



This is a repository copy of *Deriving a Preference-Based Measure for Cancer Using the EORTC QLQ-C30*.

White Rose Research Online URL for this paper:
<http://eprints.whiterose.ac.uk/110873/>

Version: Accepted Version

Article:

Rowen, D. orcid.org/0000-0003-3018-5109, Brazier, J. orcid.org/0000-0001-8645-4780, Young, T. orcid.org/0000-0001-8467-0471 et al. (4 more authors) (2011) Deriving a Preference-Based Measure for Cancer Using the EORTC QLQ-C30. *Value in Health*, 14 (5). pp. 721-731. ISSN 1098-3015

<https://doi.org/10.1016/j.jval.2011.01.004>

Article available under the terms of the CC-BY-NC-ND licence
(<https://creativecommons.org/licenses/by-nc-nd/4.0/>)

Reuse

This article is distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs (CC BY-NC-ND) licence. This licence only allows you to download this work and share it with others as long as you credit the authors, but you can't change the article in any way or use it commercially. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

Published in final edited form as:

Value Health. 2011 ; 14(5): . doi:10.1016/j.jval.2011.01.004.

DERIVING A PREFERENCE-BASED MEASURE FOR CANCER USING THE EORTC QLQ-C30

Donna Rowen^a, John Brazier^a, Tracey Young^a, Sabine Gaugris^b, Madeleine T King^c, Benjamin Craig^{d,e}, and Galina Velikova^f

^aHealth Economics and Decision Science, School of Health and Related Research, University of Sheffield, UK

^bJanssen-Cilag Ltd., High Wycombe, UK

^cPsycho-oncology Co-operative Research Group (PoCoG), School of Psychology, University of Sydney, Australia

^dHealth Outcomes & Behavior, Moffitt Cancer Center, US

^eDept of Economics, University of South Florida, US

^fCancer Research UK Clinical Centre, Leeds, UK

Abstract

Aims—The EORTC QLQ-C30 is one of the most commonly used measures in cancer but in its current form cannot be used in economic evaluation as it does not incorporate preferences. We address this gap by estimating a preference-based single index for cancer from the EORTC QLQ-C30 for use in economic evaluation.

Methods—Factor analysis, Rasch analysis and other psychometric analyses were undertaken on a clinical trial dataset of 655 patients with Multiple Myeloma to derive a health state classification from the QLQ-C30 that is amenable to valuation. A valuation study was conducted of 350 members of the UK general population using ranking and time trade-off. A series of regression models were fitted to the data, including the episodic random utility model (RUM) to derive preference weights for the classification system.

Results—The resulting health state classification system has 8 dimensions (physical functioning, role functioning, social functioning, emotional functioning, pain, fatigue and sleep disturbance, nausea, and constipation and diarrhoea) with 4 or 5 levels each. Mean and individual level additive multivariate regression models were estimated and compared. Mean absolute error ranges from 0.050 to 0.054 with no systematic errors. All models have few inconsistencies (0 to 2) in estimated preference weights.

Conclusions—It is feasible to derive a preference-based measure from the EORTC QLQ-C30 for use in economic evaluation, but this work needs to be extended to other countries and replicated across other conditions.

Keywords

Preference-based measures; QALYs; EORTC QLQ-C30

Corresponding author: Donna Rowen, Health Economics and Decision Science, School of Health and Related Research (SchARR), University of Sheffield, Regent Court, 30 Regent Street, Sheffield S1 4DA. d.rowen@sheffield.ac.uk, Tel: +44(0)114 222 0728. Fax: +44(0)114 272 4095.

The article has not been published elsewhere and there are no conflicts of interest.

Introduction

Generic preference-based measures, such as the EQ-5D (Brooks, 1996), HUI3 (Feeny et al, 2002) or SF-6D (Brazier et al, 2002) are widely used to calculate quality-adjusted life years (QALYs) (Drummond, 1994) for use in economic evaluation. The EQ-5D is the most common generic preference-based measure (PBM) and is currently recommended by NICE (NICE, 2008). However, generic measures of health have been found to be inappropriate or insensitive for some medical conditions (Brazier et al, 2007) and for cancer in particular (Garau et al, 2009). Furthermore, clinicians and researchers often choose to include condition-specific measures such as the EORTC QLQ-C30 in trials rather than generic preference-based measures because they need measures which are sensitive to the effects of interventions across a range of relevant symptoms, side effects and aspects of functioning and quality of life. Condition-specific measures, such as the EORTC's core quality of life questionnaire, the QLQ-C30, have great clinical utility because they summarise a number of symptom and domain-specific scales. However, because they are not preference-based, they provide a description rather than a valuation of health, and therefore cannot be used to estimate QALYs.

There are three ways in which researchers can estimate QALYs for a trial where a generic preference-based measure such as the EQ-5D is unavailable: undertake 'mapping'; value vignettes that describe the health states covered by patients in the trial; or derive a preference-based measure from the existing condition-specific measure.

Mapping (also known as 'cross-walking' or estimating exchange rates between instruments) involves predicting the relationship between the non-preference-based measure, for example the EORTC QLQ-30, and a generic preference-based measure, for example the EQ-5D, using statistical association. Mapping by statistical association may be considered less arbitrary than using the judgement of 'experts' to map between measures. Typically mapping by statistical association uses two datasets; an estimation dataset that contains respondents' self-reported scores for their own health using, for example, EORTC QLQ-C30 and EQ-5D and the study dataset that contains only EORTC QLQ-C30. A statistical relationship between the EORTC QLQ-C30 and EQ-5D is estimated using regression techniques on the estimation dataset and the results are applied to the study dataset to obtain predicted EQ-5D health state utility values. Mapping is a second best alternative to the use of a preference-based measure directly in the study as mapped estimates can have large errors, most noticeably when mapping from condition-specific measures to generic preference-based measures (Brazier et al, 2009). Mapping requires a degree of overlap between the descriptive systems of both measures, that the relationship estimated in the estimation dataset is generalisable to the study dataset, and that both measures are administered on the same population. Yet this means that mapping is valid only if both measures are appropriate for the patient population, which is unlikely to be the case for generic preference-based measures administered to cancer patients (Garau et al, 2009).

An alternative solution to estimate QALYs when generic preference-based measures are unavailable or inappropriate is to value a selection of bespoke descriptions or 'vignettes' that describe the health states covered by patients in the trial. However patients in the trial experience a variety of different health states and as only a selection will be valued this will not fully take into account variation across individuals and across treatments. It is also unlikely that identical health states are experienced in trials for different treatments and hence vignettes need to be created and valued for each trial which reduces comparability across trials and treatments.

It has therefore been argued that a better approach in many cases would be to develop a preference-based measure from the condition-specific questionnaire specifically designed for that condition (Brazier and Dixon, 1995). Typically this requires reducing the length of the questionnaire to obtain a health state classification system that remains responsive and valid for the condition but that is amenable to valuation. Preference weights for the health state classification are then obtained from a representative sample of the general population.

This study applies the methods originally developed in the estimation of a generic preference-based measure of health from the SF-36 (Brazier et al, 1998; 2002) and subsequently used with condition specific measures in urinary incontinence (Brazier et al, 2008), asthma (Yang et al, 2007; Young et al, 2007) and overactive bladder (Yang et al, 2009; Young et al, 2009a). The study involved three stages. Firstly, a health state classification system was derived from the EORTC QLQ-C30 that is amenable to valuation using a recognised preference elicitation technique. The classification system was derived using psychometric analysis, factor analysis and Rasch analysis on a dataset of patients with Multiple Myeloma. Secondly a valuation survey was conducted asking members of the general population to value a sample of states defined by the classification. Thirdly regression models were estimated on the results of the valuation survey to estimate the preference weights to produce a utility estimate for every health state defined by the classification system.

EORTC QLQ-C30

The EORTC QLQ-C30 is one of the most commonly used measures in cancer (Aaronson et al, 1993) and dominates cancer clinical trials in Europe and Canada. The EORTC QLQ-C30 contains 30 questions that cover the most common symptoms of cancer (such as pain, fatigue, nausea and vomiting) and various aspects of function (including physical, role, social, emotional and cognitive functioning). The EORTC QLQ-C30 is summarized using fourteen scales, each representing a particular symptom or aspect of function, plus one global quality of life scale (based on two global questions). Its validity has been well established for many conditions in cancer.

While the EORTC QLQ-C30 has proved to be a useful instrument for demonstrating treatment benefits, it cannot be used in economic evaluation in its current form because it does not incorporate preference information. While it generates a profile of scores representing a range of symptoms and aspects of functioning and a global quality of life score, it does not currently generate a single preference-based index of quality of life required for economic evaluation using QALYs. Further, the number of items and scales is too large to be amenable to valuation using preference elicitation techniques such as time trade off and standard gamble.

EORTC QLQ-C30 Multiple Myeloma patient dataset

The dataset used to derive the health state classification system contains patients with Multiple Myeloma collected in a randomized controlled trial, Phase III VISTA trial (FDA number...) undertaken in 2007 (Sabine can you provide correct info please). Patients were asked to complete the EORTC QLQ-C30 at screening, each cycle of treatment (cycles 1–9), end of treatment and post treatment (cycles 1–17). The screening phase of the dataset (n=655) is used to select items for the health state classification. Data from cycle 5 of treatment in the trial is then used to validate the choice of items. This study forms part of a wider cross-country study across different cancer patient groups and the health state classification will be validated using datasets of patients with different types of cancers.

Methods

Methodology used to derive a health state classification from the QLQ-C30

The aim was to produce a multidimensional health state classification from the EORTC QLQ-C30 that is amenable to valuation by respondents with a minimum loss of information and subject to the constraint that responses to the original instrument can be unambiguously mapped onto it. This implies that the text of the items should be altered as little as possible. The task is therefore to determine the dimensions, items and the levels of the EORTC QLQ-C30 to be included in the new classification. The methodology outlined here uses a combination of Rasch and classical psychometric analysis (Young et al. 2009a). SPSS version 15 (SPSS Inc, 2005) was used for the factor analysis and Rasch Unidimensional Measurement Models (RUMM2020) (RUMM Laboratory Pty Ltd, 1997–2004) was used for the Rasch analysis.

Dimensional structure—Multidimensional health state classifications should have structural independence between dimensions to avoid nonsensical states (Feeny et al, 2002). In other words, there must be little correlation between the dimensions in our classification system. The large literature on the EORTC QLQ-C30 focuses upon its use as a profile measure of health, but here we wish to determine the dimensions across all items, ignoring whether these are functions or symptoms.

Factor analysis can be used for a set of observed variables to identify structurally independent dimensions by highlighting underlying factors that explain patterns of correlation (Chatfield and Collins, 1980). We applied factor analysis to 27 of the 30 items of the EORTC QLQ-C30, (excluding global quality of life and financial impact items as these are inappropriate for a PBM of HRQoL) to explore the dimension structure of the EORTC QLQ-C30. The dimensions were determined using a varimax component matrix and eigenvalues. The extent to which items belong to a single dimension can also be examined using Rasch analysis (see below). The results were discussed with our team's clinical expert (GV) to make sure that they made sense clinically before making a final decision on the dimensionality of the health state classification.

Item selection—Each dimension of a health state classification system of a preference-based measure is usually represented by just one or two items to render the system amenable to preference elicitation methods. We used the following conventional psychometric criteria to help select items from the QLQ-C30: distribution of responses across categories of response (including floor effects and ceiling effects), the percentage of missing data, correlation of item to dimension, and responsiveness to change over two points in time using the standardised response mean (SRM).

A further technique often used (Young, 2007, 2009a, 2009b; Mavranouzouli, 2009) to select items is Rasch analysis (Rasch, 1960). This is a mathematical technique that converts qualitative (categorical) responses to points on a continuous (unmeasured) latent scale using a logit model (Tesio, 2003). It can be used to assess whether an item fits the model, the severity of health problem being covered by each item, the extent to which items have response choices that are appropriately ordered as responders can distinguish between the response levels for a given item and whether items perform differently between populations (known as differential item functioning (DIF)) (Young et al, 2009). Items that fit the Rasch model, cover the full range of severity, have ordered response choices and do not suffer from DIF were considered as candidates for inclusion in the health state classification. Items that did not fit the Rasch model (using criteria that item level Chi-square P-value<0.01) were removed; all other items were retained and the Rasch analysis was re-estimated. We used the following criteria to assess the Rasch model for each dimension: item-trait interaction

(whether data fitted the Rasch model for groups of responders with similar underlying health); person separation index (PSI) (whether the Rasch model could discriminate between responders); item fit and person fit residuals (the divergence between expected and observed responses per respondent); and item range and spread at logit zero (whether items covered a wide range of severity).

The final selection of dimensions, items and levels for the health state classification was based on what appeared to perform best using the psychometric tests and Rasch analysis and at the same time ensuring that health states made clinical sense and were amenable to valuation by respondents. The process involved judgment by our clinical expert (GV) and consideration of other factors such as wording of the health state classification system.

Valuation study to obtain preferences for the health state classification

The second stage of the research was to obtain valuations of states defined by the health state classification system. Key methodological issues were the choice of technique for eliciting preferences, the sample of states valued, sampling of respondents and overall size of sample.

The main valuation technique used to value health states was Time Trade-off (TTO). States were sampled using an orthogonal array in order to enable the estimation of an additive model for the preference weights. Use of an orthogonal array to sample states is common when dimensions are independent. SPSS version 15 was used to produce a sample of 81 states using orthogonal array and this was supplemented by 4 additional states. We chose to value 85 states in order to enable each respondent to value the worst state, an equal number of responses per state, and an equal number of states to be valued per respondent.

At the interview, respondents read the descriptive system and self-completed the system for their own health. The EORTC QLQ-C30 does not mention cancer in its descriptive system and therefore respondents were not aware that the health states were used to describe the health-related quality of life of cancer patients. Respondents were asked to rank 8 states alongside 'full health' and 'dead' to help familiarize them with the states. Respondents then valued the same 8 states using the Measurement and Valuation of Health (MVH) study version of TTO that involves the use of a visual prop designed by the MVH group (University of York) (see Dolan, 1997). Respondents were initially taken through a hypothetical TTO to help them understand the task. For each health state respondents were first asked whether they would prefer the given health state for 10 years or to die immediately. For health states considered better than being dead (BTD), respondents were asked to choose between (a) health state h for w years, after which they will die, or (b) full health for z years ($z \leq w$), after which they will die. While w is fixed at 10 years, years in full health, z , is varied to determine the point where respondents are indifferent between the two options. For health states considered worse than being dead (WTD) respondents were asked to choose between (a) health state h for w years followed by full health for z years after which they will die, or (b) immediate death. Both years in optimal health, z , and years in health state, $w=10-z$, are varied to determine the point where respondents are indifferent between the two options. After the collection of trade-off responses (w, z), respondents were asked a number of background questions covering demographic and socio-economic characteristics.

A representative sample of the general population was interviewed in their own home by trained and experienced interviewers who had worked on numerous previous valuation surveys, such as the HUI2 (McCabe et al, 2005) and OAB-5D (Yang et al, 2009).

The valuation data

Respondents were from the geographical areas in the North of England including urban and rural areas with a mix of socio-economic characteristics. There were 350 successfully conducted interviews, a response rate of 40.3% for suitable respondents answering their door at time of interview. Six respondents were excluded from the analysis; three were excluded for valuing all states as identical and less than one; two respondents were excluded for valuing the worst possible health state higher than every other state; one respondent was excluded for valuing all states as worse than dead. All other responses were used in the analysis reported here. The remaining 344 respondents had a health state completion rate of 98.5% in the TTO tasks. Characteristics of the included respondents are compared to the general population in South Yorkshire and England (Table 1). The valuation sample contained a higher proportion of people aged 41–65, retired people, more females and more people in poorer health than in the population at large.

Modelling to obtain preference weights for the health state classification

The TTO responses (z , w) were analysed using a range of different specifications. The standard specification is based on the approach first used for the UK EQ-5D preference weights (Dolan, 1997):

$$y_{ij}^s = f(\mathbf{X}_{\delta\lambda}\boldsymbol{\beta}) + \varepsilon_{ij}^s, \quad y_{ij}^s = \begin{cases} 1 - z_{ij}/w_{ij} & \text{if state better than dead} \\ 1 + z_{ij}/10 & \text{if state worse than dead} \end{cases} \quad (1)$$

Where $i=1,2 \dots n$ represents individual health states and $j=1,2 \dots m$ represents respondents. The dependent variable y_{ij}^s is disvalue for health state i valued by respondent j and $\mathbf{X}_{\delta\lambda}$ is a vector of dummy explanatory variables for each level λ of dimension δ of the health state classification. Level $\lambda = 1$ acts as a baseline for each dimension.

The second specification is the episodic random utility model (ERUM), where the value of the health state depends on its duration, w_{ij} :

$$y_{ij}^e = f(w_{ij}\mathbf{X}_{\delta\lambda}\boldsymbol{\beta}) + \varepsilon_{ij}^e, \quad y_{ij}^e = \begin{cases} w_{ij} - z_{ij} & \text{if state better than dead} \\ w_{ij} + z_{ij} & \text{if state worse than dead} \end{cases} \quad (2)$$

In order to produce error terms on the same scale as the standard specification in equation (1), both z_{ij} and w_{ij} are divided by 10 before estimation.

The standard approach transforms WTD TTO responses to be bound at -1 . This has been criticised as there is little empirical evidence for why values should be bound at -1 and arguably the transformed responses cannot be interpreted as being measured on the same utility scale as states BTM (Patrick et al, 1994). The ERUM model was developed to deal with this criticism by changing the way TTO responses are modelled. Under the ERUM, WTD TTO responses are not transformed and are therefore modelled in a way that is consistent with BTM responses (Craig and Busschbach, 2009).

Each model was estimated using ordinary least squares (OLS) estimation, which assumes that each observation is independent, regardless of the fact that there will be on average 8 observations per individual. Random and fixed effects were estimated to take into account within and between respondent variation, allowing for the fact that there are multiple observations per individual (Brazier et al, 2002). For the random and fixed effects model the error term, ε_{ij} is subdivided as follows:

$$\varepsilon_{ij}=u_j+e_{ij} \quad (3)$$

Where u_j is individual random effect and e_{ij} is the usual random error term for the i th health state valuation of the j th individual, assumed to be random across observations. Maximum likelihood estimation was used for the estimation.

Instead of examining individual responses, the TTO responses for each state may be condensed into mean estimates and the standard model (equation 1) can also be estimated using mean level data:

$$\bar{y}_i=f(\mathbf{X}_{\delta\lambda}\boldsymbol{\beta})+\varepsilon_i \quad (4)$$

Where \bar{y}_i represents mean TTO disvalue of the i th health state, $i=1,2 \dots, 86$ health states averaged across all individuals who valued the i th health state, and ε_i is the error term. Estimation was via OLS on a dataset of 86 observations. The impact of adding interaction terms and socio-demographic variables is explored for all models.

Several alternative criteria to indicate model performance are reported. The difference between actual and predicted values is assessed using mean absolute error (MAE) calculated at the health state level. MAE is an indicator of how large the prediction errors are and reporting the number of health states valued with errors greater than 5% and 10% indicates whether the errors are of a minimal important difference. Inconsistencies in parameter estimates for adjacent levels of an item were noted as these indicated that worsening health did not lead to a lower utility value. Models in which these inconsistencies were removed by merging levels were also estimated. The number of main effects with insignificant coefficients is also reported. Performance of all regression models is reported using inconsistencies, significant coefficients, mean absolute error of health state predictions and MAE greater than 5% and 10% and by examining plots of actual and predicted health state values.

Results

Health state classification

Step 1: Instrument dimensionality—A four factor model on all 27 items accounted for 58.7% of the variance. Items were divided into the four factors according to their ‘loadings’, the correlation between the item and the factor. All items loaded >0.35 on a factor, but some items cross loaded. Factor 1 contained the majority of items (14), covering physical functioning, role functioning, social functioning, pain and items ‘need a rest’, ‘felt weak’ and ‘difficulty concentrating’. Factor 2 contained all items covering emotional functioning plus ‘trouble sleeping’. Factor 3 contained all items covering eating and digestion and factor 4 contained items ‘short of breath’, ‘were you tired’ and ‘difficulty remembering’. Items loading into factors 2 and 3 were conceptually clear and in accordance with the grouping of the original EORTC instrument. Items loading into factors 1 and 4 (17 items), on the other hand, covered a large range of concepts and further factor analysis was done on these items to determine whether further differentiation was possible. The additional item ‘trouble sleeping’ was added (18 items in total) over a concern that the initial factor analysis captured some causality rather than correlation.

A four factor model on the 18 items explained 67.6% of the variance. All items loaded >0.4 on a factor, but some items cross loaded (difference between cross loadings <0.2). Cross loading items were included in the factor that was clinically meaningful. The four factors

were: physical functioning, role functioning and pain; social functioning; fatigue, trouble sleeping and short of breath; cognitive functioning (principal-component rotated factor loadings are available from authors on request).

Table 2 shows the potential items categorised according to the dimensions for consideration for the measure amenable to valuation. Overall the 27 items can be divided into 6 factors or dimensions: (1) physical functioning, role functioning and pain, (2) social functioning, (3) emotional functioning, (4) digestion, (5) fatigue, trouble sleeping and shortness of breath and (6) cognitive functioning. After consultation with our clinical specialist (GV), cognitive functioning (items 20 and 25) and shortness of breath (item 8) were excluded from the PBM on the grounds that they are neither a symptom nor side effect of treatment for Multiple Myeloma, leaving 24 items remaining. The remaining five dimensions are used for the following analysis to determine the PBM, where items were fitted to Rasch models for each of the five dimensions.

Step 2: Selecting items by dimension—Table 2 presents psychometric analysis and goodness of fit for the Rasch models for each dimension.

Item-level ordering and differential item functioning: Only item 15 (digestion dimension) was disordered and for this item ‘a little’, ‘quite a bit’ and ‘very much’ were collapsed into one level prior to proceeding with further Rasch modelling. Four items demonstrated differential item functioning by sex and were split according to sex: items 1 and 3 (physical functioning, role functioning and pain dimension), item 15 (digestion) and item 22 (emotional functioning). Item 21 (emotional functioning) demonstrated differential item functioning by age and was split according to age. This indicated that these items were not ideal for the health state classification.

Rasch model goodness of fit: Three items demonstrated poor item fit (Chi-square P-value < 0.01) and were excluded from subsequent Rasch models estimated for each dimension: item 6 (physical functioning, role functioning and pain dimension), item 14 (digestion), item 11 (fatigue).

Physical functioning, role functioning and pain: This dimension covers three separate attributes of quality of life: physical functioning, role functioning and pain. Each item covers only one attribute and it is unlikely that an item on physical functioning, for example, will capture or reflect role functioning or pain. In order to accurately represent the entire dimension, and in accordance with the EORTC QLQ-C30 scaling and with clinical opinion, we decided that a minimum of one item each for physical functioning, role functioning and pain was required.

Out of the 5 items that capture physical functioning, items 4 and 5 did not capture the full range of severity but had the highest SRM suggesting better responsiveness and ability to capture severe health problems, though in all cases they were low according to Cohen’s criteria (Cohen, 1978). Item 2 overall performed well with relatively high item fit but had floor effects, relatively low SRM and limited range of severity. Figure 1 shows the item map for all items in this dimension, reporting three thresholds between response categories “not at all” and ‘a little’, ‘a little’ and ‘quite a bit’, and ‘quite a bit’ and ‘very much’. No item captured the full severity range in terms of coverage as items 3, 4 and 5 captured severe health whereas items 1 and 2 captured less severe health. Given that none of the items were ideal, items 2 and 3 were chosen in order to cover the full severity range and there is a coherence since both measure trouble walking. Items 2 and 3 were merged to a five level item in the health state classification, with levels 1 to 4 taken directly from item 2 (no/a

little/quite a bit/very much trouble taking a long walk) and level 5 (very much trouble taking a short walk outside of the house) taken from level 4 of item 3.

Of the two role functioning items, item 7 had ceiling effects, but also had good Rasch model item fit and relatively good item range (Figure 1). Item 6 performed similarly to item 7 for the psychometric analysis, but had poor Rasch model item fit. Therefore item 7 was selected for the health state classification to represent role functioning.

Items 9 and 19 capture pain. Item 9 had a large severity range, but poor item fit, very high item fit residual and low item level p-value suggesting the item contributes poorly to the dimension. Item 19 had a more limited range and higher ceiling effects but better item fit. Item 19 was chosen due to better item fit and due to its wording, since it measures the extent to which pain interferes with daily activities rather than the existence of pain per se, which is likely to be more important to respondents.

Social functioning: Of the two social functioning items, item 26 had a slightly higher degree of ceiling effects, higher item fit residual and lower p-value. Further, it may not be applicable to all respondents as it captures whether physical condition or medical treatment interferes with family life. Item 27 was chosen as it performs marginally better psychometrically and is arguably applicable to a higher number of respondents as it measures interference with social activities.

Emotional functioning: Items 21 and 22 suffered from DIF. Items 23 and 24 performed similarly across all criteria. Item 23 had a larger severity range than item 24, but also had a higher item fit residual, suggesting greater divergence between expected and observed responses. Item 24 was chosen as it overall performs best, had a higher SRM and was felt to be more clinically relevant to patients with cancer.

Digestion: None of the digestion items performed well in psychometric or Rasch analysis. Despite this they were included in the health state classification as it was felt that they were clinically important for patients with Multiple Myeloma. These items capture multiple symptoms of digestion-related problems: lack of appetite, nausea, vomiting, constipation and diarrhoea. Lack of appetite, nausea and vomiting (items 13, 14, 15) are all closely related symptoms, and constipation and diarrhoea (items 16 and 17) are bowel symptoms, suggesting the items can be separated into two attributes. Item 13 was the only item from appetite, nausea and vomiting that did not suffer from DIF, item level disordering or poor item fit. However, this item performed poorly in the Rasch model with small coverage at logit zero and high item fit residuals. After consultation amongst our research team, including our clinical specialist, it was decided that item 13, which captures lack of appetite, would not be chosen because it maybe thought to be a desirable (positive) symptom by some people and may be a symptom due to the age of the population rather than the condition. Therefore, despite suffering from problems in the Rasch model items 14 and 15 were considered as it was felt to be important to capture appetite, nausea or vomiting in the health state classification. Both items suffer from high ceiling effects and have low SRM, with item 14 performing marginally better. Therefore item 14 (nausea) was chosen for the health state classification.

Items 16 and 17 on constipation and diarrhoea perform similarly, both suffering from extreme ceiling effects, small spread at logit zero and high residuals. It was decided to combine these items as they are both bowel symptoms where respondents rarely suffer from both during a weekly period. The items were combined such that level 1 of the merged item captures no bowel (constipation or diarrhoea) problems, levels 2, 3 and 4 capture 'a little' 'quite a bit' and 'very much' constipation and/or diarrhoea.

Fatigue and sleep disturbance: Items 10, 12 and 18 performed similarly in the Rasch and psychometric analyses, with large severity range and no large ceiling or floor effects. However item 10 had a relatively high item fit residual and item 12 has a low item p-value indicating it does not contribute well to the dimension in Rasch models. Item 18 performed marginally better, with a relatively high p-value, low residual, large coverage at logit zero and large range, and so was selected for the health state classification.

Health state classification: The classification system has 8 dimensions made up of 10 items. Table 4 summarises the items chosen from the EORTC QLQ-C30 for each dimension of the health state classification. Table 5 presents the final health state classification system of dimensions and their levels. A health state is made up of 8 sentences and hence has an 8 digit identifier, from best state 11111111 to worst state 54444444. This system generates a total of 81,920 health states.

Valuation survey

Table 6 presents mean descriptive statistics for observed TTO values (generated using the formula from Dolan, 1997) for all 86 states included in the valuation survey. Mean TTO values varied from 0.95 for best state to 0.13 for worst state, meaning that on average all states were valued as better than being dead. Of the total 2710 TTO observations, 514 observations (19%) were equal to 1 (equivalent in value to full health) and 271 (10%) were less than or equal to 0 (valued as the same or worse than being dead).

Modelling health state values

Table 7 presents the preference weights estimated using a variety of regression models. All coefficients have the expected sign (i.e., level 1 on each dimension is the reference point and higher levels increase TTO disvalue), their size is consistent with the severity scale in all but two cases (i.e., higher levels have larger coefficients and an increasing increment on TTO disvalue except physical functioning levels 4 and 5 and nausea levels 2 and 3) and the majority of coefficients are statistically significant. Models (1), (2), (4) and (5) are based on the standard specification outlined in equations (1) and (4), model (3) uses the ERUM specification in equation (2). Models (1) and (3) are individual level models estimated using OLS, model (2) is a random effects model estimated using maximum likelihood (the results of the Hausman test confirmed that a fixed effects model was not appropriate) and models (4) and (5) are mean level models. Model (5) is a consistent version of model (4) where adjacent inconsistent levels are merged into a common dummy variable.

Mean absolute error was similar between models ranging from 0.046 to 0.054. The number of health states with errors greater than 5% ranges from 33 to 41 and errors greater than 10% ranges from 6 to 13. Models including interaction effects and socio-demographic were estimated (available from the authors by request) but predictive ability, inconsistencies and significant coefficients for the main effects variables were not improved.

Discussion

We have estimated a preference-based measure for the EORTC QLQ-C30 using methodology first applied in the development of the SF-6D from the SF-36 (Brazier, 2002) and a number of condition specific measures. The EORTC-8D was constructed using psychometric analysis (including Rasch) to ensure that chosen items appropriately reflected their dimension and that each dimension covered a wide range of severity. A sample of states was valued and then the results of the survey modeled using a variety of methods, including a new ERUM method. This enables utility scores to be generated directly from EORTC QLQ-C30 datasets.

An important concern is that often condition-specific measures fail to capture co-morbidities and side-effects of treatments, and hence are not strictly comparable to generic measures when used to estimate QALYs for resource allocation. One way to enhance comparability across measures is for all measures to use the same methodology to derive values (Brazier et al, 2007). Our valuation study followed the methodology used in the development of generic measures: we implemented the protocol used to derive the UK EQ-5D preference weights (Dolan, 1997); used common anchors of 1 for full health and zero for dead; and interviews were conducted using a sample of the general population. Furthermore the EORTC-8D descriptive system captures a wide range of dimensions including generic dimensions such as physical functioning and role functioning, as well as more condition-specific symptoms such as nausea, constipation and diarrhoea. This means that the descriptive system is likely to capture overall health-related quality of life including both comorbidities and side-effects. Indeed, as the EORTC-8D has few condition-specific functionings and symptoms the measure may not be sufficiently sensitive to measure health-related quality of life of cancer patients. An area of future research will be to compare the measure to generic measures in terms of sensitivity and validity.

Given that the health state classification measures cancer it is somewhat surprising that all states have a positive mean TTO value and hence at the aggregate level all states are valued as being better than being dead. This is contrast to other valuation studies such as the UK EQ-5D valuation study where 38% (16/42) of health states valued were on average valued as being worse than dead (with mean TTO below zero). One hypothesis is that respondents would view the states differently if cancer was included in the health state description, and this is a topic for future research.

The final preferred model estimating preference weights should have no inconsistencies; health state utility values should always decrease as health states become more severe. The main inconsistency was observed for physical functioning levels 4 and 5. This may have been due to the merging of 2 items in the EORTC QLQ-C30 to form this dimension. Models (3) and (5) perform best overall according to the criteria of predictive ability, inconsistencies and significant coefficients, with model (3) performing better using all criteria. Model (5) has a predicted range of utilities from 1 to 0.199 whereas model (3) has a predicted range of 1 to 0.291, meaning that the worst state defined by the classification has a much lower value using preference weights estimated using model (5). Deciding upon the preferred model comes down to a choice between the mean model (5) with no inconsistencies, as chosen both in the valuation of the SF-6D (Brazier and Roberts, 2004) and the overactive-bladder-specific measure (Yang et al, 2009), or the recently developed technique of the ERUM. Although model (5) is in accordance with the recommended value set of many similar measures, the ERUM model (3) is here the preferred model as it more appropriately deals with TTO values for SWD and performs best.

The EORTC-8D was developed out of a concern that generic measures were not appropriate to measure the quality of life of cancer patients (Garau et al, 2009). The EORTC-8D enables QALYs to be directly estimated using the EORTC QLQ-C30, a questionnaire typically included in cancer trials, rather than the use of generic measures that are less appropriate (Garau et al, 2009) or mapping to generic measures that is both less appropriate and increases error around utility estimates. It is hoped that this measure will provide appropriate and useful information for cost per QALY analysis undertaken in cancer trials.

A general concern regarding the development of PBMs from existing questionnaires is that the classification system is strongly influenced by the specific patient dataset used to develop the classification. The patient dataset and valuation study used here are both UK trial datasets and hence the EORTC-8D classification and preference weights presented here

are most appropriate for UK trials. The health state classification was developed using Multiple Myeloma patients. Further testing of the classification will be undertaken across datasets of cancer patients with different types of cancers. This study forms part of a wider cross-country study that will examine the use of preference-based measures from the EORTC QLQ-C30 on a variety of countries and different patient groups.

Acknowledgments

Financial support: The studies reported in this paper were funded by Janssen-Cilag Ltd.

We would like to thank the Centre for Research and Evaluation at Sheffield Hallam University for conducting the interviews. The studies reported in this paper were funded by Janssen-Cilag Ltd.

References

- Aaronson NK, Ahmedzai S, Bergman B, Bullinger M, Cull A, et al. The European Organization for Research and Treatment of Cancer QLQ-C30: a quality-of-life instrument for use in international clinical trials in oncology. *Journal of the National Cancer Institute*. 1993 Mar 3; 85(5):365–76. [PubMed: 8433390]
- Brazier JE, Dixon S. The use of condition specific outcome measures in economic appraisal. *Health Economics*. 1995; 4:255–264. [PubMed: 8528428]
- Brazier JE, Harper R, Thomas K, Jones N, Underwood T. Deriving a preference based single index measure from the SF-36. *J Clinical Epidemiology*. 1998; 51 (11):1115–1129.
- Brazier J, Roberts J, Deverill M. The estimation a preference-based single index measure for health from the SF-36. *Journal of Health Economics*. 2002; 21(2):271–292. [PubMed: 11939242]
- Brazier J, Roberts J. The estimation of a preference-based index from the SF-12. *Medical Care*. 2004; 42:851–859. [PubMed: 15319610]
- Brazier, JE.; Ratcliffe, J.; Tsuchiya, A.; Solomon, J. *Measuring and valuing health for economic evaluation*. Oxford: Oxford University Press; 2007.
- Brazier JE, Czoski-Murray C, Roberts J, Brown M, Symonds T, Kelleher C. Estimation of a preference-based index from a condition specific measure: the King's Health Questionnaire. *Medical Decision Making*. 2008; 28(1):113–126. [PubMed: 17641139]
- Brazier J, Yang Y, Tsuchiya A, Rowen D. A review of studies mapping (or cross walking) from non-preference based measures of health to generic preference-based measures. *European Journal of Health Economics*. 2009 Forthcoming.
- Chatfield, C.; Collins, AJ. *Introduction to Multivariate Analysis*. Chapman and Hall; 1980.
- Cohen, J. *Statistical power analysis for the behavioural sciences*. New York: Academic Press; 1978.
- Craig BM, Busschbach JJV. The Episodic Random Utility Model Unifies Worse Than Death and Better Than Death TTO Responses in Health State Valuation. *Population Health Metrics*. 2009; 7(3):1–10. [PubMed: 19126218]
- Dolan P. Modelling valuations for EuroQol Health States. *Medical Care*. 1997; 35 (11):1095–1108. [PubMed: 9366889]
- Drummond, M. *Economic Evaluation Alongside Clinical Trials*. Department of Health; London: 1994.
- Drummond, MF.; O'Brien, BJ.; Stoddart, GL.; Torrance, GW. *Methods for the economic evaluation of health care programmes*. Oxford: Oxford Medical Publications; 1997.
- Feeny D, Furlong W, Torrance GW, Goldsmith CH, Zhu Z, DePauw S, Denton M, Boyle M. Multi-attribute and single-attribute utility functions for the Health Utilities Index Mark 3 system. *Medical Care*. 2002; 40 (2):113–128. [PubMed: 11802084]
- Garau M, Shah K, Towse A, Wang Q, Drummond M, Mason A. Assessment and appraisal of oncology medicines: does NICE's approach include all relevant elements? What can be learnt from international HTA experiences? Report for the Pharmaceutical Oncology Initiative (POI). 2009 Feb.
- Kind, P.; Hardman, G.; Macran, S. *UK Population Norms for EQ-5D*. University of York; 1999. Centre for Health Economics Discussion Paper Series

- Mavranzeouli I, Brazier J, Young TA, Barkham M. Using Rasch analysis to form plausible health states amenable to valuation: the development of CORE-6D from CORE-OM in order to elicit preferences for common mental health problems. *Health Economics and Decision Sciences Discussion Paper*. 2009
- McCabe C, Stevens K, Roberts J, Brazier J. Health state values for the HUI 2 descriptive system: Results from a UK survey. *Health Economics*. 2005; 14:231–244. [PubMed: 15386655]
- NICE. Guide to the methods of technology appraisal. 2004. (http://www.nice.org.uk/pdf/TAP_Methods.pdf)
- NICE. Guide to the methods of technology appraisal. NICE; London: 2008. <http://www.nice.org.uk/aboutnice/howwework/devnicetech/technologyappraisalprocessguides/guidetothemethodsoftechnologyappraisal.jsp>
- Rasch, G. Probabilistic models for some intelligence and attainment tests. Chicago: University of Chicago Press; 1960. Reprinted 1980
- Rasch Unidimensional Measurement Models (RUMM) 2020. RUMM Laboratory Pty Ltd; 1997–2004.
- SPSS for Windows. Release. 14.0.1. 2005. Chicago: SPSS Inc; 2005.
- Stevens K, Brazier J, McKenna S, Doward L, Cork M. The development of a preference-based measure of health in children with atopic dermatitis. *British Journal of Dermatology*. 2005; 153:372–377. [PubMed: 16086752]
- Tesio L. Measuring behaviours and perceptions: Rasch analysis as a tool for rehabilitation research. *Journal of rehabilitation medicine*. 2003; 35(3):105–15. [PubMed: 12809192]
- Yang Y, Tsuchiya A, Brazier J, Young T. Estimating a preference-based single index from the Asthma Quality of Life Questionnaire (AQLQ). *Health Economics and Decision Sciences Discussion Paper*. 2007
- Yang Y, Brazier JE, Tsuchiya A, Coyne K. Estimating a preference-based index from the Over Active Bladder questionnaire. *Value in Health*. 2009; 12(1):159–166. [PubMed: 18647258]
- Young, T.; Yang, Y.; Brazier, J.; Tsuchiya, A. The use of Rasch analysis as a tool in the construction of a preference based measure: the case of AQLQ. HEDS Discussion Paper Series No. 07/01. 2007. <http://www.shef.ac.uk/scharr/sections/heds/discussion.html>
- Young T, Yang Y, Brazier J, Tsuchiya A, Coyne K. The first stage of developing preference-based measures: constructing a health-state classification using Rasch analysis. *Quality of Life Research*. 2009a; 18:253–265. [PubMed: 19082759]
- Young T, Rowen D, Brazier J, Norquist J, Ambegaonkar B, Sazonov V. Developing Preference-Based Health Measures: Using Rasch Analysis to Generate Health State Values. *Health Economics and Decision Sciences Discussion Paper*. 2009b Forthcoming.

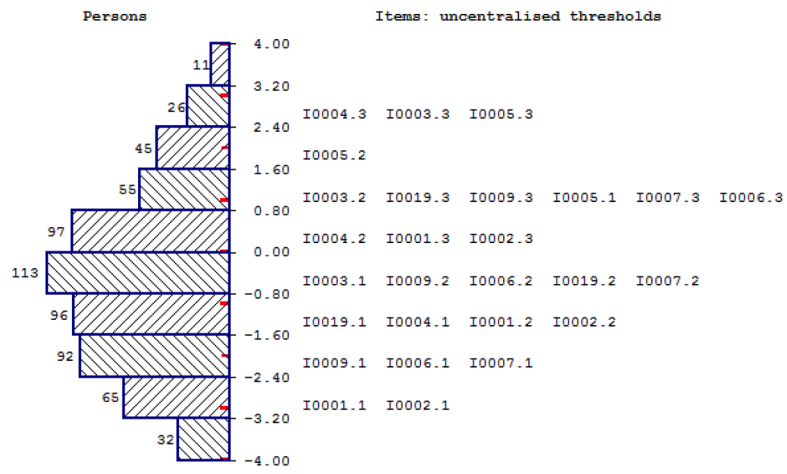


Figure 1.
Item map for physical functioning, role functioning and pain dimension

Table 1

Characteristics of respondents

	Included respondents (n=344) ¹	South Yorkshire ²	England ²
Mean age (s.d.)	47.8 (18.5)	NA	NA
Age distribution			
18–40	38.6%	41.2%	41.6%
41–65	42.7%	39.1%	39.1%
Over 65	17.2%	19.7%	19.3%
Female	61.9%	51.2%	51.3%
Married/Partner	57.0%	NA	NA
Employed or self-employed	43.6%	56.1%	60.9%
Unemployed	0.6%	4.1%	3.4%
Long-term sick	6.4%	7.7%	5.3%
Full-time student	7.3%	7.5%	7.3%
Retired	24.7%	14.4%	13.5%
Own home outright or with a mortgage	69.2%	64.0%	68.7%
Renting property	29.9%	36.0%	31.3%
Secondary school is highest level of education	41%	NA	NA
EQ-5D score (s.d.)	0.76(0.34)	NA	0.86(0.23) ³
TTO completion rate	98.5%		

¹Six respondents were excluded; three for valuing all states as identical and less than 1; two for valuing the worst possible state higher than every other state; one for valuing all states as worse than dead.

²Statistics for South Yorkshire Health Authority and for England in the Census 2001. Questions used in this study and the census are not identical. The census includes persons aged 16 and above whereas this study only surveys persons aged 18 and above. Age distribution is here reported as the percentage of all adults aged 18 and over.

³Interviews conducted in the Measurement and Valuation of Health (MVH) study in 1993 (Kind et al, 1999).

Table 2

Summary of psychometric and Rasch analysis used to select items per dimension

Items	Item summary	Psychometric analysis			Rasch analysis									
		% response at floor (very much)	% response at ceiling (not at all)	% missing data	Correlation with domain score	SRM	Item range	Residual	Item level chi-sq P-value	Spread at logit zero	Disordered	DIF characteristic	Poorly fitting item in Rasch model (Chi-sq P-value<0.01)	
Physical and role functioning and pain														
1	Trouble doing strenuous activities	28.5	14.4	0.6	-0.869	-0.206							Sex	
2	Trouble taking a long walk	26.9	16.9	1.1	-0.896	-0.195	-2.576 to 0.516	-1.540	0.385	0.93 to 0.37				
3	Trouble taking a short walk	8.1	46.7	1.7	-0.844	-0.235							Sex	
4	Need to stay in bed or a chair during the day	9.3	36	1.4	-0.788	-0.289	-1.158 to 2.385	0.874	0.120	0.76 to 0.08				
5	Need help with eating, dressing, washing or using the toilet	4.1	72.8	0.6	-0.575	-0.313	1.220 to 2.584	-0.678	0.083	0.23 to 0.07				
6	Limited in doing your work or daily activities	18.5	25.8	1.5	-0.953	-0.259								Yes
7	Limited in pursuing hobbies or other leisure time activities	17.9	30.1	1.2	-0.954	-0.237	-1.526 to -1.192	-0.629	0.600	0.82 to 0.23				
9	Pain	20	20.2	0.9	0.926	-0.499	-2.176 to 1.004	4.685	0.075	0.90 to 0.27				
19	Pain interfered with daily activities	20	30.4	0.5	0.936	-0.432	-1.427 to 0.987	-0.844	0.138	0.81 to 0.27				
Social functioning														
26	Physical condition or medical treatment interfered with family life	8.7	49.2	0.8	-0.899	-0.085	-1.247 to 2.013	1.146	0.175	0.78 to 0.12				

Items	Item summary	Psychometric analysis			Rasch analysis								
		% response at floor (very much)	% response at ceiling (not at all)	% missing data	Correlation with domain score	SRM	Item range	Residual	Item level chi-sq P-value	Spread at logit zero	Disordered	DIF characteristic	Poorly fitting item in Rasch model (Chi-sq P-value<0.01)
27	Physical condition or medical treatment interfered with social life	12.8	41.8	0.6	-0.942	-0.090	-1.869 to 1.144	0.555	0.345	0.87 to 0.24			
Emotional functioning													
21	Feel tense	6.6	37.7	0.9	-0.834	-0.264						Age	
22	Worried	12.8	26.0	0.2	-0.860	-0.477						Sex	
23	Feel irritable	4.3	44.3	0.5	-0.776	-0.223	-1.546 to 2.740	3.912	0.442	0.82 to 0.06			
24	Feel depressed	7.8	41.5	0.8	-0.834	-0.330	-1.696 to 2.017	-0.742	0.232	0.85 to 0.12			
Digestion													
13	Lacked appetite	8.9	48.9	0.5	No domain	-0.348	-1.349 to -0.314	-2.244	0.027	0.79 to 0.58			
14	Felt nauseated	2.4	73.9	0.3	No domain	-0.124							Yes
15	Vomited	1.2	87.9	0.3	No domain	-0.046					Yes	Sex	
16	Constipated	7.3	52.7	0.5	No domain	-0.225	-1.258 to -0.449	-1.220	0.231	0.78 to 0.61			
17	Had diarrhoea	0.8	79.5	0.9	No domain	0.058	-0.081 to 1.083	1.029	0.027	0.52 to 0.25			
Fatigue and trouble sleeping													
10	Needed to rest	15.9	14	1.4	No domain	-0.326	-3.113 to 2.212	1.409	0.881	0.96 to 0.10			
11	Trouble sleeping	10.1	41.4	0.2	No domain	-0.276							Yes
12	Felt weak	12.7	22.6	0.6	No domain	-0.296	-2.155 to 2.641	-0.433	0.043	0.90 to 0.07			
18	Tired	12.8	15.9	0.3	No domain	-0.245	-2.933 to 2.805	-0.727	0.540	0.95 to 0.06			

Note: DIF by sex is split male/female and DIF by age is split <65/65+.

Table 3

Goodness of fit for the Rasch model for each dimension

Item-trait interaction					
Dimension	Chi-sq (degrees of freedom)	P-value	Item fit (SD)	Person fit (SD)	Person separation index
Physical and role functioning and pain	121.36 (88)	0.01	-0.58(2.12)	-0.32(1.02)	0.90
Social functioning	10.82 (8)	0.21	0.85(0.42)	-0.44(0.80)	0.79
Emotional functioning	65.34 (54)	0.14	-0.18(2.14)	-0.35(0.96)	0.85
Digestion	72.08 (35)	0.00	-0.67(1.20)	-0.30(0.75)	0.47
Fatigue and trouble sleeping	22.05 (20)	0.34	0.08(1.16)	-0.65(1.27)	0.83

Table 4

Summary of EORTC QLQ-C30 items included in the EORTC-8D descriptive system

EORTC-8D dimension	EORTC QLQ-C30 items	Question
Physical functioning	2	Trouble taking a long walk
	3	Extra level added from 'trouble taking a short walk'
Role functioning	7	Limited in pursuing hobbies or other leisure time activities
Pain	19	Pain interfere with daily activities
Social functioning	27	Physical condition or medical treatment interfered with social life
Emotional functioning	24	Feel depressed
Nausea	14	Felt nauseated
	16	Constipated
Constipation and diarrhoea	17	Diarrhoea
	18	Tired

Table 5

EORTC-8D descriptive system

During the past week:

Physical functioning

You had no trouble taking a long walk
 You had a little trouble taking a long walk
 You had quite a bit of trouble taking a long walk
 You had very much trouble taking a long walk
 You had very much trouble taking a short walk outside of the house

Role functioning

You were not limited in pursuing your hobbies or other leisure time activities
 You were limited a little in pursuing your hobbies or other leisure time activities
 You were limited quite a bit in pursuing your hobbies or other leisure time activities
 You were limited very much in pursuing your hobbies or other leisure time activities

Pain

Pain did not interfere with your daily activities
 Pain interfered a little with your daily activities
 Pain interfered quite a bit with your daily activities
 Pain interfered very much with your daily activities

Emotional functioning

You did not feel depressed
 You felt a little depressed
 You felt quite a bit depressed
 You felt depressed very much

Social functioning

Your physical condition or medical treatment did not interfere with your social activities
 Your physical condition or medical treatment interfered a little with your social activities
 Your physical condition or medical treatment interfered quite a bit with your social activities
 Your physical condition or medical treatment interfered very much with your social activities

Fatigue and sleep disturbance

You were not tired
 You were a little tired
 You were quite a bit tired
 You were tired very much

Nausea

You did not feel nauseated
 You felt a little nauseated
 You felt nauseated quite a bit
 You felt nauseated very much

Constipation and diarrhoea

You were not constipated and did not have diarrhoea
 You were constipated and/or had diarrhoea a little
 You were constipated and/or had diarrhoea quite a bit
 You were constipated and/or had diarrhoea very much

Table 6

Observed and modelled TTO values for the 85 health states valued in the valuation survey

State	Observed				Modelled	
	Mean (SD)	Median	Range	Inter-quartile range	Model (3)	Model (6)
11111111	0.95 (0.17)	1.00	0.80	0.95 to 1.00	1.00	1.00
11122222	0.80 (0.27)	0.95	0.98	0.61 to 1.00	0.84	0.85
11132244	0.83 (0.20)	0.91	0.68	0.73 to 1.00	0.76	0.73
11221112	0.94 (0.18)	1.00	0.93	0.93 to 1.00	0.89	0.89
11323424	0.68 (0.40)	0.88	1.75	0.50 to 0.93	0.70	0.66
11341112	0.75 (0.29)	0.83	1.33	0.61 to 1.00	0.86	0.77
11413321	0.85 (0.17)	0.93	0.55	0.71 to 1.00	0.75	0.80
11432343	0.69 (0.31)	0.73	1.28	0.48 to 1.00	0.63	0.68
12142131	0.64 (0.45)	0.78	1.85	0.58 to 1.00	0.82	0.77
12222122	0.81 (0.26)	0.88	1.00	0.75 to 1.00	0.81	0.81
12334212	0.61 (0.28)	0.66	1.00	0.38 to 0.83	0.65	0.61
12341441	0.63 (0.28)	0.70	0.95	0.39 to 0.88	0.70	0.61
13114313	0.65 (0.38)	0.74	1.50	0.48 to 1.00	0.74	0.69
13124141	0.67 (0.30)	0.73	1.33	0.48 to 0.90	0.71	0.71
13211434	0.77 (0.29)	0.81	1.48	0.69 to 0.94	0.75	0.70
13423411	0.69 (0.28)	0.73	0.98	0.49 to 1.00	0.63	0.71
14131234	0.71 (0.34)	0.88	1.03	0.49 to 1.00	0.74	0.71
14214133	0.68 (0.30)	0.73	1.33	0.53 to 0.91	0.70	0.62
14243313	0.59 (0.45)	0.71	1.75	0.48 to 0.91	0.69	0.55
14411222	0.70 (0.35)	0.81	1.38	0.44 to 0.97	0.72	0.77
21122243	0.68 (0.47)	0.83	1.93	0.63 to 1.00	0.74	0.72
21123234	0.80 (0.22)	0.83	1.00	0.73 to 1.00	0.72	0.70
21231112	0.88 (0.15)	0.93	0.58	0.79 to 1.00	0.83	0.79
21241113	0.60 (0.39)	0.64	1.88	0.41 to 0.90	0.82	0.67
21314331	0.66 (0.39)	0.83	1.48	0.48 to 1.00	0.68	0.63
21414321	0.62 (0.36)	0.58	1.43	0.41 to 0.98	0.62	0.63
22113314	0.72 (0.32)	0.88	1.03	0.68 to 0.93	0.75	0.72

State	Observed				Modelled	
	Mean (SD)	Median	Range	Inter-quartile range	Model (3)	Model (6)
22114314	0.54 (0.49)	0.60	1.98	0.28 to 1.00	0.67	0.61
22212111	0.82 (0.28)	0.93	1.03	0.80 to 1.00	0.84	0.82
22222122	0.81 (0.23)	0.93	0.65	0.55 to 1.00	0.75	0.74
22441431	0.62 (0.27)	0.63	0.95	0.49 to 0.83	0.58	0.57
23141224	0.59 (0.41)	0.65	1.88	0.46 to 0.97	0.73	0.60
23314142	0.49 (0.49)	0.61	1.72	0.38 to 0.83	0.64	0.59
23333133	0.63 (0.43)	0.70	1.83	0.51 to 0.98	0.67	0.62
23431441	0.65 (0.33)	0.73	1.28	0.50 to 0.93	0.56	0.61
24141223	0.51 (0.49)	0.60	1.98	0.36 to 0.90	0.73	0.58
24213142	0.70 (0.23)	0.73	0.85	0.54 to 0.90	0.71	0.67
24322411	0.60 (0.42)	0.68	1.75	0.48 to 0.93	0.66	0.64
24432411	0.61 (0.33)	0.70	1.03	0.34 to 0.93	0.56	0.61
31143432	0.44 (0.44)	0.50	1.85	0.28 to 0.76	0.66	0.58
31144423	0.58 (0.29)	0.48	0.88	0.33 to 0.88	0.60	0.44
31212241	0.69 (0.39)	0.80	1.78	0.63 to 0.99	0.77	0.73
31312231	0.77 (0.28)	0.88	0.98	0.60 to 1.00	0.77	0.76
31421113	0.81 (0.24)	0.89	1.00	0.69 to 1.00	0.71	0.78
31431114	0.73 (0.27)	0.88	0.98	0.50 to 0.98	0.68	0.72
32121343	0.62 (0.43)	0.73	1.73	0.50 to 0.93	0.69	0.70
32234211	0.51 (0.42)	0.61	1.88	0.31 to 0.78	0.62	0.56
32344211	0.60 (0.29)	0.58	0.98	0.36 to 0.88	0.59	0.47
32413124	0.56 (0.49)	0.70	1.53	0.30 to 1.00	0.63	0.67
33121342	0.74 (0.28)	0.83	1.00	0.60 to 0.98	0.68	0.72
33231321	0.67 (0.33)	0.75	1.28	0.52 to 0.93	0.72	0.67
33312124	0.71 (0.27)	0.71	0.98	0.52 to 1.00	0.73	0.69
33333133	0.59 (0.39)	0.68	1.63	0.33 to 0.93	0.65	0.61
34112412	0.72 (0.28)	0.83	1.00	0.43 to 1.00	0.71	0.71
34113412	0.72 (0.29)	0.75	1.23	0.56 to 0.93	0.69	0.67
34221331	0.66 (0.25)	0.61	0.78	0.48 to 0.92	0.68	0.67
34444144	0.36 (0.46)	0.50	1.78	0.03 to 0.65	0.40	0.34

State	Observed				Modelled	
	Mean (SD)	Median	Range	Inter-quartile range	Model (3)	Model (6)
41111111	0.94 (0.15)	1.00	0.75	0.93 to 1.00	0.90	0.87
41134422	0.45 (0.59)	0.68	1.98	0.23 to 0.93	0.58	0.51
41213241	0.69 (0.37)	0.83	1.65	0.50 to 0.98	0.72	0.65
41234423	0.25 (0.59)	0.43	1.90	-0.36 to 0.70	0.55	0.41
41321114	0.64 (0.29)	0.68	1.00	0.38 to 0.88	0.75	0.71
41343342	0.50 (0.44)	0.55	1.58	0.25 to 0.94	0.59	0.48
42131332	0.77 (0.21)	0.80	0.83	0.66 to 0.93	0.68	0.68
42211444	0.52 (0.37)	0.53	1.50	0.23 to 0.83	0.63	0.56
42412133	0.65 (0.31)	0.68	0.98	0.43 to 1.00	0.63	0.67
42423213	0.57 (0.45)	0.60	1.68	0.43 to 0.95	0.55	0.59
43112413	0.74 (0.28)	0.75	0.98	0.53 to 1.00	0.70	0.66
43142121	0.65 (0.29)	0.68	1.25	0.48 to 0.88	0.73	0.63
43243211	0.41 (0.50)	0.50	1.98	0.28 to 0.73	0.66	0.51
43411232	0.60 (0.45)	0.70	1.93	0.48 to 0.93	0.62	0.67
44124131	0.67 (0.27)	0.63	1.23	0.51 to 0.89	0.60	0.58
44321321	0.56 (0.42)	0.64	1.88	0.33 to 0.96	0.66	0.62
44332314	0.42 (0.56)	0.58	1.93	0.35 to 0.83	0.58	0.51
44444144	0.36 (0.47)	0.44	1.73	0.06 to 0.79	0.38	0.29
51111111	0.82 (0.35)	0.95	1.48	0.81 to 1.00	0.90	0.87
51224434	0.41 (0.42)	0.49	1.97	0.21 to 0.69	0.51	0.43
51442332	0.56 (0.25)	0.53	0.98	0.40 to 0.73	0.55	0.54
52133121	0.77 (0.24)	0.80	0.95	0.64 to 0.98	0.72	0.69
52311443	0.65 (0.33)	0.76	1.20	0.42 to 0.94	0.64	0.58
53242314	0.53 (0.32)	0.58	1.13	0.26 to 0.78	0.61	0.47
53424212	0.58 (0.43)	0.70	1.73	0.43 to 0.93	0.47	0.51
54133141	0.57 (0.40)	0.68	1.33	0.38 to 1.00	0.64	0.61
54311223	0.69 (0.23)	0.68	0.78	0.50 to 0.95	0.70	0.61
54444444	0.13 (0.51)	0.13	2.00	-0.23 to 0.50	0.29	0.20

Table 7

Estimated preference weights

Dimensions and levels	(1) OLS	(2) MLE RE	(3) ERUM OLS	(4) Mean model		(5) Consistent mean model
PF2	0.061	0.052	0.052	0.065	PF2	0.065
PF3	0.076	0.079	0.077	0.078	PF3	0.078
PF4	0.135	0.134	0.103	0.139	PF45	0.127
PF5	0.121	0.127	0.104	0.105		
RF2	0.026	0.023	0.044	0.032	RF2	0.032
RF3	0.042	0.052	0.050	0.045	RF3	0.045
RF4	0.082	0.090	0.076	0.079	RF4	0.078
PAIN2	0.059	0.041	0.054	0.059	PAIN2	0.059
PAIN3	0.060	0.060	0.064	0.062	PAIN3	0.062
PAIN4	0.070	0.083	0.070	0.065	PAIN4	0.064
EF2	0.028	0.027	0.032	0.030	EF2	0.030
EF3	0.063	0.072	0.053	0.066	EF3	0.066
EF4	0.157	0.160	0.132	0.150	EF4	0.149
SF2	0.025	0.022	0.029	0.027	SF2	0.027
SF3	0.059	0.065	0.046	0.059	SF3	0.059
SF4	0.173	0.174	0.132	0.163	SF4	0.163
FAT2	0.046	0.026	0.038	0.046	FAT2	0.047
FAT3	0.052	0.031	0.052	0.054	FAT3	0.054
FAT4	0.104	0.064	0.084	0.093	FAT4	0.092
NAU2	0.031	0.036	0.025	0.032	NAU23	0.026
NAU3	0.015	0.037	0.027	0.019		
NAU4	0.062	0.079	0.052	0.057	NAU4	0.056
CD2	0.012	0.022	0.011	0.016	CD2	0.016
CD3	0.050	0.037	0.035	0.052	CD3	0.052
CD4	0.078	0.070	0.059	0.073	CD4	0.072
Observations	2710	2710	2710	85		85
R-squared	0.60		0.56	0.97		0.97
Number of id		344				
Inconsistencies	2	1	0	2		0

Dimensions and levels	(1) OLS	(2) MLE RE	(3) ERUM OLS	(4) Mean model	(5) Consistent mean model
Insignificant level coefficients	5	5	3	6	5
MAE	0.052	0.054	0.046	0.050	0.051
MAE>0.05	41	41	33	37	39
MAE>-0.10	9	13	6	8	9

Note: Figures in bold have t-statistics significant at the 5% level

PF=Physical functioning, RF=Role functioning, PAIN=Pain, EF=Emotional functioning, SF=Social functioning, FAT=Fatigue, NAU=Nausea, CD=Constipation and/or diarrhoea. PF23=Physical functioning at level 2 or 3, NAU23=Nausea at level 2 or 3.