



ELSEVIER

Available online at www.sciencedirect.com

ScienceDirect

journal homepage: www.elsevier.com/locate/jval

Preference-Based Assessments

Learning and Satisficing: An Analysis of Sequence Effects in Health Valuation



Benjamin M. Craig, PhD^{1,*}, Shannon K. Runge, MA¹, Kim Rand-Hendriksen, PhD²,
Juan Manuel Ramos-Goñi, BSc³, Mark Oppe, PhD⁴

¹Health Outcomes and Behavior, Moffitt Cancer Center and University of South Florida, Tampa, FL, USA; ²Department of Health Management and Health Economics, University of Oslo, Oslo, Norway; ³Canary Islands Health Service (SESCS), Santa Cruz de Tenerife, Spain; ⁴Institute for Medical Technology Assessment, Erasmus University Rotterdam, Rotterdam, The Netherlands

ABSTRACT

Objective: To estimate the effect of sequence on response precision and response behavior in health valuation studies. **Methods:** Time trade-off (TTO) and paired comparison responses from six health valuation studies—four US, one Spanish, and one Dutch—were examined (22,225 respondents) to test whether task sequence influences response precision (e.g., rounding), response changes, and median response times. Each study used a computer-based instrument that randomized task sequence among a national sample of adults, age 18 years or older, from the general population. **Results:** For both TTO and paired comparisons, median response times decreased with sequence (i.e., learning), but tended to flatten after the first three tasks. Although the paired comparison evidence demonstrated that sequence had no effect on

response precision, the frequency of rounded TTO responses (to either 1-year or 5-year units) increased with sequence. **Conclusions:** Based on these results, randomizing or reducing the number of paired comparison tasks does not appear to influence response precision; however, generalizability, practicality, and precautionary considerations remain. Overall, participants learned to respond efficiently within the first three tasks and did not resort to satisficing, but may have rounded their TTO responses. **Keywords:** health valuation, paradata, preferences, QALY, response precision, sequence effects, time trade-off.

Copyright © 2015, International Society for Pharmacoeconomics and Outcomes Research (ISPOR). Published by Elsevier Inc.

Introduction

Most economic evaluations summarize effectiveness using preference weights on a quality-adjusted life-year (QALY) scale, as recommended by numerous health technology assessment agencies. Such QALY weights may be from societal or patient perspectives and derived using a wealth of preference elicitation tasks (e.g., best-worst scaling). Although valuation research has a well-established history, the use of online computer-based surveys for health valuation offers an array of new capabilities, such as quota-sampling at the task level; paradata on respondent behavior, device, and browser; and other interactive technologies. Compared with interview, postal, or telephone surveys, online computer-based experiments increase control in the randomization of tasks, while reducing cognitive burden and minimizing missing data and other data collection errors and biases.

Although online instruments typically randomize the order of presentation of tasks, response precision and behavior may change with sequence. For example, when a respondent is shown two alternatives and asked, “Which do you prefer?” he or she may take longer or change his or her responses on initial pairs

while becoming acquainted with the valuation task as compared with later pairs. Furthermore, a respondent’s attention may wane in later pairs, leading to satisficing (i.e., expediting selection among alternatives to minimize effort), reducing response precision [1,2]. This article examines whether response precision and response behavior vary with the number of tasks completed (i.e., sequence effect) in health valuation studies for two types of valuation tasks, time trade-off (TTO) and paired comparisons.

Understanding the relationship between response precision and task sequence guides the number of tasks to be included in a valuation study, informs weights that place a greater emphasis on earlier or later tasks, and justifies the randomization of task sequence. Although studies have attempted to identify respondents who randomize all responses (i.e., shufflers and satisficers) [3], few studies to date have examined the effect of sequence on response precision in health valuation [4].

Sequence effects have been identified in other forms of discrete choice experiments (DCEs) as a type of ordering effect specifically related to the order in which choice sets are presented (i.e., position-dependent order effects) [5]. This type of order effect differs from those related to the order or position of

* Address correspondence to: Benjamin M. Craig, Moffitt Cancer Center, 12902 Magnolia Drive, MRC-CANCONT, Tampa, FL 33612.

E-mail: benjamin.craig@moffitt.org.

1098-3015/\$36.00 – see front matter Copyright © 2015, International Society for Pharmacoeconomics and Outcomes Research (ISPOR).

Published by Elsevier Inc.

<http://dx.doi.org/10.1016/j.jval.2014.11.005>

attributes within a choice set [5–7]. Experimental design, such as the layout of questions, the number of attributes, and the number of tasks, can influence ordering effects and response time [8–10]. A key example in survey research is the primacy effect or the tendency for respondents to choose the first reasonable answer to a survey question (e.g., first response option in a list of potential answers) [6,11]. This weak form of satisficing leads to nonrandom response; expedites response with minimum effort; reduces response quality and time; and is commonly cited by experimenters to justify randomization and reduction in the number of attributes, scenarios, and tasks [12].

A wealth of studies have examined order effects in terms of perception and salience [5,7,9,10,13–17], although the results have been somewhat inconsistent. For example, some evidence suggests that the order of attributes affects choice [5,7], yet other studies did not find this effect [9,14,18]. In addition, the number and complexity of task sets within an experiment may induce order effects through respondent fatigue or boredom [19]. Evaluating the association between participant response behaviors (i.e., response times and changes) and task sequence has the potential to provide valuable insight regarding the influence of study design.

In complement to evidence on response precision, we examine response behaviors (i.e., response times and changes) that may indicate learning and added deliberative effort beyond that which is needed to satisfy the task requirements. Typically, response behavior is examined at the questionnaire level (e.g., the amount of time it takes a respondent to complete all tasks). In addition to evaluating response behavior at the questionnaire level, computerized software offers a unique opportunity to examine response behaviors at the level of individual questions (e.g., the amount of time it takes to complete a single task set or a series of different task sets). A better understanding of response behavior at each of these levels can aid in the interpretation of the empirical association between sequence and response precision and in the improvement of survey design (e.g., cognitive burden).

The present study contributes to an innovative evaluation of client-side paradata. Client-side paradata is the information recorded in Web surveys by the respondent's computer (e.g., the number of times and locations of mouse clicks on a computer screen). Unlike server-side paradata, which refers to data management processes, client-side information allows researchers to interpret participant response behaviors in terms of changed responses (CRs) and response time at the level of

individual questions [20]. Evaluating response behavior patterns at such a specific level contributes to our knowledge of how sequence influences preferences. In this secondary analysis of health valuation data, we examine sequence effects, specifically whether response precision and response behavior vary with the number of tasks completed.

Methods

Preference Elicitation

In a paired comparison, respondents are asked, “Which do you prefer?” given two health episodes, and their choices define the relative value between these episodes. An original TTO task is more involved, using an adaptive series of paired comparisons based on either time with no health problems or “immediate death.” Specifically, each TTO begins with a paired comparison in which the respondent must first decide whether the health episode is preferred to immediate death. If so, an adaptive series of paired comparisons is presented to determine the number of years with no health problems that is equivalent to the health episode (i.e., better-than-dead indifference statement). If the respondent prefers immediate death, an alternative series of paired comparisons is completed to identify a worse-than-dead indifferent statement. The original adaptation procedure [21–23] is like a dose-response study in that it increases the duration of problems within an episode until it is equivalent to immediate death (e.g., how much poison is needed until it kills you). Thus, the TTO exercise is a matching task that produces an equivalence statement regardless of whether the original paired comparison response is better or worse than death.

Data

To test the effect of sequence on response precision and behavior, we examined paired comparisons and TTO responses from six health valuation studies—four US, one Spanish, and one Dutch—totaling 259,318 responses from 22,225 respondents who completed 17 to 37 tasks [2,24–27]. Table 1 summarizes the characteristics of these six studies. All studies used a computerized instrument that randomized task sequence using national samples of adults from the general population. For the US-based studies, respondents completed a set of paired comparisons trading improvements in health-related quality of life (HRQOL)

Table 1 – Health valuation studies*.

Study title	Dates	No.	First set of tasks	Second set of tasks
Patient Reported Outcomes Measurement Information System (PROMIS) Valuation Study - United States [2]	March–July 2012	7557	6 lifespan pairs	24 health pairs
EQ-5D-5L Valuation Study - Spanish	May–July 2012	986	10 time trade-offs	7 health state pairs [†]
Child Health Valuation Study - US, Wave 1 [24]	July–August 2012	2008	6 lifespan pairs	31 health pairs
EQ-5D-5L Valuation Study - Dutch [27]	September–October 2012	1052	10 time trade-offs	7 health state pairs [†]
Child Health Valuation Study - United States, Wave 2 [24]	January–February 2013	2147	12 lifespan pairs	18 health pairs
Women's Health Valuation Study - United States [25]	April 2013	3397	8 lifespan pairs	22 health pairs
Measurement and Valuation of Health Study - United States [26]	November–December 2013	5078	8 lifespan pairs	22 health pairs

EQ-5D-5L, five-level EuroQol five-dimensional questionnaire.
* Each wave of the US Child Health Valuation Study is shown separately because of changes in the valuation tasks.
[†] Unlike health and lifespan pairs, health state pairs do not describe duration in the health state.

for reduced lifespan (i.e., lifespan pairs) before completing a second set that traded alternative HRQOL scenarios of a common duration (i.e., health pairs). For the valuation of the five-level EuroQol five-dimensional questionnaire, respondents completed a set of TTOs before completing a set of paired comparisons that traded alternative HRQOL scenarios without a description of duration (i.e., health state pairs). Further description of the protocol of each study is provided online [2,24–27].

The TTO task in the Spanish and Dutch studies was an adaptive hierarchy of steps known as the composite TTO (Fig. 1) [27]. The composite TTO is derived from both the original and lead-time TTO [21–23]. Each step displayed two scenarios, and the respondent was asked, “Which is better?” If the respondent did not wish to choose, the respondent may instead state indifference (i.e., the scenarios were “about the same”).

In this adaptive process, the task began with the choice between 10 years in full health and 10 years in the health state (i.e., step 1). If the respondent chose the health state scenario or stated indifference, the TTO response was +10 and the task ended. If the respondent chose the full health scenario in step 1, the task continue on to step 2 and displayed 0 years in full health (i.e., immediate death) instead of 10 years in full health.

If the respondent chose the full health scenario in step 2, the task continued to step 3 and displayed 5 years in full health instead of 0 years in full health. If the respondent chose the health state scenario in step 2, the task continued to step 3 and displayed –5 years in full health instead of 0 years in full health. If the respondent stated indifference in step 2, the TTO response was 0 and the task ended. This task continued for up to nine steps until the respondent expressed indifference between the two scenarios (Fig. 1).

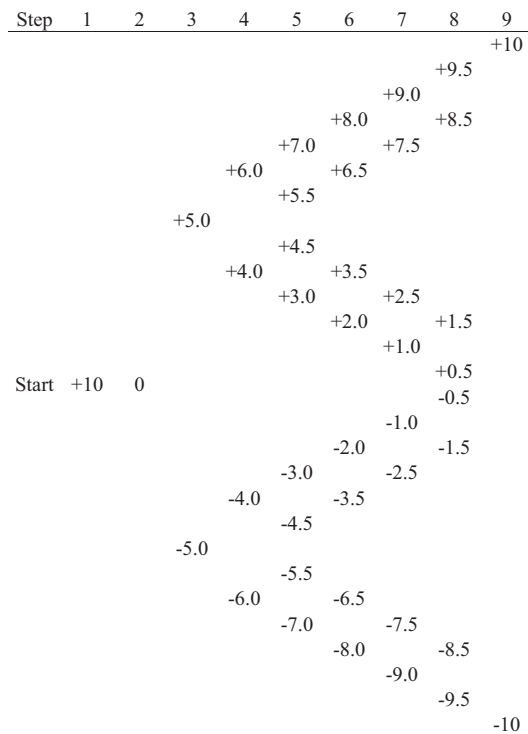


Fig. 1 – Minimum number of steps involved in each composite time trade-off response. Numbers in the time trade-off represent the value of 10 years in health state on a quality-adjusted life-year (QALY) scale based on a statement of indifferences (e.g., 10 years in health state = +5 QALYs).

Aside from the highest possible response (+10), which required either one or nine steps, each TTO response required a minimum number of steps (i.e., some TTO responses required more steps than did others). The lowest possible response (–10) required the most effort (i.e., nine steps). By construction, about half of any TTO sample should have been in half-year units.

Paired comparison tasks differed by the studies. The US-based paired comparisons began with three examples and asked “Which do you prefer?” showing two health scenarios with only two attributes and their durations. The Spanish and Dutch paired comparisons had no examples. Respondents completed between 7 and 37 paired comparisons. Unlike the TTO task, indifference was not allowed in any of these paired comparison tasks.

Econometrics

For each study, we graphed the median response time and the relative risk of a CR and a modal response (MR) by sequence. Response time was measured in seconds from the time that the task was first shown until the final response to the task.

A CR is when multiple responses were registered in the paradata for the task (e.g., a respondent may choose the first scenario in a paired comparison as the preferred scenario and then change his or her response to the second scenario). Changing a response may be related to the difficulty of the choice. For example, if two scenarios seemed similar, the probability of changing a response is higher than for a pair with dissimilar scenarios. Nevertheless, we hypothesize that sequence is unrelated to CRs when pairs are randomly sequenced. Specifically, the relative risk of a CR is the risk of a CR at the location in the sequence divided by the overall risk of a CR. We did, however, investigate the impact of the difficulty of the choice on the CR in a sensitivity analysis.

To identify half-year unit responses in the TTO tasks, respondents may be required to complete additional steps to achieve the final response. These steps include overshooting the point of indifference by half a year and backtracking half a year. For example, for a respondent to achieve a final TTO of 6.5, he or she would first be presented with additional scenarios comparing 7 years in a health state (overshooting) and 6 years in a health state (backtracking). Therefore, a TTO CR requires added steps and responses, and a DCE CR implies just added responses. In either case, we hypothesize that sequence is unrelated to the relative risk of CR.

An MR is whether the respondent provided the same response as the modal response for the task. For example, in a choice between mild pain and mild depression, 80% may choose mild pain and this MR should not vary by sequence. If respondent attention waned, however, the frequency of MRs should diminish until just 50% prefer mild pain. Specifically, the relative risk of an MR is the risk of an MR at the location in the sequence divided by the overall risk of an MR.

For a TTO task, the responses are not binary but are integer and half-integer values on a scale ranging from +10 to –10. Therefore, the risk of a TTO MR may be lower than a risk of a DCE MR. In either case, we hypothesize that sequence is unrelated to the relative risk of MR, the relative risk of CR, or median response times.

As ancillary measures of TTO response precision, we illustrated the frequency of 5-year and half-year unit TTO responses by sequence. A half-year unit response requires that the respondent complete at least one more step than a 1-year unit response. The frequency of half-unit responses represented a trade-off between added effort and greater precision, which may have varied by sequence. Likewise, a respondent may have stopped the task early (i.e., within three steps: +10, 0, +5 or –5) and responded in 5-year units. Rounding to 1-year or to 5-year units was a tacit way to avoid added effort in the TTO task (i.e., satisficing).

All analyses were repeated using varying levels of difficulty (i.e., comparing different levels of severe health states) on the basis of the assumption that decision difficulty increases as respondents compare health scenarios with similar levels of severity. For the TTO tasks, decision difficulty is assumed to peak at the point of respondent indifference between health scenarios. For the DCE tasks, this point occurs when the choice probability of two health scenarios is approximately 50%. Subsequently, we used posterior information about DCE pair probabilities to describe subgroups.

Results

Figure 2 illustrates median response times by sequence. At the beginning of each sequence, response times were reduced substantially. Each line exhibits the same downward sloping shape (i.e., learning) and shows a flattening out. Dutch respondents had a higher median time than did Spanish and US respondents, regardless of task. Spanish paired comparisons had a higher response time than US tasks, possibly due to differences in the number of attributes of each alternative (5 vs. 2). This pattern was also observed in the subgroup analysis, which confirmed that more time was needed when the task was more difficult. We examined, however, whether sequence effects (i.e., median response times, CR, MR, and rounding) were similar among tasks with different levels of difficulty (e.g., greater effect seen in easier tasks) and found no differences. Figure 3 illustrates the relative risk of CR by sequence, which decreases over the initial tasks. Figure 4 illustrates the relative risk of MR by sequence, and the MR lines appear flat (i.e., relative risks range from 1.1 to 0.9) aside from some wavering.

Unlike the paired comparison responses, TTO responses may be rounded to 1-year or 5-year units, possibly to reduce response effort (Fig. 1). Figure 5 illustrates the frequency of 5-year, 1-year, and half-year unit TTO responses. The results show that more than 40% of the Spanish TTO responses were either +10, +5, 0, or -5, regardless of sequence, and that the frequency of these 5-year unit responses increased from 30% to 40% in the Dutch data, representing a reduction in TTO response precision with sequence. Half-year unit responses potentially indicated a small gain in precision and should be half of each sample. The frequencies of half-year unit responses were clearly less than 50% and decreased from 19% to 12% and from 14% to 12% in the Dutch and Spanish samples, respectively. Furthermore, all 86 modal TTO responses in the Dutch and Spanish

studies were in 1-year units and most (77% and 87%, respectively) were in 5-year units. It should be noted, however, that even though the proportion of 5-year values and 1-year values is large across respondents, only a small number of respondents give only 5-year values (2% and 36% in the Dutch study and 6% and 47% in the Spanish study).

Discussion

Using data from 22,225 respondents, we found that sequence had no effect on paired comparison response precision, but may induce greater rounding in TTO responses. The CR lines (Fig. 3) decrease over the initial tasks, illustrating that those respondents may be learning the task or establishing heuristics that govern their responses of all similar tasks. The first six tasks for each US study were lifespan pairs that involved the trade-off between reduced lifespan and HRQOL. This emphasis on a single attribute (i.e., lifespan) may have induced the formation of time-specific heuristics compared with latter pairs that traded two losses in HRQOL with common duration. Aside from some wavering, the MR lines (Fig. 4) appear flat (i.e., relative risks range from 1.1 to 0.9), illustrating that response precision was not associated with sequence. The greater variability seen in the TTO MR is likely attributable to its use of nonbinary responses.

With TTO, it can be argued that the proportion of half-year responses, theoretically, should be similar to integer-year responses (1- or 5-year units), given the assumption that the distribution of preferences could be considered continuous. The results show that the proportion of half-year responses was less than half and decreased with sequence, although at a different rate in the Spanish data than in the Dutch data. Nevertheless, such TTO rounding had no effect on the relative risk of a modal TTO response. This absence of effect may be attributed to the fact that most modal TTO responses are in 5-year units (i.e., rounding increases the likelihood of MR).

The low and falling proportion of half-year unit TTO responses is striking, but the correct interpretation is not straightforward. The procedure used to identify a half-year unit response requires overshooting the point of indifference and backtracking half a year. For example, a respondent who has a TTO value of 6.5 for a health scenario would be offered 10 years of perfect health, followed by 0, 5, 6, and 7 (overshoot) before stating indifference at 6.5 years. Similarly, a respondent who has a TTO value of 3.5 for a health scenario is offered 10 years of full health, followed by 0, 5, 4, 3 (overshoot) before stating indifference at 3.5 years. The reduction in elicited half-year unit response could

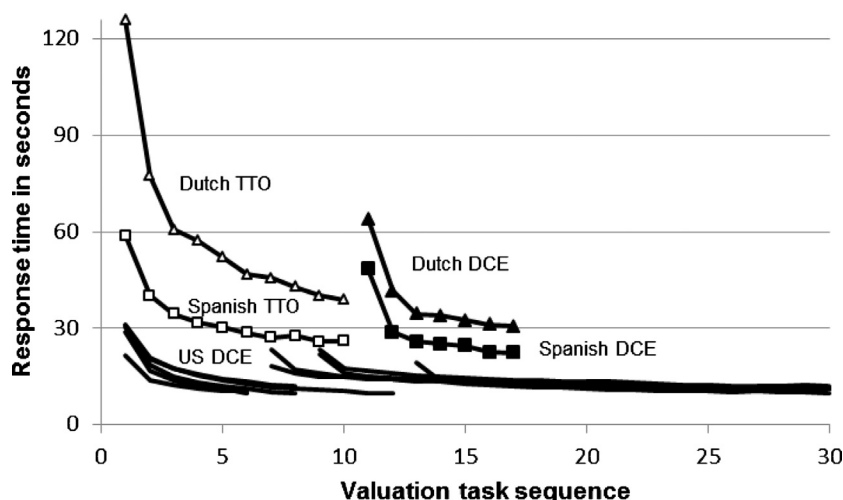


Fig. 2 – Median response time by sequence. DCE, discrete choice experiment; TTO, time trade-off.

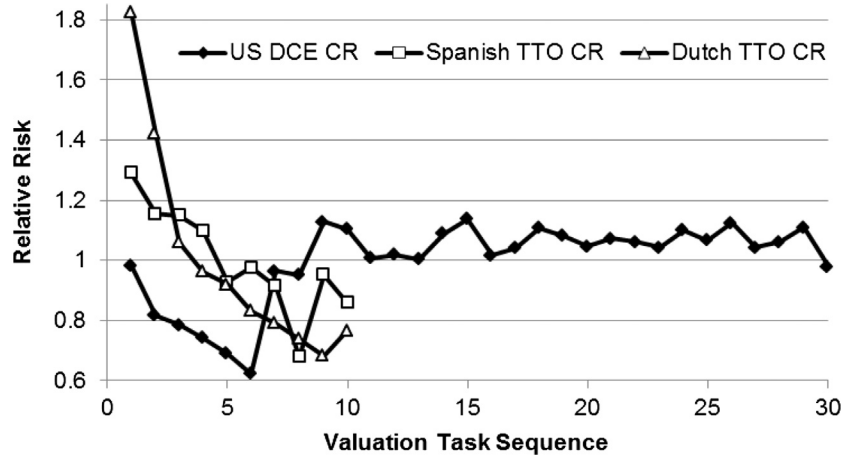


Fig. 3 – Relative risk of changed responses (CRs) by sequence. The Dutch and Spanish studies did not collect CR data on paired comparisons. DCE, discrete choice experiment; TTO, time trade-off.

represent satisficing, but it could also reflect a reluctance to backtrack, reluctance to overshoot, or a genuine satisfaction with the level of precision offered by sticking to whole years. Which of these explanations is at play could possibly be determined through strategic manipulation of the routing, such as removing the half-year correction, altering the step size to half a year, or giving respondents multiple alternatives (i.e., more than two scenarios in a choice set) at each step. Regardless of explanation, the results show that sequence influences the frequency of half-year responses; however, the infrequency of half-year responses suggests that the potential loss of information is limited.

The apparent and increasing frequency of 5-year unit responses is more troubling because the loss of information is large. The results suggest that most of the respondents are attracted to these 5-year unit responses, increasing the risk of bias. The extent of these primacy effects and their attraction may be caused by digit preference, satisficing, or cognitive biases, such as anchoring, and should be investigated further. Based on the paired comparison results, randomizing or reducing the number of paired comparison tasks does not appear to influence response precision; however, generalizability, practicality, and precautionary considerations remain.

These considerations are largely related to the design of DCEs: What is the optimal number of tasks that should be included in a survey? Should later tasks be downweighted? Should tasks be randomized? It has been proposed that certain variations in

survey design (e.g., increases in the number of tasks, scenarios, and attributes) increase respondent burden and fatigue, thus contributing to ordering effects and response variability [10,19,28]. Despite a growing interest in identifying the optimal design for DCEs, the existing literature remains inconclusive and the results of this study failed to identify any benefits from decreasing the number of tasks, downweighting later tasks, or randomizing tasks.

Shortening a health preference survey may limit the breadth of the results (e.g., too few attributes) and collect insufficient data to calculate preferences on attributes, particularly if sample size is small [18,29]. In their widely cited article, Hensher et al. [18] found 4 and 8 tasks to be insufficient to estimate preferences for attributes that were selected less often but concluded that this could be remedied by presenting 24 to 32 tasks without overburdening respondents. Similarly, Carlsson and Martinsson [29] compared the results of 12 and 24 tasks and found no evidence of sequence effects, but they did report a significantly higher dropout rate for the longer survey. The results of these studies, however, contradict other findings. In a valuation of travel time, Hensher [28] reported that increasing the number of tasks significantly decreased participant response time and significantly affected the outcome of the study. These results were echoed by Chung et al. [30], who concluded the ideal number of tasks to be six per survey. Although it has been noted that researchers should use careful pretesting to identify the optimal

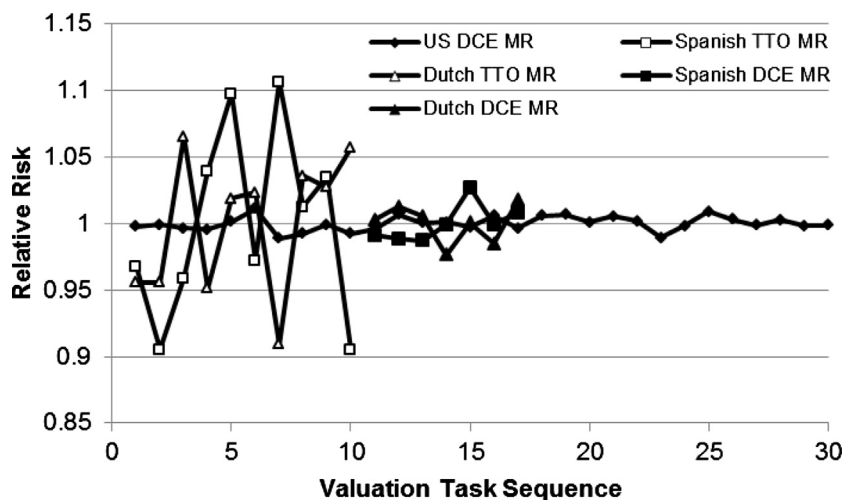


Fig. 4 – Relative risk of modal response (MR) by sequence. DCE, discrete choice experiment; TTO, time trade-off.

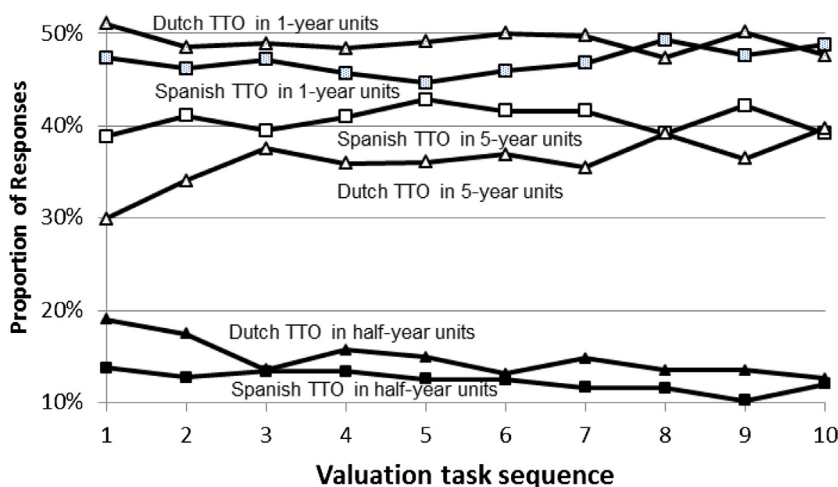


Fig. 5 – TTO rounding by sequence. TTO, time trade-off.

number of tasks to include in a DCE [30], our results did not find any sequence effects in the DCE, possibly due to their simplicity (two alternatives with two attributes). Still, additional research is needed to rectify these discrepancies.

The primary limitation of this study is that each study included a maximum of 37 tasks because these components were designed to be completed in less than 30 minutes. Evidence, however, from health valuation studies with more than 40 tasks will be explored in future work. In fact, Craig et al. are currently in the beginning stages of a study that will allow respondents to complete hundreds of pairs. Our sensitivity analyses on the time it takes to complete a task by difficulty indicated, however, that the time needed to answer a task is shorter for easy tasks than for difficult tasks. This should be taken into account in the design of a study.

Another limitation of the present study is that trends in the relative risk of MR may underrepresent losses in TTO precision due to rounding, because most TTO MR are in 5-year units. The use of MR allowed for a uniform summary of trends in TTO and paired comparison response precision, but did not compensate rounding. The proportion of 5-year units and 1-year units is quite large across respondents. Only a few respondents, however, use only 5-year responses or 1-year responses. Future studies may investigate whether rounding is a greater concern in subgroups of respondents, particularly those with low numeracy. The conclusion from this analysis is that sequence effects are present more in TTOs than in DCEs, but both show some learning effect. In summary, the results of this study failed to identify any benefits from decreasing the number of DCE tasks, downweighting later DCE tasks, or randomizing DCE tasks.

Acknowledgments

We thank Carol Templeton and Michelle Owens at Moffitt Cancer Center for their contributions to the research and creation of this article.

Source of financial support: Funding support for this research was provided by a National Cancer Institute R01 grant (1R01CA160104), the EuroQol Group (EQ Project 2013130), and Dr. Craig's support account at Moffitt Cancer Center.

REFERENCES

- [1] Barge S, Gehlbach H. Using the theory of satisficing to evaluate the quality of survey data. *Res High Educ* 2012;53:182–200.
- [2] Craig B, Reeve BB. Methods Report on the Promis Valuation Study: Year 1. Available from: <http://labpages.moffitt.org/craig/Publications/Report120928.pdf>. [Accessed October 11, 2012].
- [3] Craig BM, Ramachandran S. Relative risk of a shuffled deck: a generalizable logical consistency criterion for sample selection in health state valuation studies. *Health Econ* 2006;15:835–48.
- [4] Augestad LA, Rand-Hendriksen K, Kristiansen IS, et al. Learning effects in time trade-off based valuation of EQ-5D health states. *Value Health* 2012;15:340–5.
- [5] Day B, Bateman JJ, Carson RT, et al. Ordering effects and choice set awareness in repeat-response stated preference studies. *J Environ Econ Manage* 2012;63:73–91.
- [6] Malhotra N. Completion time and response order effects in web surveys. *Public Opin Q* 2008;72:914–34.
- [7] Kjaer T, Bech M, Gyrd-Hansen D, et al. Ordering effect and price sensitivity in discrete choice experiments: need we worry? *Health Econ* 2006;15:1217–28.
- [8] Christian LM, Parsons NL, Dillman DA. Designing scalar questions for web surveys. *Sociol Methods Res* 2009;37:393–425.
- [9] Farrar S, Ryan M. Response-ordering effects: a methodological issue in conjoint analysis. *Health Economics* 1999;8:75–9.
- [10] Savage SJ, Waldman DM. Learning and fatigue during choice experiments: a comparison of online and mail survey modes. *J Appl Econ* 2008;23:351–71.
- [11] Krosnick JA. Response strategies for coping with the cognitive demands of attitude measures in surveys. *Appl Cogn Psychol* 1991;5:213–36.
- [12] Schwarz N, Sudman S, Schuman H. *Context Effects in Social and Psychological Research*. New York, NY: Springer-Verlag, 1992.
- [13] Blumenschein K, Johannesson M. An experimental test of question framing in health state utility assessment. *Health Policy* 1998;45:187–93.
- [14] Boyle KJ, Ozdemir S. Convergent validity of attribute-based, choice questions in stated-preference studies. *Environ Resour Econ* 2009;42:247–64.
- [15] Howard K, Salkeld G. Does attribute framing in discrete choice experiments influence willingness to pay? Results from a discrete choice experiment in screening for colorectal cancer. *Value Health* 2009;12:354–63.
- [16] Kamoen N, Holleman B, Mak P, et al. Agree or disagree? Cognitive processes in answering contrastive survey questions. *Discl Process* 2011;48:355–85.
- [17] Yan T, Tourangeau R. Fast times and easy questions: the effects of age, experience and question complexity on web survey response times. *Appl Cogn Psychol* 2008;22:51–68.
- [18] Hensher DA, Stophor PR, Louviere JJ. An exploratory analysis of the effect of numbers of choice sets in designed choice experiments: an airline choice application. *J Air Trans Manag* 2001;7:373–9.
- [19] De Palma A, Myers GM, Papageorgiou YY. Rational choice under an imperfect ability to choose. *Am Econ Rev* 1994;84:419–40.
- [20] Heerwegh D. Explaining response latencies and changing answers using client-side paradata from a web survey. *Soc Sci Comp Rev* 2003;21:360–73.
- [21] Devlin NJ, Tsuchiya A, Buckingham K, et al. A uniform time trade off method for states better and worse than dead: feasibility study of the 'Lead Time' approach. *Health Econ* 2011;20:348–61.

- [22] Gudex C. Time Trade-Off User Manual: Props and Self-Completion Methods. Report of the Centre for Health Economics. York, United Kingdom: University of York, 1994.
- [23] Torrance GW, Thomas WH, Sackett DL. A utility maximization model for evaluation of health care programs. *Health Serv Res* 1972;7:118–33.
- [24] Craig B, Owens MA. Methods Report on the Child Health Valuation Study (CHV): Year 1. Available from: http://labpages.moffitt.org/craigb/Publications/CHVMethods_130917.pdf. [Accessed November 5, 2013].
- [25] Craig B, Owens MA. Methods Report on the Women's Health Valuation Study (WHV): Year 1. Available from: http://labpages.moffitt.org/craigb/Publications/WHVMethods_140106.pdf. [Accessed January 10, 2014].
- [26] Craig B, Owens MA. 2013 United States Measurement and Valuation of Health Study (2013 US MVH): Methods Report. Available from: http://labpages.moffitt.org/craigb/Publications/MVHMethods_140116.pdf. [Accessed January 27, 2014].
- [27] Janssen BMF, Oppe M, Versteegh MRN, et al. Introducing the composite time trade-off: a test of feasibility and face validity. *Eur J Health Econ* 2013;14(Suppl):S5–13.
- [28] Hensher DA. Revealing differences in willingness to pay due to the dimensionality of stated choice designs: an initial assessment. *Environ Resour Econ* 2006;34:7–44.
- [29] Carlsson F, Martinsson P. How much is too much? *Environ Resour Econ* 2008;40:165–76.
- [30] Chung C, Boyer T, Han S. How many choice sets and alternatives are optimal? Consistency in choice experiments. *Agribusiness* 2011;27:114–25.