of survey design and reduces the time that is available for collecting other information. Of course, anchoring vignettes would not be necessary in a survey if one is willing to apply the response scale correction from a different survey under the assumption that the DIF problem is the same.

Vignette questions have been applied to a variety of problems including the comparison of health (Salomon *et al.*, 2004; King and Wand, 2007; Bago d'Uva, O'Donnell and van Doorslaer, 2008; Bago d'Uva, van Doorslaer, Lindeboom and O'Donnell, 2008; Peracchi and Rossetti, 2009), health system responsiveness (Rice *et al.*, 2012), political efficacy (King *et al.*, 2004), work disability (Kapteyn *et al.*, 2007), life satisfaction (Angelini *et al.*, 2008) and job satisfaction (Kristensen and Johansson, 2008). In most cases, evidence of reporting heterogeneity is found and corrections on the comparisons of interest are made by using the vignettes.

Although vignettes are increasingly employed by researchers in various fields, reliability of this approach hinges crucially on the validity of two key assumptions (King *et al.*, 2004). The first assumption ('response consistency') is that individuals use the available response categories in the same way when assessing their own situation and the hypothetical situations in the vignettes. The second assumption ('vignette equivalence') is that the hypothetical situation in a vignette is perceived by all respondents in the same way and on the same unidimensional scale, apart from random error. As pointed out by Deaton (2011), the vignette approach replaces the assumption that there are no differences in the way that people rank themselves on a subjective scale with the alternative assumption (response consistency) that there are no differences in their capacity for empathy with other people's conditions. In addition, vignette equivalence assumes that there are no systematic differences in the way that people perceive the situations that are represented in each vignette. This is also a very strong assumption, e.g. because of problems with translation of the same vignette in different languages. Hence, testing these two key assumptions turns out to be a critical step in evaluating the validity of the vignette approach.

One approach to testing for response consistency relies on the availability of some objective measure of the concept of interest. This approach, which rests on the maintained assumption of vignette equivalence, was used by King *et al.* (2004) and van Soest *et al.* (2011) to provide evidence supporting the assumption of response consistency. Other evidence, however, is less supportive (Datta Gupta *et al.*, 2010; Bago d'Uva *et al.*, 2011; Voňková and Hullelgie, 2011). The main problem with this approach is that objective measures of the concept of interest are typically only available in *ad hoc* studies. Recently, Kapteyn *et al.* (2011) proposed a different test of response consistency based on longitudinal data where respondents are shown vignettes that

are descriptions of their own health collected in a previous interview. They found that response consistency is satisfied only for one of the five health domains that are considered, namely sleep.

Far less attention has been paid to vignette equivalence. King *et al.* (2004) suggested an informal test based on the ordering of the answers to different vignette questions on the same domain. A more formal approach was adopted by Bago d'Uva *et al.* (2011), who tested the necessary condition of no systematic variation across individuals by allowing vignette evaluations to depend on observed personal characteristics. This test does not require objective measures but maintains the assumption of response consistency and needs at least two vignettes questions for each concept of interest.

In this paper we propose a simple joint test of the two key assumptions of response consistency and vignette equivalence. The test proposed exploits the fact that, as pointed out by Deaton (2011), the statistical model is overidentified under these two assumptions. Our test offers several advantages. First, it does not require the availability of some objective measures and can be carried out by using any data set containing at least one vignette question for each concept of interest. Second, it does not require embedding the restricted model that imposes response consistency and vignette equivalence in a larger encompassing model. Third, it requires only a consistent and asymptotically normal estimator of the estimable parameters in the model. This is an advantage, both computationally and because the test can easily be extended to models with sample selection and to semiparametric settings where strong distributional assumptions are relaxed. Fourth, because it exploits the mapping between the estimable parameters and the full set of model parameters, imposing additional restrictions on the model is particularly transparent and simple. Of course, as typical with tests of parametric or semiparametric models, our test is conditional on some other assumptions. Thus, it may reject the overidentifying restrictions for other reasons than failure of response consistency and vignette equivalence, e.g. because of failure of parametric restrictions or because relevant variables have been omitted from the model.

We investigate the finite sample performance of the proposed test through a Monte Carlo study. We find that the test has good size and power properties in finite samples. Specifically, the test has no size distortion and no overrejection is reported when the number of overidentifying restrictions increases.

Finally, we apply our test to self-assessment of pain by using data from release 2 of the first (2004) wave of the SHARE. Release 2 of the data also includes the answers to vignettes questions in a self-administered questionnaire submitted to a randomly selected subsample of respondents. We find that the overidentifying restrictions are less likely to be rejected when using only one of the three available vignettes, or when the test is carried out separately by subgroups of respondents.

The remainder of this paper is organized as follows. Section 2 presents the heterogeneous thresholds ordered response model, discusses its identification and proposes a test of the overidentifying restrictions that are implied by the assumptions of response consistency and vignette equivalence. Section 3 presents the results of a Monte Carlo study to assess the finite sample performance of the test proposed. Section 4 illustrates the use of our test through an empirical application to self-assessment on various health domains. Finally, Section 5 offers some conclusions.

2. Heterogeneous thresholds ordered response model

Let Y_0 denote the answer by a randomly chosen individual on some concept of interest, and let Y_1, \dots, Y_J denote the answers given by the same individual to J vignette questions on the given

concept. For concreteness, we think of Y_0 as the assessment of own health on some domain and of $Y_j, j = 1, \dots, J$, as the assessment of health on the same domain in the j th vignette. We assume that the elements of the $(J + 1)$ -vector of observed responses $Y = (Y_0, Y_1, \dots, Y_J)$ are all categorical and take values $r = 0, 1, \dots, R$.

Each observed categorical response Y_j is assumed to depend on an underlying continuous latent variable Y_j^* through the observation rule

$$Y_j = \sum_{r=0}^R r \mathbf{1}(\xi_{j,r-1} < Y_j^* \leq \xi_{jr}), \quad j = 0, 1, \dots, J,$$

where $\mathbf{1}(\cdot)$ is the indicator function, and the ξ_{jr} are $R + 2$ individual-specific thresholds or cut-off points satisfying $\xi_{j,r-1} < \xi_{jr}$, with $\xi_{j,-1} = -\infty$ and $\xi_{jR} = \infty$. We refer to Greene and Hensher (2010) for a history and an extensive review of this type of models.

The statistical problem is how to use the sample information to learn about the conditional distribution of Y_0^* given a vector of observable regressors. The vignette information is not of direct interest but is used instrumentally to control for the fact that the cut-offs ξ_{jr} may vary across individuals depending on observable regressors and, possibly, unobservable individual effects.

2.1. Model specification

We assume that the continuous latent variables Y_j^* obey linear models of the form

$$Y_j^* = \alpha_j + \beta_j^T X_j + \sigma_j U_j, \quad j = 0, 1, \dots, J, \tag{1}$$

where X_j is a vector of observable exogenous regressors, which are possibly specific to the j th latent variable, α_j, β_j and σ_j are unknown parameters and U_j is an unobservable random error distributed independently of X_j with mean 0 and distribution function F . We could easily generalize this model by representing Y_j^* as additively separable in X_j and U_j , as in Cunha *et al.* (2007), i.e. by assuming that $Y_j^* = \varphi_j(X_j) + \sigma_j U_j$, where φ_j is an unknown function.

To account for observed heterogeneity in response scales, we let the thresholds depend on a vector W_j of observable exogenous regressors, which are possibly specific to the j th latent variable, i.e.

$$\xi_{jr} = \begin{cases} -\infty, & \text{if } r = -1, \\ \kappa_{jr}(W_j), & \text{if } r = 0, 1, \dots, R - 1, \\ \infty, & \text{if } r = R, \end{cases}$$

for $j = 0, 1, \dots, J$, where the κ_{jr} are unknown functions. To guarantee monotonicity of the thresholds, i.e. $\xi_{j,r-1} < \xi_{jr}$ for all r , the functions κ_{jr} must be monotonically increasing. Unobserved heterogeneity may easily be accommodated by including in W_j an unobserved individual effect, as in Rossi *et al.* (2001). This offers a simple way of allowing for correlation between self-assessment and vignette responses conditional on the observed regressors.

A parametric specification of the κ_{jr} -functions is the so-called compound hierarchical ordered response model of King *et al.* (2004), where

$$\kappa_{jr}(W_j) = \begin{cases} \gamma_{j0} + \delta_{j0}^T W_j, & \text{if } r = 0, \\ \kappa_{j,r-1} + \exp(\gamma_{jr} + \delta_{jr}^T W_j), & \text{if } r = 1, \dots, R - 1. \end{cases} \tag{2}$$

This specification guarantees monotonicity of the thresholds, i.e. $\xi_{j0} < \dots < \xi_{j,R-1}$. In addition, the non-linearities in model (2) provide weak (through functional form) identification of the model when W_j includes the same variables as X_j . An alternative parametric specification, which was originally proposed by Terza (1985), is

$$\kappa_{jr}(W_j) = \gamma_{jr} + \delta_{jr}^T W_j, \quad r = 0, 1, \dots, R - 1. \tag{3}$$

This specification does not guarantee monotonicity of the thresholds but is computationally simpler than model (2) and has the advantage of making the identification issues more transparent.

To avoid identification via functional form restrictions, we adopt the linear model (3) for the cut-offs and consider the extreme but very relevant case of no exclusion restrictions, where $X_j = W_j = X$ for all j , with X containing k exogenous regressors. Pudney and Shields (2000) also specified the thresholds as linear functions of observed regressors but achieved identification through exclusion restrictions, by excluding some of the variables in the threshold equations from those in the latent linear model (1). Since model (3) puts no constraints on the threshold parameters, we cannot ensure monotonicity of the thresholds. As a result, although the probabilities sum to 1 by construction, there is no guarantee that they are positive.

Under this model specification, the likelihood contribution of the self-assessment component is

$$\mathcal{L}_1(\theta_1; X, Y_0) \propto \prod_{r=0}^R \left\{ F\left(\frac{\xi_{0r} - \alpha_0 - \beta_0^T X}{\sigma_0}\right) - F\left(\frac{\xi_{0,r-1} - \alpha_0 - \beta_0^T X}{\sigma_0}\right) \right\}^{Y_{0r}},$$

where $Y_{0r} = \mathbf{1}(Y_0 = r)$ and the vector θ_1 consists of the parameters in $\alpha_0, \beta_0, \sigma_0, \gamma_0 = (\gamma_{00}, \dots, \gamma_{0,R-1})$ and $\delta_0 = (\delta_{00}, \dots, \delta_{0,R-1})$. The total number of parameters in θ_1 is equal to $(k + 1) \times (R + 1) + 1$. The likelihood contribution of the vignette component is

$$\mathcal{L}_2(\theta_2; X, Y_1, \dots, Y_J) \propto \prod_{j=1}^J \prod_{r=0}^R \left\{ F\left(\frac{\xi_{jr} - \alpha_j - \beta_j^T X}{\sigma_j}\right) - F\left(\frac{\xi_{j,r-1} - \alpha_j - \beta_j^T X}{\sigma_j}\right) \right\}^{Y_{jr}},$$

where $Y_{jr} = \mathbf{1}(Y_j = r)$ and the vector θ_2 consists of the parameters in all the $\alpha_j, \beta_j, \sigma_j, \gamma_j = (\gamma_{j0}, \dots, \gamma_{j,R-1})$ and $\delta_j = (\delta_{j0}, \dots, \delta_{j,R-1})$. The total number of parameters in θ_2 is equal to $J\{(k + 1)(R + 1) + 1\}$. The full likelihood for a single observation is

$$\begin{aligned} \mathcal{L}(\theta; X, Y) &= \mathcal{L}_1(\theta_1; X, Y_0) \mathcal{L}_2(\theta_2; X, Y_1, \dots, Y_J) \\ &\propto \prod_{j=0}^J \prod_{r=0}^R \left\{ F\left(\frac{\xi_{jr} - \alpha_j - \beta_j^T X}{\sigma_j}\right) - F\left(\frac{\xi_{j,r-1} - \alpha_j - \beta_j^T X}{\sigma_j}\right) \right\}^{Y_{jr}}, \end{aligned} \tag{4}$$

where $\theta = (\theta_1, \theta_2) = \{(\alpha_j, \beta_j, \sigma_j, \gamma_j, \delta_j), j = 0, \dots, J\}$ and we write (θ_1, θ_2) as a shorthand for $(\theta_1^T, \theta_2^T)^T$. The total number of parameters in θ is equal to $\{(k + 1)(R + 1) + 1\}(J + 1)$.

2.2. Identification

Identification of the model parameters requires location and scale restrictions, plus restrictions linking the self-assessment and the vignette contributions to the likelihood. After substituting the model for the cut-offs (3) into expression (4), the full likelihood for a single observation becomes

$$\mathcal{L}(\theta; X, Y) \propto \prod_{j=0}^J \prod_{r=0}^R \left[F\left\{ \frac{(\gamma_{jr} - \alpha_j) + (\delta_{jr} - \beta_j)^T X}{\sigma_j} \right\} - F\left\{ \frac{(\gamma_{j,r-1} - \alpha_j) + (\delta_{j,r-1} - \beta_j)^T X}{\sigma_j} \right\} \right]^{Y_{jr}}.$$

In the absence of prior restrictions, the parameters in θ are clearly not separately identifiable. The identifiable parameters are the following functions of the parameters in θ :

$$\gamma_{jr}^* = (\gamma_{jr} - \alpha_j) / \sigma_j,$$

$$\delta_{jr}^* = (\delta_{jr} - \beta_j) / \sigma_j,$$

with $r = 0, 1, \dots, R - 1$ and $j = 0, 1, \dots, J$. We shall refer to these parameters as the reduced form parameters. The reduced form of the model corresponds to a set of $J + 1$ ordered response models with outcome-specific parameters, which is a model that was first proposed by Pudney and Shields (2000) and referred to as the generalized ordered response model. Because the total number of parameters in the reduced form is equal to $R(k + 1)(J + 1)$, the number of restrictions that are needed to identify the parameters in θ from the identifiable reduced form parameters exactly is equal to

$$\{(k + 1)(R + 1) + 1\}(J + 1) - R(k + 1)(J + 1) = (k + 2)(J + 1).$$

Standard location and scale restrictions, namely the $2(J + 1)$ restrictions $\gamma_{j0} = 0$ and $\sigma_j = 1$, $j = 0, \dots, J$, are not enough to identify the parameters in θ , so $k(J + 1)$ additional restrictions are needed.

In the absence of vignette information ($J = 0$), the $(k + 1)(R + 1) + 1$ parameters of the model for the self-assessment cannot be obtained from the $R(k + 1)$ identifiable parameters of the reduced form because we have only two normalization restrictions ($\gamma_{00} = 0$ and $\sigma_0 = 1$). In this case, k additional restrictions would be needed to identify the model exactly. This means that we cannot separately identify the coefficients β_0 on the exogenous regressors in the latent regression for Y_0^* model from the coefficients δ_{0r} in the thresholds.

One way of achieving exact identification of the model is to exclude exogenous regressors from one threshold (Terza, 1985). This gives the k additional restrictions that are needed. In this case, however, only deviations from the cut-off from which the regressors are arbitrarily excluded can be identified. Alternatively, a standard practice in ordered response models is to assume homogeneous thresholds, i.e. $\delta_{0r} = 0$, $r = 0, 1, \dots, R - 1$, which corresponds to a set of Rk restrictions. Because only k restrictions would be needed to identify the model, when there are more than two response categories ($R > 1$) we have $(R - 1)k$ overidentifying restrictions that allow us to test the assumption of homogeneous thresholds. This test corresponds to the Wald test that was proposed by Brant (1990) for testing the proportional odds restriction in the ordered logistic regression.

If vignette information is available ($J > 0$), King *et al.* (2004) proposed to identify the model by linking the self-assessment and the vignettes through the following assumptions.

Assumption 1 (response consistency). $\gamma_{jr} - \gamma_{0r} = \delta_{jr} - \delta_{0r} = 0$, $r = 0, 1, \dots, R - 1$, $j = 1, \dots, J$.

Assumption 2 (vignette equivalence). $\beta_j = 0$, $j = 1, \dots, J$.

The first assumption is that each individual uses the response categories for a particular survey question in the same way when providing self-assessment and when assessing each of the hypothetical situations in the vignettes. The second assumption is that the level of the variable that is represented in each vignette is perceived by all respondents in the same way and on the same unidimensional scale, apart from random measurement error. Imposing assumptions 1 and 2 provides $\{R(k + 1) + k\}J$ restrictions. Because in this case self-assessment and vignettes are linked together, location and scale can be fixed by setting the constant term of the first (common) threshold $\gamma_{00} = 0$ and the variance of the self-assessment $\sigma_0 = 1$. Alternatively, location and scale can be fixed by setting the constant terms of the extreme vignettes $\alpha_1 = 0$ and $\alpha_J = 1$ (King *et al.*, 2009). Imposing assumptions 1 and 2, together with location and scale normalization of the self-assessment ($\gamma_{00} = 0$ and $\sigma_0 = 1$), gives a total of $\{R(k + 1) + k\}J + 2$ restrictions.

To illustrate, in the special case of three response categories ($R = 2$) and one exogenous

regressor ($k = 1$), the model contains $7(J + 1)$ parameters, namely $\{(\alpha_j, \beta_j, \gamma_{j0}, \delta_{j0}, \gamma_{j1}, \delta_{j1}, \sigma_j), j = 0, 1, \dots, J\}$. The reduced form parameters are only $4(J + 1)$, namely

$$\gamma_{j0}^* = (\gamma_{j0} - \alpha_j) / \sigma_j,$$

$$\delta_{j0}^* = (\delta_{j0} - \beta_j) / \sigma_j,$$

$$\gamma_{j1}^* = (\gamma_{j1} - \alpha_j) / \sigma_j,$$

$$\delta_{j1}^* = (\delta_{j1} - \beta_j) / \sigma_j,$$

with $j = 0, 1, \dots, J$. In this case, $3(J + 1)$ restrictions are needed to identify the model exactly.

Without vignettes ($J = 0$), the seven parameters in the model ($\alpha_0, \beta_0, \sigma_0, \gamma_{00}, \delta_{00}, \gamma_{01}, \delta_{01}$) cannot be obtained from the four reduced form parameters ($\gamma_{00}^*, \delta_{00}^*, \gamma_{01}^*, \delta_{01}^*$) because we have only two normalization restrictions ($\gamma_{00} = 0$ and $\sigma_0 = 1$). The model is exactly identified under the additional assumption that $\delta_{00} = 0$. Nonetheless, in this case only deviations from δ_{00} can be identified. Another possibility to achieve identification is to assume homogeneous thresholds ($\delta_{00} = 0$ and $\delta_{01} = 0$). In this case, there is one overidentifying restriction that would allow testing the homogeneous thresholds hypothesis.

With vignettes ($J > 0$), the assumption of response consistency gives $4J$ restrictions

$$\gamma_{j0} - \gamma_{00} = \gamma_{j1} - \gamma_{01} = \delta_{j0} - \delta_{00} = \delta_{j1} - \delta_{01} = 0, \quad j = 1, \dots, J,$$

whereas the assumption of vignette equivalence gives J restrictions

$$\beta_j = 0, \quad j = 1, \dots, J.$$

Because these two sets of restrictions, together with location and scale normalization ($\gamma_{00} = 0$ and $\sigma_0 = 1$), provide a total of $5J + 2$ restrictions, we have a total of $2J - 1$ overidentifying restrictions.

For example, with only one vignette ($J = 1$) we have 14 model parameters ($\alpha_0, \beta_0, \sigma_0, \gamma_{00}, \delta_{00}, \gamma_{01}, \delta_{01}, \alpha_1, \beta_1, \sigma_1, \gamma_{10}, \delta_{10}, \gamma_{11}, \delta_{11}$) and eight reduced form parameters ($\gamma_{00}^*, \delta_{00}^*, \gamma_{01}^*, \delta_{01}^*, \gamma_{10}^*, \delta_{10}^*, \gamma_{11}^*, \delta_{11}^*$). Under the two normalization restrictions ($\gamma_{00} = 0$ and $\sigma_0 = 1$) and the five restrictions that are implied by assumptions 1 and 2 ($\gamma_{10} = \gamma_{00}, \gamma_{11} = \gamma_{01}, \delta_{10} = \delta_{00}, \delta_{11} = \delta_{01}$ and $\beta_1 = 0$) the model is overidentified (it has one overidentifying restriction). With two vignettes ($J = 2$) we have 21 model parameters and 12 reduced form parameters. In this case, with two normalization restrictions and 10 restrictions implied by assumptions 1 and 2, we have three overidentifying restrictions. Finally, with three vignettes ($J = 3$) we have 28 model parameters and 16 reduced form parameters. In this case, with two normalization restrictions and 15 restrictions implied by assumptions 1 and 2, we have five overidentifying restrictions.

2.3. Inference

With vignettes ($J \geq 1$) and more than two response categories ($R \geq 2$), overidentification of the restricted model that imposes assumptions 1 and 2 and the location and scale normalizations provides the basis for testing the key assumptions 1 and 2.

One way of approaching the problem of testing is to use a minimum distance approach. Let θ be the vector of $s = \{(k + 1)(R + 1) + 1\}(J + 1)$ model parameters and let π be the vector of $q = R(k + 1)(J + 1)$ reduced form parameters. Also let ψ be the subvector of θ containing the ‘free’ parameters, namely those which are not subject to the restrictions that are implied by assumptions 1 and 2 and the location and scale normalizations. Since the number of these restrictions is equal to $\{R(k + 1) + k\}J + 2$, the number of free parameters in ψ is equal to

$p = k + R(k + 1) + 2J$, so the number of overidentifying restrictions is equal to

$$q - p = R(k + 1)(J + 1) - \{k + R(k + 1) + 2J\} = k(JR - 1) + J(R - 2).$$

When there are more than two response categories ($R \geq 2$) and at least one vignette ($J \geq 1$), we have that $q - p \geq 1$ (assuming that $k \geq 1$). In the binary response case ($R = 1$), we still have overidentifying restrictions if either $J = 2$ and $k \geq 3$, or $J \geq 3$ and $k \geq 2$.

Let π_0 and ψ_0 be the values of π and ψ in the population. Because ψ_0 includes the scale parameters σ_j , for $j = 1, \dots, J$, the relationship between π_0 and ψ_0 is non-linear. We write this relationship as

$$\pi_0 = g(\psi_0),$$

where $g: \mathfrak{R}^p \rightarrow \mathfrak{R}^q$ is a differentiable function with Jacobian matrix G . For (local) identifiability, we need $G(\psi)$ to be of full rank in an open neighbourhood of ψ_0 . Appendix A presents the structure of g and G .

Given a sample of size n from the joint distribution of (X, Y) , let $\hat{\pi}_n$ denote the estimator of π_0 that is obtained by fitting $J + 1$ generalized ordered response models, one for each categorical variable in Y . This estimator is very easy to compute and is \sqrt{n} consistent and asymptotically normal under general conditions. Given $\hat{\pi}_n$, the minimum distance method suggests estimating the vector ψ_0 of free parameters by picking the element in the parameter space Ψ such that the difference $\hat{\pi}_n - g(\psi)$ is the smallest possible. The resulting estimator of ψ_0 is consistent and asymptotically normal under general conditions (Ferguson, 1996).

An asymptotically optimal minimum distance estimator of ψ_0 is the solution $\hat{\psi}_n$ to the problem

$$\min_{\psi \in \Psi} Q_n(\psi) = (\hat{\pi}_n - g(\psi))^T \hat{V}_n^{-1} (\hat{\pi}_n - g(\psi)), \tag{5}$$

where the $q \times q$ matrix \hat{V}_n is a positive definite estimate of the asymptotic variance of $\hat{\pi}_n$. Under general conditions,

$$(\hat{\psi}_n - \psi_0) \sqrt{n} \Rightarrow \mathcal{N}\{0, (G_0 V_0^{-1} G_0^T)^{-1}\}$$

as $n \rightarrow \infty$, where $G_0 = G(\psi_0)$ denotes the $p \times q$ Jacobian matrix of g evaluated at ψ_0 and V_0 denotes the asymptotic variance of $\hat{\pi}_n$.

Computation of $\hat{\psi}_n$ is straightforward by using an iterative procedure. Starting from an initial estimate $\hat{\psi}^{(0)}$, the updated estimate at the $(h + 1)$ th iteration is given by

$$\hat{\psi}^{(h+1)} = (\hat{G}_h \hat{V}_n^{-1} \hat{G}_h^T)^{-1} \hat{G}_h \hat{V}_n^{-1} (\hat{\pi}_n - \hat{g}_h + \hat{G}_h^T \hat{\psi}^{(h)}), \quad h = 0, 1, \dots,$$

where $\hat{G}_h = G(\hat{\psi}^{(h)})$ and $\hat{g}_h = g(\hat{\psi}^{(h)})$. This corresponds to a generalized least squares regression of the transformed reduced form estimates $\hat{\pi}_n - \hat{g}_h + \hat{G}_h^T \hat{\psi}^{(h)}$ on the columns of \hat{G}_h with weighting matrix \hat{V}_n^{-1} .

When $J \geq 1$, the model that imposes conditions 1 and 2 is overidentified so, under the null hypothesis that both assumptions hold,

$$n Q_n(\hat{\psi}) \Rightarrow \chi_{q-p}^2$$

as $n \rightarrow \infty$, where $q - p = k(JR - 1) + J(R - 2)$ is the number of overidentifying restrictions. This result provides the basis for asymptotic tests that reject the key assumptions 1 and 2 for large values of the statistic $n Q_n(\hat{\psi}_n)$.

A test of this type offers several advantages. First, it can be performed with any data set containing vignette questions (one vignette is enough) on a given concept of interest and does not

require additional information like objective measures. Second, it does not require embedding the restricted model that imposes response consistency and vignette equivalence on a larger encompassing model. Third, it requires only a consistent and asymptotically normal estimator of the reduced form parameters. This is an advantage, both computationally and because the test can easily be extended to models with sample selection and to semiparametric settings where strong distributional assumptions are relaxed. Fourth, because we exploit the mapping g between the free parameters and the reduced form parameters, imposing additional restrictions is particularly simple and transparent. A potential disadvantage of our test is that it may reject the overidentifying restrictions for reasons other than failure of response consistency and vignette equivalence, e.g. because of failure of linear index restrictions or because relevant variables have been omitted from the model.

2.4. Power of the test

There are a few special cases in which the test proposed lacks power. The first case is when

$$\gamma_{jr} - \gamma_{0r} = 0$$

and

$$\delta_{jr} - \delta_{0r} - \beta_j = \delta_{ls} - \delta_{0s} - \beta_s,$$

for all vignettes j and l and all thresholds r and s . This is the unlikely case when

- (a) there is no violation of assumption 1 due to differences in the intercepts and
- (b) the violations of assumptions 1 and 2 due to the differences in the slopes are exactly the same for all thresholds and all vignettes, so they all cancel out.

For example, with three response categories ($R = 2$), one exogenous regressor ($k = 1$) and one vignette ($J = 1$), the vector of model parameters is $\theta = (\alpha_0, \beta_0, \sigma_0, \gamma_{00}, \delta_{00}, \gamma_{01}, \delta_{01}, \alpha_1, \beta_1, \sigma_1, \gamma_{10}, \delta_{10}, \gamma_{11}, \delta_{11})$ whereas the vector of reduced form parameters is $\pi = (\gamma_{00}^*, \delta_{00}^*, \gamma_{01}^*, \delta_{01}^*, \gamma_{10}^*, \delta_{10}^*, \gamma_{11}^*, \delta_{11}^*)$. In this case, if

$$\delta_{10} - \delta_{00} - \beta_1 = \delta_{11} - \delta_{01} - \beta_1 = \Delta \neq 0,$$

then the vector $\tilde{\psi} = (\alpha_0, \tilde{\beta}_0, \tilde{\delta}_{00}, \gamma_{01}, \tilde{\delta}_{01}, \alpha_1, \sigma_1)$, with $\tilde{\beta}_0 = \beta_0 + \Delta$, $\tilde{\delta}_{00} = \delta_{00} + \Delta$ and $\tilde{\delta}_{01} = \delta_{01} + \Delta$, also solves the minimization problem (5) and satisfies restrictions 1 and 2.

The second case is when

$$\gamma_{jr} - \gamma_{0r} = \gamma_{js} - \gamma_{0s} \neq 0,$$

for any vignette j and all thresholds r and s . This is the unlikely case when the violations of assumption 1 due to differences in the intercepts are exactly the same for all thresholds, so they all cancel out. In this case the violation of response consistency affects only the intercepts α_j in the vignette equations but does not affect the parameters of interest α_0 and β_0 .

Consider again the example with three response categories ($R = 2$), one exogenous regressor ($k = 1$) and one vignette ($J = 1$). In this case, if

$$\gamma_{10} - \gamma_{00} = \gamma_{11} - \gamma_{01} = \Delta \neq 0,$$

then the vector $\tilde{\psi} = (\alpha_0, \beta_0, \delta_{00}, \gamma_{01}, \delta_{01}, \tilde{\alpha}_1, \sigma_1)$, where $\tilde{\alpha}_1 = \alpha_1 - \Delta$, is also a solution to the minimization problem (5) and satisfies restrictions 1 and 2. Note that in this case the parameters of interest α_0 and β_0 are not affected.

3. Monte Carlo results

In this section, we investigate the finite sample performance of our test of the overidentifying restrictions that are implied by assumptions 1 (response consistency) and 2 (vignette equivalence) through a Monte Carlo study. Our set-up is as follows.

Step 1: we set the number of thresholds or cut-offs to $R = 2$.

Step 2: we set the number of exogenous regressors to $k = 1, 2$.

Step 3: we set the number of vignettes to $J = 1, 2$.

Step 4: we set the sample size to $n = 250, 500, 1000$.

Step 5: for all j , we draw the errors U_j from a standard normal distribution.

Step 6: the first regressor X_1 is drawn from a $U(0, 1)$ distribution, whereas the second regressor X_2 is a 0–1 indicator equal to 1 with probability 0.50. Considering the case of a binary regressor is useful because researchers are often interested in comparing subjective assessments across groups.

Step 7: the null hypothesis H_0 corresponds to the case when assumption 1 and 2 both hold. As for the alternatives, we consider three cases:

- (a) assumption 1 holds but assumption 2 fails (hypothesis H_1),
- (b) assumption 2 holds but assumption 1 fails (hypothesis H_2) and
- (c) both assumption 1 and assumption 2 fail (hypothesis H_3).

Step 8: we choose the model parameters to have an approximately even distribution of reports in each category under the null hypothesis H_0 .

Step 9: each Monte Carlo experiment consists of 1000 runs using antithetic pseudorandom numbers.

The reduced form parameters are estimated by maximizing the log-likelihood of $J + 1$ generalized ordered probit models by using the Newton–Raphson method with analytical first and second derivatives. The routines that compute the estimates of the reduced form and the free parameters are all written in Mata, which is the matrix programming language of the statistical package Stata (version 11).

Table 1 shows the Monte Carlo rejection frequencies for tests of asymptotic 5% level. The row that is labelled H_0 reports the observed size of our test, which should be compared with its asymptotic value of 5%. Already for $n = 250$, rejection frequencies are close to nominal under the null hypothesis. Thus, our test shows no evidence of size distortion in finite samples. In contrast, the size of the test remains stable when the number of overidentifying restrictions increases from 1 to 6.

The block that is labelled H_1 reports the power of our test when response consistency holds but vignette equivalence fails. The rejection frequencies are presented for increasing values of β_1 , which is the coefficient on the first regressor in the linear index for the first vignette. As discussed in Section 2.4, our test has essentially no power in the case of only one vignette, but its power increases with β_1 in the case of two vignettes. The block that is labelled H_2 reports the power of our test when vignette equivalence holds but response consistency fails. The first four rows present rejection frequencies for increasing values of the difference $\delta_{11} - \delta_{01}$, whereas the last four rows present rejection frequencies for increasing values of the differences $\delta_{11} - \delta_{01}$ and $\gamma_{11} - \gamma_{01}$. In this case, the power curves are always increasing except when $J = k = 1$ and the shift is only in the slope ($\delta_{11} - \delta_{01}$ is different from 0). Finally, the block that is labelled H_3 reports the power of our test when both assumptions fail. In this case, the results are qualitatively similar to the case of H_2 but now the power of our test is higher in all experiments.

With $n = 500$ and $n = 1000$, the results are qualitatively similar to the case of $n = 250$, but the power increases with the sample size in most experiments.

4. Empirical application

Women tend to report worse health more than men at all ages, although they are less likely to die than men and are less likely to be hospitalized than men at ages when pregnancy-related hospitalization is no longer an issue. As argued by Case and Deaton (2005), ‘this pattern . . . by gender is close to universal around the world’. This paradox could have various explanations, which are not necessarily mutually exclusive. One is that gender differences in self-assessment of health reflect systematic differences in the prevalence of chronic conditions, for either biological or behavioural reasons. For example, Case and Paxson (2005) showed that, in the USA, gender differences in self-rated general health are almost entirely due to the differences in the distribution of reported chronic conditions, with hardly any role for gender differences in the mapping from chronic conditions to reported poor health. Another explanation is that gender differences in self-assessment of health reflect systematic differences in the way that respondents locate themselves on subjective scales (Lindeboom and van Doorslaer, 2004). Anchoring vignettes offer one way of controlling for such differences.

4.1. Data

Our data are from release 2 of the first (2004–2005) wave of the SHARE, which is a multidisciplinary and cross-national biannual household panel survey, that is nationally representative of the population aged 50 years or older living in private households in Europe. The first wave covers about 19 500 households and about 28 500 individuals in 11 European countries (Austria, Belgium, Denmark, France, Germany, Greece, the Netherlands, Spain, Sweden and Switzerland). For a detailed description, see Börsch-Supan and Jürges (2005).

The SHARE collects detailed information on demographic and economic variables health, psychological variables and social support variables. In particular, respondents are asked to use a five-point ordered scale to rate their own health in general and to assess their health on six domains, namely pain, sleeping problems, mobility problems, concentration problems, shortness of breath and depression. In eight countries (Belgium, France, Germany, Greece, Italy, the Netherlands, Spain and Sweden), a random subsample of the respondents is also asked to answer vignette questions on the six health domains. The vignette questions are presented in a random order after the self-assessment questions. For each domain, respondents are presented three hypothetical situations, corresponding to people with low, moderate and serious health problems. They are instructed to evaluate the hypothetical people on exactly the same five-point ordered scale as is used for the self-assessments and to assume that the hypothetical people in the vignettes have their same age and background.

4.2. Descriptive statistics

We restrict attention to men and women aged 50–80 years for whom the vignette information is available and there are no missing data on any of the variables that we use. Because the fraction with missing data is small (less than 3% for self-assessment questions and less than 5% for vignette questions), we work with the subsample with complete data and ignore selection issues. This gives a sample of 3458 observations (1631 men and 1827 women) that represents about 16% of the full SHARE sample in the relevant age group. Table 2 compares the composition of our working sample with that of the full SHARE sample and the vignette sample for the age

Table 1. Monte Carlo rejection frequencies for tests of asymptotic 5% level†

| Hypothesis | Results for $k = 1, J = 1,$ $q - p = 1$ and the following values of n : | | | Results for $k = 2, J = 1,$ $q - p = 2$ and the following values of n : | | | Results for $k = 3, J = 2,$ $q - p = 3$ and the following values of n : | | | Results for $k = 2, J = 2,$ $q - p = 6$ and the following values of n : | | |
|---|---|-------|-------|---|-------|-------|---|-------|-------|---|-------|-------|
| | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 |
| H_0 | 0.057 | 0.059 | 0.050 | 0.050 | 0.042 | 0.053 | 0.056 | 0.059 | 0.043 | 0.053 | 0.052 | 0.052 |
| H_1 | 0.062 | 0.052 | 0.051 | 0.051 | 0.047 | 0.052 | 0.056 | 0.065 | 0.053 | 0.052 | 0.059 | 0.052 |
| $\beta_1 = 0.1$ | 0.070 | 0.060 | 0.063 | 0.043 | 0.051 | 0.052 | 0.055 | 0.057 | 0.079 | 0.053 | 0.071 | 0.052 |
| $\beta_1 = 0.2$ | 0.057 | 0.079 | 0.055 | 0.068 | 0.056 | 0.052 | 0.080 | 0.128 | 0.185 | 0.070 | 0.081 | 0.148 |
| $\beta_1 = 0.4$ | 0.067 | 0.062 | 0.053 | 0.051 | 0.057 | 0.057 | 0.107 | 0.207 | 0.428 | 0.076 | 0.158 | 0.328 |
| $\beta_1 = 0.6$ | 0.055 | 0.057 | 0.050 | 0.051 | 0.053 | 0.054 | 0.172 | 0.337 | 0.663 | 0.142 | 0.291 | 0.587 |
| $\beta_1 = 0.8$ | 0.068 | 0.073 | 0.070 | 0.065 | 0.046 | 0.045 | 0.254 | 0.543 | 0.903 | 0.188 | 0.463 | 0.774 |
| H_2 | 0.055 | 0.054 | 0.054 | 0.054 | 0.059 | 0.057 | 0.052 | 0.068 | 0.050 | 0.054 | 0.058 | 0.056 |
| $\delta_{11} - \delta_{01} = 0.1$ | 0.059 | 0.067 | 0.067 | 0.064 | 0.076 | 0.100 | 0.047 | 0.046 | 0.059 | 0.057 | 0.067 | 0.059 |
| $\delta_{11} - \delta_{01} = 0.2$ | 0.059 | 0.074 | 0.047 | 0.102 | 0.139 | 0.244 | 0.061 | 0.091 | 0.143 | 0.065 | 0.088 | 0.129 |
| $\delta_{11} - \delta_{01} = 0.4$ | 0.074 | 0.067 | 0.080 | 0.124 | 0.257 | 0.425 | 0.103 | 0.137 | 0.263 | 0.071 | 0.130 | 0.210 |
| $\delta_{11} - \delta_{01} = 0.6$ | 0.059 | 0.066 | 0.065 | 0.181 | 0.372 | 0.654 | 0.104 | 0.227 | 0.369 | 0.096 | 0.190 | 0.378 |
| $\delta_{11} - \delta_{01} = 0.8$ | 0.056 | 0.064 | 0.055 | 0.256 | 0.501 | 0.808 | 0.151 | 0.260 | 0.518 | 0.103 | 0.240 | 0.518 |
| $\delta_{11} - \delta_{01} = \gamma_{11} - \gamma_{01} = 0.1$ | 0.099 | 0.125 | 0.204 | 0.051 | 0.062 | 0.062 | 0.049 | 0.054 | 0.061 | 0.046 | 0.057 | 0.066 |
| $\delta_{11} - \delta_{01} = \gamma_{11} - \gamma_{01} = 0.2$ | 0.181 | 0.254 | 0.407 | 0.067 | 0.096 | 0.129 | 0.054 | 0.062 | 0.088 | 0.054 | 0.069 | 0.129 |
| $\delta_{11} - \delta_{01} = \gamma_{11} - \gamma_{01} = 0.4$ | 0.296 | 0.495 | 0.754 | 0.108 | 0.146 | 0.317 | 0.065 | 0.078 | 0.142 | 0.091 | 0.153 | 0.341 |
| $\delta_{11} - \delta_{01} = \gamma_{11} - \gamma_{01} = 0.6$ | 0.377 | 0.600 | 0.875 | 0.155 | 0.273 | 0.458 | 0.092 | 0.140 | 0.224 | 0.135 | 0.309 | 0.622 |
| $\delta_{11} - \delta_{01} = \gamma_{11} - \gamma_{01} = 0.8$ | 0.344 | 0.629 | 0.911 | 0.148 | 0.318 | 0.632 | 0.087 | 0.186 | 0.331 | 0.187 | 0.418 | 0.805 |
| $\delta_{11} - \delta_{01} = \gamma_{11} - \gamma_{01} = 1$ | 0.281 | 0.595 | 0.905 | 0.174 | 0.360 | 0.671 | 0.088 | 0.223 | 0.438 | 0.214 | 0.559 | 0.904 |
| H_3 | 0.065 | 0.059 | 0.061 | 0.050 | 0.053 | 0.063 | 0.049 | 0.061 | 0.059 | 0.057 | 0.061 | 0.059 |
| $\beta_1 = \delta_{11} - \delta_{01} = 0.1$ | 0.074 | 0.063 | 0.064 | 0.070 | 0.073 | 0.096 | 0.059 | 0.062 | 0.104 | 0.065 | 0.074 | 0.065 |
| $\beta_1 = \delta_{11} - \delta_{01} = 0.2$ | 0.068 | 0.078 | 0.074 | 0.110 | 0.144 | 0.260 | 0.095 | 0.170 | 0.270 | 0.089 | 0.123 | 0.232 |
| $\beta_1 = \delta_{11} - \delta_{01} = 0.4$ | 0.081 | 0.078 | 0.075 | 0.145 | 0.261 | 0.464 | 0.149 | 0.298 | 0.562 | 0.119 | 0.239 | 0.490 |
| $\beta_1 = \delta_{11} - \delta_{01} = 0.6$ | 0.063 | 0.072 | 0.092 | 0.189 | 0.381 | 0.705 | 0.243 | 0.458 | 0.786 | 0.202 | 0.423 | 0.780 |
| $\beta_1 = \delta_{11} - \delta_{01} = 0.8$ | 0.057 | 0.084 | 0.098 | 0.291 | 0.555 | 0.854 | 0.312 | 0.628 | 0.941 | 0.269 | 0.590 | 0.899 |

(continued)

Table 1 (continued)

| Hypothesis | Results for $k = 1, J = 1,$ $q - p = 1$ and the following values of n : | | | Results for $k = 2, J = 1,$ $q - p = 2$ and the following values of n : | | | Results for $k = 1, J = 2,$ $q - p = 3$ and the following values of n : | | | Results for $k = 2, J = 2,$ $q - p = 6$ and the following values of n : | | |
|--|---|-------|-------|---|-------|-------|---|-------|-------|---|-------|-------|
| | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 | 250 | 500 | 1000 |
| H_3 | | | | | | | | | | | | |
| $\beta_1 = \delta_{11} - \delta_{01} = \gamma_{11} - \gamma_{01}$ $= 0.1$ | 0.094 | 0.125 | 0.186 | 0.058 | 0.053 | 0.071 | 0.052 | 0.062 | 0.078 | 0.048 | 0.059 | 0.074 |
| $\beta_1 = \delta_{11} - \delta_{01} = \gamma_{11} - \gamma_{01}$ $= 0.2$ | 0.168 | 0.262 | 0.413 | 0.066 | 0.103 | 0.137 | 0.069 | 0.099 | 0.147 | 0.070 | 0.103 | 0.195 |
| $\beta_1 = \delta_{11} - \delta_{01} = \gamma_{11} - \gamma_{01}$ $= 0.4$ | 0.307 | 0.525 | 0.809 | 0.105 | 0.148 | 0.324 | 0.125 | 0.218 | 0.444 | 0.139 | 0.275 | 0.548 |
| $\beta_1 = \delta_{11} - \delta_{01} = \gamma_{11} - \gamma_{01}$ $= 0.6$ | 0.419 | 0.641 | 0.915 | 0.175 | 0.286 | 0.524 | 0.184 | 0.409 | 0.719 | 0.254 | 0.572 | 0.879 |
| $\beta_1 = \delta_{11} - \delta_{01} = \gamma_{11} - \gamma_{01}$ $= 0.8$ | 0.467 | 0.758 | 0.963 | 0.192 | 0.367 | 0.681 | 0.251 | 0.591 | 0.880 | 0.382 | 0.724 | 0.985 |
| $\beta_1 = \delta_{11} - \delta_{01} = \gamma_{11} - \gamma_{01}$ $= 1$ | 0.485 | 0.801 | 0.972 | 0.244 | 0.459 | 0.760 | 0.324 | 0.654 | 0.958 | 0.474 | 0.860 | 0.996 |

†The number of thresholds is $R = 2$ and the number of runs is 1000 per experiment.

Table 2. SHARE sample size by country and gender (people aged 50–80 years)†

| <i>Country</i> | <i>Full sample</i> | | <i>Vignette sample</i> | | <i>Working sample</i> | |
|----------------|--------------------|--------------|------------------------|--------------|-----------------------|--------------|
| | <i>Men</i> | <i>Women</i> | <i>Men</i> | <i>Women</i> | <i>Men</i> | <i>Women</i> |
| Belgium | 1602 | 1791 | 234 | 291 | 201 | 244 |
| France | 1270 | 1484 | 352 | 451 | 301 | 368 |
| Germany | 1323 | 1448 | 211 | 264 | 168 | 210 |
| Greece | 1154 | 1277 | 317 | 298 | 285 | 254 |
| Italy | 1077 | 1295 | 189 | 229 | 149 | 184 |
| Netherlands | 1272 | 1402 | 242 | 257 | 213 | 216 |
| Spain | 900 | 1194 | 185 | 238 | 173 | 202 |
| Sweden | 1285 | 1439 | 186 | 203 | 141 | 149 |
| Total | 9883 | 11330 | 1916 | 2231 | 1631 | 1827 |

†The full sample includes all 50–80-year-old respondents, the vignette sample includes all 50–80-year-old respondents who answer the vignette questions and the working sample includes the respondents in the vignette sample with no missing data on any of the variables used in our analysis.

Table 3. Correlation between self-rated general health and self-assessments on the various health domains†

| | <i>Self-rated health</i> | <i>Pain</i> | <i>Sleeping problems</i> | <i>Mobility problems</i> | <i>Concentration problems</i> | <i>Shortness of breath</i> | <i>Depression</i> | <i>Ordered probit coefficient</i> |
|------------------------|--------------------------|-------------|--------------------------|--------------------------|-------------------------------|----------------------------|-------------------|-----------------------------------|
| Self-rated health | 1.000 | | | | | | | |
| Pain | 0.443 | 1.000 | | | | | | 0.513 (0.040) |
| Sleeping problems | 0.292 | 0.415 | 1.000 | | | | | 0.109 (0.033) |
| Mobility problems | 0.446 | 0.537 | 0.371 | 1.000 | | | | 0.487 (0.039) |
| Concentration problems | 0.245 | 0.340 | 0.304 | 0.339 | 1.000 | | | 0.046 (0.037) |
| Shortness of breath | 0.281 | 0.306 | 0.241 | 0.383 | 0.298 | 1.000 | | 0.178 (0.038) |
| Depression | 0.291 | 0.378 | 0.399 | 0.353 | 0.391 | 0.329 | 1.000 | 0.112 (0.035) |

†The last column shows estimated coefficients from an ordered probit model for self-rated general health on self-assessments of health on the six domains.

group 50–80 years. Country differences in the importance of the vignette sample are mainly due to differences in sampling design and availability of funding. We account for such differences by using the survey weights that were specifically provided for the vignette sample.

Table 3 shows the correlation between self-rated general health and self-assessments on the six health domains. Self-reported problems on these domains are all positively correlated with general health and with each other. The correlation with general health is highest for pain and mobility problems (0.44), whereas the correlation between domains is highest for pain and mobility problems (0.53). The last column of Table 3 shows the estimated coefficients from an ordered probit model for general health on self-assessments of health in the six domains. Because the estimated coefficient is highest for pain (0.513), we focus on this particular health domain. (Results for the other five health domains are available from the authors on request.)

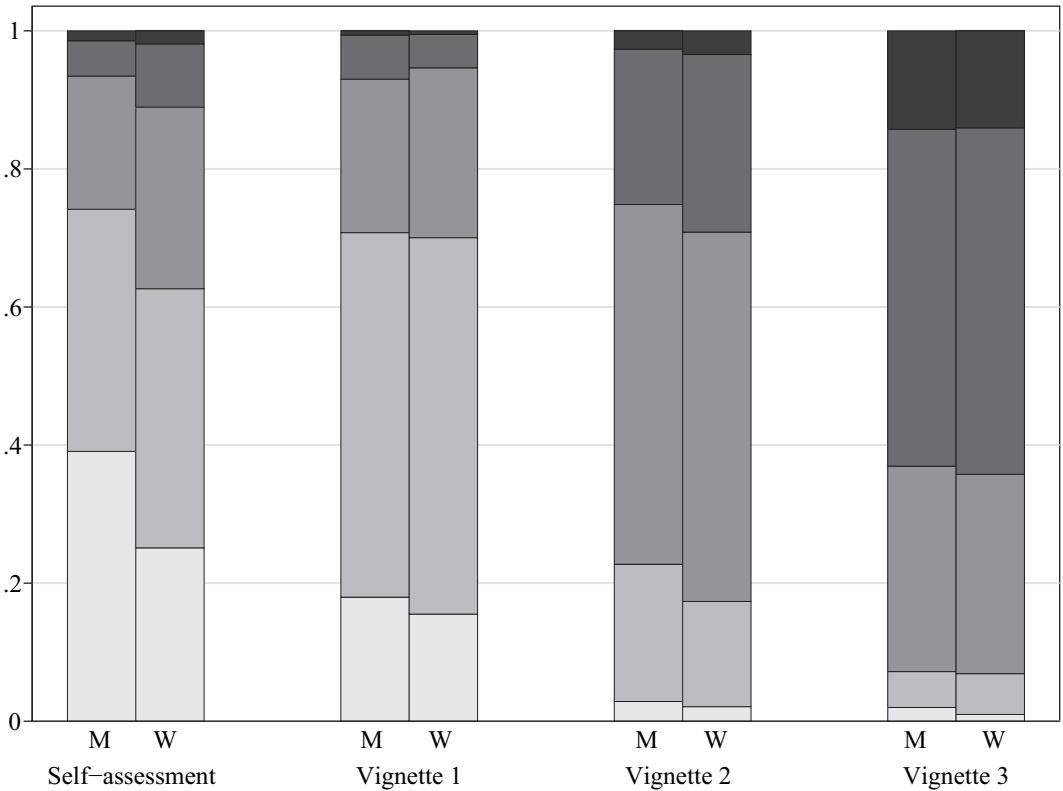


Fig. 1. Histograms of self-assessments and answers to the vignette question on pain by gender (M, men; W, women): □ none; □, mild; □, moderate; □, severe; □, extreme

Appendix B reports the various vignettes for pain, where the labels ‘vignette 1’, ‘vignette 2’ and ‘vignette 3’ do not represent the order in which the three vignettes are presented (which is random) but instead refer to the severity of the hypothetical situation (low, moderate and serious).

Fig. 1 shows the histograms of the self-assessment of pain and the answers to the vignette questions by gender. Women are more likely to report severe or extreme pain than men. The distribution of the answers to the vignette questions confirms that on average respondents tend to rank the three vignettes from least to most severe pain.

4.3. Results

We estimate a fully parametric version of our model by assuming normality of the latent errors in model (1). The reduced form of our model corresponds to a set of $J + 1 = 4$ ordered probit models with outcome-specific parameters. To avoid increasing too much the number of over-identifying restrictions, we merge the response category ‘mild’ with ‘moderate’, and ‘severe’ with ‘extreme’. This gives $R = 2$ thresholds.

Because no credible exclusion restriction is available, we allow W_j to contain exactly the same regressors as X_j for all j . In our baseline specification, the regressors include a females indicator, age, an indicator for college education completed, the logarithm of *per capita* household income, an indicator for reporting at least one diagnosed chronic condition and hand grip strength. The last is to consider an objective measure of health and is known to be a good predictor of future

Table 4. Minimum distance estimates of the coefficients of the ordered response model with heterogeneous thresholds for pain under the assumptions of response consistency and vignette equivalence

| | <i>Self-assessment</i> | <i>Threshold 1</i> | <i>Threshold 2</i> |
|--------------------------|------------------------|--------------------|--------------------|
| Any condition | 0.527† | -0.085† | -0.033 |
| Grip strength 34.9 | -0.018† | -0.004‡ | -0.001 |
| Age 55 years | 0.001 | 0.000 | 0.003‡ |
| Post-secondary education | -0.177† | -0.145† | -0.047 |
| Log-household-income | -0.098† | -0.063† | -0.067† |
| Female | 0.006 | -0.181† | -0.013 |
| Constant | -0.084 | 0.000 | 1.776† |
| | <i>Vignette 1</i> | <i>Vignette 2</i> | <i>Vignette 3</i> |
| Constant | 0.560† | 1.337† | 2.256† |
| ln(σ) | -0.283† | -0.286† | 0.057 |

†Significant at 1%.
‡Significant at 5%.

Table 5. Tests of response consistency and vignette equivalence: all respondents ($n = 3458$)

| | <i>k</i> | <i>R</i> | <i>J</i> | <i>q</i> | <i>p</i> | <i>q - p</i> | χ^2 | <i>p-value</i> |
|------------------|----------|----------|----------|----------|----------|--------------|----------|----------------|
| All vignettes | 6 | 2 | 3 | 56 | 26 | 30 | 67.4 | 0.000 |
| Only vignette 1 | 6 | 2 | 1 | 28 | 22 | 6 | 8.7 | 0.188 |
| Only vignette 2 | 6 | 2 | 1 | 28 | 22 | 6 | 8.6 | 0.200 |
| Only vignette 3 | 6 | 2 | 1 | 28 | 22 | 6 | 12.6 | 0.051 |
| Vignette 1 and 2 | 6 | 2 | 2 | 42 | 24 | 18 | 43.2 | 0.001 |
| Vignette 1 and 3 | 6 | 2 | 2 | 42 | 24 | 18 | 52.0 | 0.000 |
| Vignette 2 and 3 | 6 | 2 | 2 | 42 | 24 | 18 | 20.5 | 0.305 |
| All 5 categories | 6 | 4 | 3 | 112 | 40 | 72 | 468.6 | 0.000 |

medical problems (Rantanen *et al.*, 1999). It is measured here as the maximum of up to four measurements taken by the interviewer: two for each hand. The baseline specification includes $k = 6$ regressors, so the number of reduced form parameters is equal to $q = R(k + 1)(J + 1) = 56$, the number of free parameters is equal to $p = k + R(k + 1) + 2J = 26$ and the number of overidentifying restrictions is equal to $q - p = k(JR - 1) + J(R - 2) = 30$. Table 4 presents the minimum distance estimates of the model parameters under the assumptions of response consistency and vignette equivalence. Note that predicted probabilities are always positive in this specification of the model.

Table 5 presents the results of the χ^2 -test of the overidentifying restrictions that are implied by our two key assumptions. Using the full sample, three vignettes ($J = 3$ and three response categories ($R = 2$)), the overidentifying restrictions are rejected at any conventional level of significance. The remainder of Table 5 shows the results that are obtained when the test is carried out by using different subsets of the vignettes ($J = 1$ or $J = 2$) or using all five original response categories ($R = 4$). The overidentifying restrictions are not rejected at the 5% level when using only one vignette, especially when using the first or the second. They are also not rejected when using the second and the third vignette together, but not when using the first and

Table 6. Tests of response consistency and vignette equivalence: subgroups of respondents

| | <i>n</i> | <i>k</i> | <i>R</i> | <i>J</i> | <i>q</i> | <i>p</i> | <i>q - p</i> | χ^2 | <i>p-value</i> |
|--|----------|----------|----------|----------|----------|----------|--------------|----------|----------------|
| Men | 1631 | 5 | 2 | 3 | 48 | 23 | 25 | 33.1 | 0.129 |
| Women | 1827 | 5 | 2 | 3 | 48 | 23 | 25 | 40.7 | 0.025 |
| Aged 50–64 years | 2152 | 6 | 2 | 3 | 56 | 26 | 30 | 48.0 | 0.020 |
| Aged 65–80 years | 1306 | 6 | 2 | 3 | 56 | 26 | 30 | 32.4 | 0.347 |
| No conditions | 995 | 5 | 2 | 3 | 48 | 23 | 25 | 38.2 | 0.044 |
| Any condition | 2463 | 5 | 2 | 3 | 48 | 23 | 25 | 31.3 | 0.180 |
| Less than secondary education | 1834 | 5 | 2 | 3 | 48 | 23 | 25 | 49.8 | 0.002 |
| Secondary and post-secondary education | 1624 | 5 | 2 | 3 | 48 | 23 | 25 | 22.8 | 0.591 |
| Mediterranean countries | 1247 | 6 | 2 | 3 | 56 | 26 | 30 | 36.1 | 0.204 |
| Non-Mediterranean countries | 2211 | 6 | 2 | 3 | 56 | 26 | 30 | 49.2 | 0.015 |
| Mediterranean men | 607 | 5 | 2 | 3 | 48 | 23 | 25 | 28.3 | 0.294 |
| Mediterranean women | 640 | 5 | 2 | 3 | 48 | 23 | 25 | 12.1 | 0.985 |
| Non-Mediterranean men | 1024 | 5 | 2 | 3 | 48 | 23 | 25 | 21.4 | 0.671 |
| Non-Mediterranean women | 1187 | 5 | 2 | 3 | 48 | 23 | 25 | 35.1 | 0.087 |

either the second or the third. Our results are consistent with those of Voňková and Hulleigie (2011), who also found that the vignette method is sensitive to the choice of the vignette. When the test is carried out by using all five original response categories (the last row of Table 5), the overidentifying restrictions are again strongly rejected.

Table 6 shows the results that are obtained when the test is carried out separately for various subgroups of respondents. Specifically, we group respondents by gender, age group (50–64 years *versus* 65–80 years), health status (no self-reported chronic condition *versus* some conditions), educational attainments (less than secondary *versus* secondary or post secondary) and region of residence (Mediterranean *versus* non-Mediterranean country). Now the overidentifying restrictions are rejected for women, people aged 50–64 years, people reporting no chronic condition, people with less than secondary education and residents in non-Mediterranean countries, but are not rejected for men, people aged 65–80 years, people reporting some chronic conditions, more educated people and for residents in Mediterranean countries. The fact that, when splitting the sample in two subgroups of similar size, the results of the test may be quite different suggests three things. First, failure to reject is not simply due to a smaller sample size. Second, since response consistency is a within-respondent property whereas vignette equivalence is a between-respondent property, our evidence suggests that the assumption of vignette equivalence is perhaps more problematic. Third, some of our subgroups may still be too heterogeneous for vignette equivalence to hold. In fact, when we further distinguish by region and gender (the last four rows of Table 6), the overidentifying restrictions are never rejected at the 5% level and are rejected at the 10% level only for non-Mediterranean women.

5. Conclusions

Vignette questions have been introduced in several household surveys (the SHARE, the Health and Retirement Study, the English Longitudinal Study on Ageing and the World Health Survey) and are increasingly used in various fields as an instrument to anchor response scales and to allow comparisons across individuals. Reliability of this approach hinges crucially on the validity of the key assumptions of response consistency and vignette equivalence (King *et al.*, 2004). In this paper we introduce a simple joint test of these two assumptions by exploiting the fact that, as pointed out by Deaton (2011), the statistical model is overidentified under these

two assumptions. Our Monte Carlo results show that the test proposed has good size and power properties in finite samples.

Using data from the first wave of the SHARE, we apply our test to self-assessment of pain. We find that, in several cases, the overidentifying restrictions that are imposed by the assumptions of response consistency and vignette equivalence are rejected. This typically occurs when we use more than one vignette question or, as also argued by Rice *et al.* (2012), when the model specification is not sufficiently rich to account fully for individual heterogeneity. These results suggest that the assumption of vignette equivalence is perhaps more problematic, but also that care is needed with model specification because vignette equivalence may be violated because of failure to control properly for heterogeneity across respondents. In fact, when we carry out the test separately for subgroups of respondents who are distinguished by gender, age group, health status, education and region, the evidence against the overidentifying restrictions becomes weaker, especially for men and for people who are less healthy, more educated or live in Mediterranean countries.

Overall, our results confirm the importance of testing the validity of the vignette approach that is used for identifying and correcting interpersonal incomparability of answers to subjective survey questions. Our results also point to the fruitfulness of exploring new research directions. One direction is vignette design, in particular how to minimize the risk that the vignettes may be interpreted differently. Another direction is extensions to semiparametric or non-parametric settings. Relaxing distributional assumptions will also avoid the risk that the test rejects because of problems with the parametric specification assumed.

Acknowledgements

We thank Karim Abadir, Anne Case, Valentino Dardanoni, Angus Deaton, Chris Paxson, Frank Vella, the Associate Editor, three referees and seminar participants at the European Center for Advance Research in Economics and Statistics, the 2011 Italian Congress of Econometrics and Empirical Economics, the National Bureau of Economic Research, Princeton University, the University of Alicante and the University of Naples for helpful comments. Franco Peracchi also thanks the Center for Health and Wellbeing at Princeton University for generous hospitality during the autumn of 2010.

Appendix A: Structure of the function g and its Jacobian matrix

Write the vector of $p = k + R(k + 1) + 2J$ 'free' parameters in the model as $\psi = (\rho, \sigma)$, where ρ is the $(p - J)$ -subvector of ψ containing the parameters entering the function g linearly and $\sigma = (\sigma_1, \dots, \sigma_J)$ is the J -subvector of ψ containing the scale parameters entering g non-linearly. Then, the relationship between the reduced form parameters in π and the free parameters in ψ may be written

$$\pi = g(\psi) = A(\sigma)\rho,$$

where $A(\sigma)$ is a $q \times (p - J)$ matrix that does not depend on ρ . The $p \times q$ Jacobian matrix of $g(\psi)$ is then

$$G(\psi) = \frac{\partial g(\psi)}{\partial \psi} = \begin{pmatrix} \frac{\partial g(\psi)}{\partial \rho} \\ \frac{\partial g(\psi)}{\partial \sigma} \end{pmatrix} = \begin{pmatrix} A(\sigma)^T \\ \rho^T A_1(\sigma)^T \\ \vdots \\ \rho^T A_J(\sigma)^T \end{pmatrix},$$

where $A_j(\sigma) = \partial A(\sigma) / \partial \sigma_j$ is a $q \times (p - J)$ matrix.

To illustrate, in the special case of three response categories ($R = 2$), one exogenous regressor ($k = 1$) and one vignette ($J = 1$), the vector of $q = 8$ reduced form parameters is

$$\pi = (\gamma_{00}^*, \delta_{00}^*, \gamma_{01}^*, \delta_{01}^*, \gamma_{10}^*, \delta_{10}^*, \gamma_{11}^*, \delta_{11}^*).$$

Let $\psi = (\rho, \sigma)$ be the vector of $p = 7$ free parameters, where $\rho = (\alpha_0, \beta_0, \delta_{00}, \gamma_{01}, \delta_{01}, \alpha_1)$ and $\sigma = \sigma_1$. In this case, the relationship between π and ψ can be rewritten as $\pi = g(\psi) = A(\sigma_1)\rho$, where

$$A(\sigma_1) = \begin{pmatrix} -1 & 0 & 0 & 0 & 0 & 0 \\ 0 & -1 & 1 & 0 & 0 & 0 \\ -1 & 0 & 0 & 1 & 0 & 0 \\ 0 & -1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & -1/\sigma_1 \\ 0 & 0 & 1/\sigma_1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1/\sigma_1 & 0 & -1/\sigma_1 \\ 0 & 0 & 0 & 0 & 1/\sigma_1 & 0 \end{pmatrix}.$$

The 7×8 Jacobian matrix of $g(\psi)$ is then

$$G(\psi) = \begin{pmatrix} A(\sigma_1)^T \\ \rho^T A'(\sigma_1)^T \end{pmatrix},$$

where

$$A'(\sigma_1) = \begin{pmatrix} 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1/\sigma_1^2 \\ 0 & 0 & -1/\sigma_1^2 & 0 & 0 & 0 \\ 0 & 0 & 0 & -1/\sigma_1^2 & 0 & 1/\sigma_1^2 \\ 0 & 0 & 0 & 0 & -1/\sigma_1^2 & 0 \end{pmatrix}.$$

Appendix B: Vignette questions for pain

The vignette questions for pain are as follows.

- ‘1. Paul/Karen has a headache once a month that is relieved after taking a pill. During the headache he/she can carry on with his/her day-to-day affairs.’
- ‘2. Henri/Maria has pain that radiates down his/her right arm and wrist during his/her day at work. This is slightly relieved in the evenings when he/she is no longer working on his/her computer.’
- ‘3. Charles/Alice has pain in his/her knees, elbows, wrists and fingers, and the pain is present almost all the time. Although medication helps, he/she feels uncomfortable when moving around, holding and lifting things.’

References

Angelini, V., Cavapozzi, D., Corazzini, L. and Paccagnella, O. (2008) Do Danes and Italians rate life satisfaction in the same way?: using vignettes to correct for individual-specific scale biases. *Mimeo*. University of Padua, Padua.

Bago d’Uva, T., van Doorslaer, E., Lindeboom, M. and O’Donnell, O. (2008) Does reporting heterogeneity bias the measurement of health disparities? *Hlth Econ.*, **17**, 351–375.

Bago d’Uva, T., Lindeboom, M., O’Donnell, O. and van Doorslaer, E. (2011) Slipping anchor?: testing the vignettes approach to identification and correction of reporting heterogeneity. *J. Hum. Resour.*, **46**, 872–903.

Bago d’Uva, T., O’Donnell, O. and van Doorslaer, E. (2008) Differential health reporting by education level and its impact on the measurement of health inequalities among older Europeans. *Int. J. Epidem.*, **37**, 1375–1383.

Börsch-Supan, A. and Jürges, H. (2005) *The Survey of Health, Aging, and Retirement in Europe: Methodology*. Mannheim: Mannheim Research Institute for the Economics of Aging.

Brant, R. (1990) Assessing proportionality in the proportional odds model for ordinal logistic regression. *Biometrics*, **46**, 1171–1178.

Case, A. and Deaton, A. (2005) Broken down by work and sex: how our health declines. In *Analyses in the Economics of Aging* (ed. D. A. Wise). Chicago: University of Chicago Press.

- Case, A. and Paxson, C. (2005) Sex differences in morbidity and mortality. *Demography*, **42**, 189–214.
- Cunha, F., Heckman, J. J. and Navarro, S. (2007) The identification and economic content of ordered choice models with stochastic thresholds. *Technical Working Paper 340*. National Bureau of Economic Research, Cambridge.
- Datta Gupta, N., Kristensen, N. and Pozzoli, D. (2010) External validation of the use of vignettes in cross-country health studies. *Econ. Modelling*, **27**, 854–865.
- Deaton, A. (2011) Comment on ‘Work disability, work, and justification bias in Europe and the U.S.’. In *Explorations in the Economics of Aging* (ed. D. A. Wise), pp. 312–314. Chicago: University of Chicago Press.
- Ferguson, T. S. (1996) *A Course in Large Sample Theory*. London: Chapman and Hall.
- Greene, W. H. and Hensher, D. A. (2010) *Modeling Ordered Choices: a Primer*. New York: Cambridge University Press.
- Holland, P. W. and Wainer, H. (1993) *Differential Item Functioning*. Hillsdale: Erlbaum.
- Kapteyn, A., Smith, J. and van Soest, A. (2007) Vignettes and self-reports of work disability in the United States and the Netherlands. *Am. Econ. Rev.*, **97**, 461–473.
- Kapteyn, A., Smith, J., van Soest, A. and Voňková, H. (2011) Anchoring vignettes and response consistency. *Working Paper 840*. RAND Corporation, Santa Monica.
- King, G., Lau, O. and Wand, J. (2009) Anchors: software for anchoring vignette data. *J. Statist. Softw.*, to be published.
- King, G., Murray, C. J. L., Salomon, J. A. and Tandon, A. (2004) Enhancing the validity and cross-cultural comparability of measurement in survey research. *Am. Polit. Sci. Rev.*, **98**, 191–207.
- King, G. and Wand, J. (2007) Comparing incomparable survey responses: evaluating and selecting anchoring vignettes. *Polit. Anal.*, **15**, 46–66.
- Kristensen, N. and Johansson, E. (2008) New evidence on cross-country differences in job satisfaction using anchoring vignettes. *Lab. Econ.*, **15**, 96–117.
- Lindeboom, M. and van Doorslaer, E. (2004) Cut-point shift and index shift in self-reported health. *J. Hlth Econ.*, **23**, 1083–1099.
- Peracchi, F. and Rossetti, C. (2009) Gender and regional differences in self-rated health in Europe. *Working Paper 142*. Centro per l’Economia Internazionale e lo Sviluppo, Rome.
- Pudney, S. and Shields, M. (2000) Gender, race, pay and promotion in the British nursing profession: estimation of a generalized ordered probit model. *J. Appl. Econometr.*, **15**, 367–399.
- Rantanen, T., Guralnik, J. M., Foley, D., Masaki, K., Leveille, S. G., Curb, J. D. and White, L. (1999) Midlife hand grip strength as a predictor of old age disability. *J. Am. Med. Ass.*, **281**, 558–560.
- Rice, N., Robone, S. and Smith, P. C. (2012) Vignettes and health systems responsiveness in cross-country comparative analyses (with discussion). *J. R. Statist. Soc. A*, **175**, 337–369.
- Rossi, P. E., Gilula, Z. and Allenby, G. M. (2001) Overcoming scale usage heterogeneity: a Bayesian hierarchical approach. *J. Am. Statist. Ass.*, **96**, 20–31.
- Salomon, J. A., Tandon, A. and Murray, C. J. L. (2004) Comparability of self rated health: cross sectional multi-country survey using anchoring vignettes. *Br. Med. J.*, **328**, 258–260.
- Sen, A. (2002) Health: perception versus observation. *Br. Med. J.*, **324**, 860–861.
- van Soest, A., Delaney, L., Harmon, C., Kapteyn, A. and Smith, J. P. (2011) Validating the use of anchoring vignettes for the correction of response scales differences in subjective questions. *J. R. Statist. Soc. A*, **174**, 575–595.
- StataCorp (2009) *Statistical Software: Release 11*. College Station: Stata Corp.
- Terza, J. (1985) Ordered probit: a generalization. *Commun. Statist.*, **14**, 1–11.
- Voňková, H. and Hulle, P. (2011) Is the anchoring vignette method sensitive to the domain and choice of the vignette? *J. R. Statist. Soc. A*, **174**, 597–620.