

Medical Decision Making

<http://mdm.sagepub.com>

Health Utility Bias: A Systematic Review and Meta-Analytic Evaluation

Jason N. Doctor, Han Bleichrodt and H. Jill Lin

Med Decis Making 2010; 30; 58 originally published online Jun 12, 2008;

DOI: 10.1177/0272989X07312478

The online version of this article can be found at:
<http://mdm.sagepub.com/cgi/content/abstract/30/1/58>

Published by:



<http://www.sagepublications.com>

On behalf of:



<http://www.smdm.org>
Society for Medical Decision Making

Additional services and information for *Medical Decision Making* can be found at:

Email Alerts: <http://mdm.sagepub.com/cgi/alerts>

Subscriptions: <http://mdm.sagepub.com/subscriptions>

Reprints: <http://www.sagepub.com/journalsReprints.nav>

Permissions: <http://www.sagepub.com/journalsPermissions.nav>

Citations <http://mdm.sagepub.com/cgi/content/refs/30/1/58>

Health Utility Bias: A Systematic Review and Meta-Analytic Evaluation

Jason N. Doctor, Han Bleichrodt, H. Jill Lin

Background. A common assertion is that rating scale (RS) values are lower than both standard gamble (SG) and time tradeoff (TTO) values. However, differences among these methods may be due to method specific bias. Although SG and TTOS suffer systematic bias, RS responses are known to depend on the range and frequency of other health states being evaluated. Over many diverse studies this effect is predicted to diminish. Thus, a systematic review and data synthesis of RS-TTO and RS-SG difference scores may better reveal persistent dissimilarities. **Purpose.** The purpose of this study was to establish through systematic review and meta-analysis the net effect of biases that endure over many studies of utilities. **Methods.** A total of 2206 RS and TTO and 1318 RS and SG respondents in 27 studies of utilities participated. MEDLINE was searched for data from 1976 to 2004, complemented by a hand search of full-length articles and conference abstracts for 9 journals known to publish utility studies, as well as review of results and additional recommendations by 5 outside experts in the field. Two

investigators abstracted the articles. We contacted the investigators of the original if required information was not available. **Results.** No significant effect for RS and TTO difference scores was observed: effect size (95% confidence interval [CI]) = 0.04 (−0.02, 0.09). In contrast, RS scores were significantly lower than SG scores: effect size (95% CI) = −0.23 (−0.28, −0.19). Correcting SG scores for 3 known biases (loss aversion, framing, and probability weighting) eliminated differences between RS and SG scores: effect size (95% CI) = 0.01 (−0.03, 0.05). Systematic bias in the RS method may exist but be heretofore unknown. Bias correction formulas were applied to mean not individual utilities. **Conclusions.** The results of this study do not support the common view that RS values are lower than TTO values, may suggest that TTO biases largely cancel, and support the validity of formulas for correcting SG bias. **Key words:** utility measurement; rating scale; category scale; time tradeoff; standard gamble. (*Med Decis Making* 2010;30:58–67)

The purpose of this study is to establish through systematic review and meta-analysis the net effect of health utility biases that occur under different elicitation methods. Health utilities play an important role in cost-effectiveness analysis. Through health utility assessment, to each health state in the

analysis a presumably unique quality weight is assigned. The standard gamble (SG), time tradeoff (TTO), and rating scale (RS) are the most common preference assessment methods for assigning such weights. However, when more than 1 elicitation method is used it is often the case that more than 1 quality weight may be assigned to any particular health state.^{1,2} One negative implication of this is that treatment recommendations may be sensitive to the method of preference assessment.³ Differences among health state valuation methods may be due to biases that lead to errors in measurement and result in health state utilities that are too high or too low. By seeking to understand the net effect of bias, we may be in a better position to recommend certain methods that minimize the occurrence of errors.

Received 25 May 2006 from Department of Pharmaceutical Economics and Policy, School of Pharmacy, University of Southern California, Los Angeles, CA (JND); Department of Economics and iMTA/iBMG, Erasmus University, Rotterdam, The Netherlands (HB); and Department of Radiology, School of Medicine, Stanford University, Menlo Park, CA (HJL). Jason Doctor's research was made possible by a grant from the United States Department of Health and Human Services, National Institutes of Health, National Center for Medical Rehabilitation Research (NIH-NCMRR: K01HD01221). Han Bleichrodt's research was made possible by a grant from the Netherlands Organization for Scientific Research (NWO). We thank Les Lenert, Dennis Revicki, Sean Sullivan, Anne Stiggelbout, and Jose-Luis Pinto Prades in assisting with the review and identification of articles for the analysis. We also thank the two anonymous reviewers for their helpful comments on the paper. Revision accepted for publication 10 September 2007.

Address correspondence to Jason N. Doctor, PhD, Department of Clinical Pharmacy & Pharmaceutical Economics & Policy, School of Pharmacy, University of Southern California, Los Angeles, CA 90089-9004; e-mail: jdoctor@u.washington.edu.

DOI: 10.1177/0272989X07312478

Table 1 Known Predominantly Upward (+) and Downward (–) Causes of Systematic Error in SG, TTO, and RS Values

Type of Effect	SG	TTO	RS
Loss aversion	+	+	No effect
Scale compatibility	Ambiguous	+	No effect
Utility curvature	No effect	–	No effect
Probability weighting	+	No effect	No effect

Note: SG = standard gamble; TTO = time tradeoff; RS = rating scale.

Errors that affect measurement may be divided into 2 classes: 1) *systematic error*—misestimation of a measurement value that is persistent both in direction and magnitude, and 2) *nonsystematic error*—misestimation of a measurement value that is variable in magnitude and direction. Over many observations, systematic error endures and nonsystematic error abates. We capitalize on this fact, to study within a meta-analytic framework the net effect of health utility bias. As we will explain next, the TTO and the SG are affected by systematic biases and the RS by nonsystematic biases. Consequently, over many studies the bias in the RS may decrease whereas the bias in the TTO and the SG remains. By pooling the results from many studies the comparison of the TTO and the SG with the RS can, therefore, give insight in the direction of the bias in the TTO and the SG. It is important to emphasize that we do not claim that the RS is the gold standard in health utility measurement. Any single RS measurement will be affected by biases. Our point is that over many studies these biases will be reduced and this property provides a benchmark with which to compare the TTO and the SG.

SYSTEMATIC ERROR IN HEALTH STATE VALUATIONS

The TTO and SG methods are susceptible to several known effects that lead to persistent, or systematic, errors. These effects are: loss aversion, scale compatibility, utility curvature over life duration, and probability weighting. A review of these effects is beyond the scope of this study and can be found elsewhere (see Bleichrodt⁴). These biases alter scores such that they deviate from a value that best characterizes preference for a health state, thus making scores too high or too low. They generally increase SG scores, have both upward and downward effects on TTO scores, and are predicted to have no effect on RS scores. Table 1 provides a summary of the aforementioned known predominantly

upward (+) and downward (–) causes of systematic error in SG, TTO, and RS values.

NONSYSTEMATIC ERROR IN HEALTH STATE VALUATION

Although the RS method is not susceptible to known systematic biases, individual observations are well known to be influenced by nonsystematic error resulting from contextual bias. With the RS method, the respondent's task is to assign categories (typically integer numbers) to health state stimuli such that succeeding categories represent equal steps in value. However, empirical research has demonstrated that characteristics of an RS response depend on the range and frequency of other health states being rated.^{5–7} Figure 1 illustrates range and frequency effects for a health state with a bias-free health state value of 0.40.

In each panel the x-axis represents bias free value and the y-axis denotes observed value. In the left panel, labeled “Range Effect,” 1 group of respondents rated the health state in context (C_1) which includes a limited range of health state values (range = 0.30 to 0.70). Because of a desire to spread responses over the full range of the response scale, the observed rating differs in C_1 as compared with C_2 , a context with a broader range of health state values (range = 0.0 – 1.0). In the right panel, labeled “Frequency Effect,” the health state is presented among a set of health states where a preponderance have either low subjective value (C_3) or high subjective value (C_4). With the frequency effect, the observed rating response is more sensitive to changes in value when most stimuli are of similar value to the state being evaluated. An important point is that range and frequency effects produce error magnitude and direction that is specific to context; hence, error is not systematic but changes with context. Schwartz⁸ applied range-frequency theory to explain with great precision contextual bias in RS scores reported elsewhere.⁵

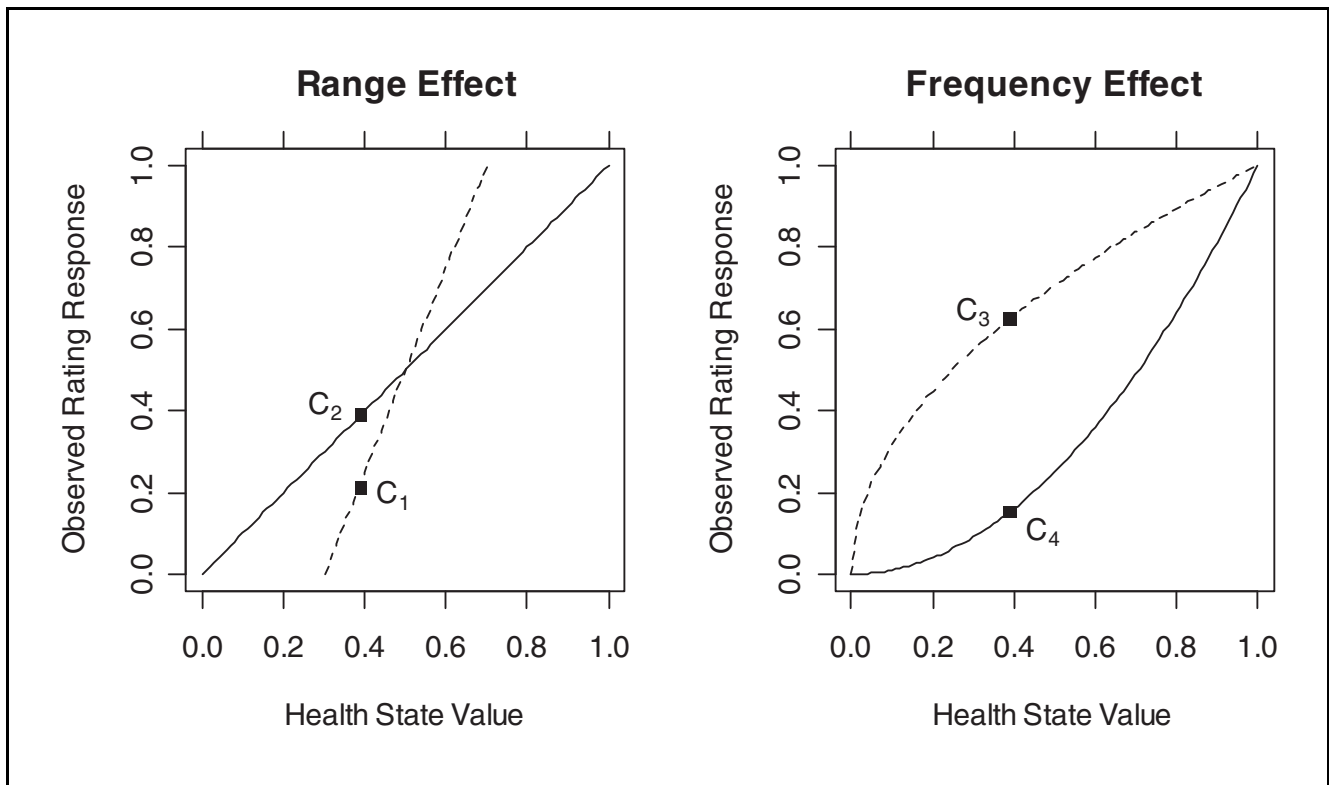


Figure 1 Observed rating responses for a hypothetical health state with “context free” value of 0.40 presented in 4 between-subject contexts: restricted stimulus range (C_1), broad stimulus range (C_2), positively skewed stimulus set (C_3), and negatively skewed stimulus set (C_4). The left panel shows a range effect on observed rating response (C_1 v. C_2); the right panel shows a frequency effect on observed rating response (C_3 v. C_4).

Robinson et al.⁶ confirmed this finding in a separate experiment. Pollack^{9,10} demonstrated convincingly that rating scales could be unbiased when contextual factors were varied iteratively over many experiments; that is, Pollack^{9,10} identified and subsequently manipulated bias effects to neutralize bias. The non-systematic nature of rating scale context bias suggests that over many naturally occurring studies, rating scale bias may decrease in size.

Whether or not SG or TTO values are influenced by nonsystematic factors like context has received much less attention. Robinson et al.⁶ found in a context manipulation experiment that SG values were much less susceptible to context effects than were RS values. We are unaware of any studies examining context effects and TTO responses.

COMPARING RS, TTO, AND SG VALUES

Empirically, RS, TTO, and SG values do not appear to agree. A common assertion is that RS values are

lower than TTO and SG values.^{1,2} However, given that the RS is subject to a context bias, one may not conclude from any single study that RS values are lower or higher than TTO or SG values. This caveat applies even when no explicit context is given, in particular, when respondents rate only their current health. Birnbaum¹¹ has shown that when not given an explicit context, respondents choose their own contexts and choose different ones for different stimuli. He was in fact able to show through a between-subjects experiment that the number “9” achieved a higher largeness rating than the number “221.” Presumably, “9” is large in the context of 1-digit numbers and “221” is small in the context of 3-digit numbers. Such an effect is not easily alleviated by explicit use of anchors at points along the rating scale.^{11,12} Hence, conclusions about relative value differences between TTO (or SG) and RS drawn from data collected within any single study in which not every respondent rated the same health states are also not likely trustworthy. Only by comparing RS values against TTO (or SG) values in explicit contexts, across many studies

Table 2 Journals Searched by Hand for Full-Length Articles and Conference Abstracts Possibly Missed by MEDLINE Search

Journal Title	Search Interval
<i>Health Economics</i>	1993-2002
<i>Health Policy</i>	1984-1989
<i>Health Policy in Amsterdam and The Netherlands</i>	1989-2000
<i>International Journal of Technology Assessment in Health Care</i>	1985-Present
<i>Journal of Health Economics</i>	1984-2002
<i>Medical Care</i>	1978-Present
<i>Medical Decision Making</i>	1981-Present
<i>Quality of Life Research</i>	1993-Present
<i>Pharmacoeconomics</i>	1992-Present

and administered within-subject, is it likely that context effects will diminish. In this study, using a meta-analytic approach, we address the question of the overall effect of bias on TTO and SG scores. We capitalize on the fact that although the TTO and SG are susceptible to biases that result in systematic error in health state value, another method, the rating scale (RS), is susceptible to contextual effects that are non-systematic across studies. Hence, although nonsystematic error diminishes when RS data are aggregated over many studies, systematic TTO and SG method error should persist.

METHODS

Search Strategy and Inclusion Criteria

We searched (with no language restrictions) for all reports where RS and TTO measures, or RS and SG measures, were given to the same subjects evaluating the same health state at any 1 measurement interval. We performed a MEDLINE search using the following queries in all fields: 1) (rating scale OR category scale OR visual analogue scale OR visual analog scale) AND (time tradeoff OR time trade-off), and 2) (category scale OR rating scale OR visual analogue scale OR visual analog scale) AND standard gamble. These searches were thought to be general enough to contain, as a smaller subset, as many studies as possible within our inclusion criteria (listed below). The search period was January 1, 1976, through December 31, 2004. We also completed a second manual search of 9 journals that are well known to publish health utility data (see Table 2).

This second search was conducted to: 1) identify articles possibly missed by the MEDLINE search,

and 2) extract results from abstracts published from conference proceedings printed in a subset of the journals listed in Table 2. The latter was done to avoid publication bias. When findings reported in an abstract were later published as a full-length article, only the data from the full-length article were used in the meta-analysis. We complemented our search by reviewing the reference lists from original research and review articles. Finally, we circulated the list of studies we found to 5 experts in the field to see whether they could come up with more studies. Each expert was a lead or senior author on an article found on the list generated by our search methods. Four experts accepted and 1 declined on the grounds that she had not worked in the area for some time. The expert who declined did recommend a well-known replacement who agreed to serve as the 5th expert.

Inclusion criteria were: 1) studies that elicited, for the same set of subjects, multiple methods of utility assessment; 2) multiple methods had to include the RS method along with either the SG or TTO methods; 3) all subjects had to receive the same health state descriptions; 4) reported utility scores had to be elicited and could not be predicted from formulas or multiattribute questionnaires (e.g., EQ-5D, Health Utilities Index, or Quality of Well-Being Scale); and 5) for TTO studies, duration in current health had to exceed 5 years due to a documented unwillingness to trade time over short durations.¹³ After consultation with experts another inclusion criterion was added: Health states had to be evaluated by respondents as “better than death.” Studies that did not meet the inclusion criteria were excluded. We note that by our 3rd criterion, health state descriptions had to be hypothetical and could not reflect an individual’s unique current health description, nor could the health state choice set be manipulated in a between-subjects experiment.

We contacted the investigators of the original studies if information was required to establish inclusion criteria or information on utility for health state was not available in the published reports. Missing data that could not be resolved by attempts to contact the authors were median imputed. Two investigators abstracted the articles. They resolved disagreements by consensus.

Statistical Analysis

Using the *rmeta* package within the statistical computing language R,¹⁴ we conducted 2 meta-analyses on effect size data over the aforementioned studies.

The primary meta-analysis compared within-subject effect sizes for RS and TTO score differences. A secondary meta-analysis compared within-subject effect sizes for RS and SG score differences. A standard effect size (d) estimate for within-subject score differences was used¹⁵:

$$d = \frac{M_{RS} - M_z}{SD_{diff}}, \quad (1)$$

where M_{RS} is the mean RS score, M_z is the mean score for the competing method (either SG or TTO), and SD_{diff} is the standard deviation of the difference scores between the RS and competing method. In our case, the effect size estimates the average score difference (between 2 utility elicitation methods) relative to the variability in task performance in the population. To compute standard deviation of difference scores, an estimate of the population correlation between RS and TTO and RS and SG ratings is needed.¹⁶ Whereas several correlation statistics on these rating methods have been given in the early QALY literature (see Torrance,¹⁷ Wolfson et al.,¹⁸ and Read et al.¹⁹), Nickerson has differentiated among several types of correlations between utility elicitation methods and recommends use of a mean within-respondent correlation in any analysis postulating that psychological processes affect response.^{20(p494)} Such is the case with our current analysis, which considers that responses are affected by psychological biases.

Two articles provide appropriate (mean within-respondent) correlations for our meta-analytic purposes; they are Kartman et al.²¹ and Krabbe et al.²² With respect to the mean within-respondent correlation, r , between RS and TTO scores, Krabbe et al.²² report this value as $r = 0.23$, whereas Kartman et al.²¹ report a value of $r = 0.25$. For this analysis, we report our results under the assumption of the middle value between these two, $r = 0.24$. For the RS and SG difference score meta-analysis, we report our results under the assumption that $r = 0.19$. This is halfway between the value reported by Krabbe et al.,²² $r = 0.22$, and that of Kartman et al.,²¹ $r = 0.16$. For each analysis we also ran meta-analyses under the range of standard error assumptions as given by the range of published correlations between measures. This was done to determine the robustness of our findings. Context bias associated with the rating scale depends on the specific study methods but is statistically independent across studies. Therefore, to preserve this independence assumption, an average effect size—computed over utilities elicited for multiple health states within study—served as the dependent variable.

We chose to conduct random effects (as opposed to fixed effects) analyses of data because rating scale context bias would naturally produce statistically heterogeneous effect sizes across studies. The random effects model incorporates a between-study component of variance to address heterogeneity, whereas a fixed effects model does not. An effect size and confidence interval plot as well is given for the primary analysis.

In addition to analysis on raw standard gambles, we conducted 2 meta-analyses on corrected scores. A correction formula that adjusts for the effects of bias associated with prospect theory²³ (loss aversion, framing, and probability weighting) has been proposed²⁴ and applied elsewhere.²⁵ The first formula we used corrected for only probability weighting.^{26,27} We applied a 1-parameter weighting function as given in Tversky and Kahneman²³ to standard gamble scores (with the standard assumption that $\gamma = .61$ (see Wakker and Stiggelbout^{27(p309)}). This gives a standard gamble utility corrected for probability weighting. The second analysis used corrections for standard gamble bias (Table 3).²⁴ In addition to correcting for probability weighting, this table of values corrects for loss aversion and framing effects. This table has been used successfully to correct SG bias in other work.²⁴

Finally, an evaluation of study quality was considered. We evaluated the extent to which studies we examined adhered to reporting standards for studies of utilities. Each study received a quality score based on adherence to 10 components of reporting standards given in Table 1 of Stalmeier et al.²⁸ Quality score was computed as the weighted sum of these 10 components and scaled so that a score of 100 reflected complete adherence and a score of 0 reflected complete nonadherence. Component weightings were determined by mean expert importance ratings reported in Stalmeier et al.^{28(Table1, p206)} We evaluated the correlation of study quality with effect size, standard error, and year of publication. We also used quality scores as weights to determine if this influenced meta-analytic findings.

RESULTS

With regard to the RS and TTO meta-analysis, we identified 4 articles from systematic reviews; the MEDLINE search yielded 139 results, and of these 13 met the inclusion criteria and were not already identified in the systematic review articles. An additional 2 studies (conference presentations) were included from a hand search of the journals in Table 1 and

Table 3 Corrected Standard Gamble Utilities as Proposed by Bleichrodt et al.²⁴ for Standard Gamble Elicitations between 0.00 and 0.99

	0.00	0.01	0.02	0.03	0.04	0.05	0.06	0.07	0.08	0.09
0.0	0.000	0.025	0.038	0.048	0.057	0.064	0.072	0.078	0.085	0.091
0.1	0.097	0.102	0.108	0.113	0.118	<u>0.123</u>	0.128	0.133	0.138	0.143
0.2	0.148	0.152	0.157	0.162	0.166	<u>0.171</u>	0.176	0.180	0.185	0.189
0.3	0.194	0.199	0.203	0.208	0.213	0.217	0.222	0.227	0.231	0.236
0.4	0.241	0.246	0.251	0.256	0.261	0.266	0.271	0.276	0.281	0.286
0.5	0.292	0.297	0.303	0.308	0.314	0.320	0.325	0.331	0.337	0.343
0.6	0.350	0.356	0.363	0.369	0.376	0.383	0.390	0.397	0.405	0.412
0.7	0.420	0.428	0.436	0.445	0.454	0.463	0.472	0.481	0.491	0.502
0.8	0.512	0.523	0.535	0.547	0.560	0.573	0.587	0.601	0.617	0.633
0.9	0.650	0.669	0.689	0.710	0.734	0.760	0.789	0.822	0.861	0.911

Note: Row headings represent 10ths, column headings 100ths of the uncorrected standard gamble score, and table entries are corrected scores; for example, the corrected utility for a standard gamble of 0.15 is 0.123 (underlined).

known review articles. Experts were not able to identify any additional RS and TTO studies that met our criteria. A total of 19 studies were used for the RS and TTO meta-analysis. With respect to the RS and SG meta-analysis, we identified 7 articles from systematic reviews; the MEDLINE search yielded 150 results, and of these 5 met the inclusion criteria and were not already identified in the systematic review articles. An additional 3 studies (conference presentations) were included from a hand search of the journals in Table 2. After circulating our list to experts, they were able to identify 1 additional study that met our inclusion criteria and which was added. A total of 16 studies were used for RS-SG meta-analysis. We note that, as would be expected, studies used in the RS-TTO and RS-SG meta-analyses were not mutually exclusive. A total of 27 studies were used as data. Of these studies, 11 collected only RS and TTO responses,^{3,22,29-38} 9 collected only RS and SG responses,³⁹⁻⁴⁷ and 7 collected RS, TTO, and SG responses.^{17,19,48-52}

Results indicate no significant effect for RS and TTO difference scores: effect size (95% confidence interval [CI]) = 0.04 (-0.02, 0.09). Figure 2 shows the plot of confidence intervals centered on effect size (x-axis) for each study. The “x” indicates an overall effect; the line through it is the confidence interval. Although there is a small overall effect of 0.04, the confidence interval around this estimate crosses 0.0. These results were robust over the range of reported correlations between RS and TTO values.

As mentioned previously, a quality score was determined by the extent to which studies adhered to published reporting criteria for studies of utility.²⁸ Adherence was weighted by published expert ratings

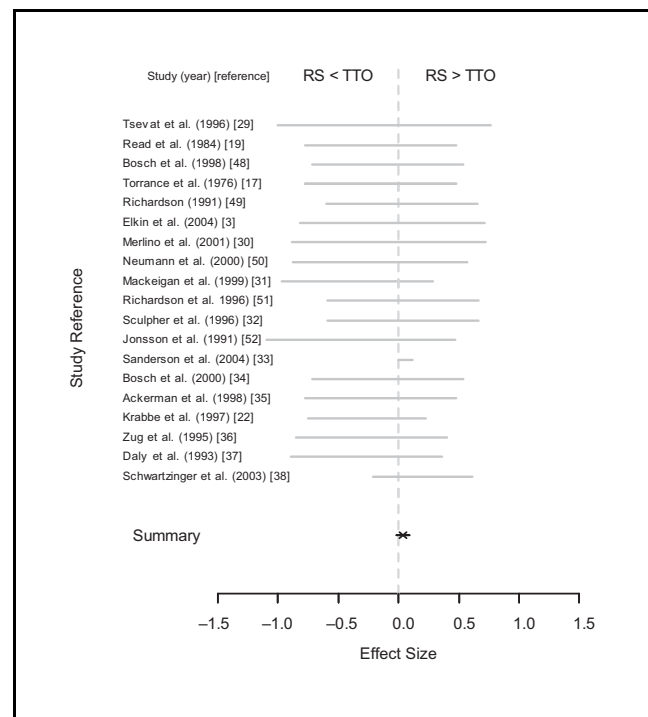


Figure 2 Plot of rating scale (RS) and time tradeoff (TTO) difference score effect sizes and confidence intervals for 19 studies.

of importance²⁸ and normalized so that a score of 100 indicates total adherence in reporting and a score of zero indicates total nonadherence. Quality scores for RS-TTO studies ranged between 21.0 and 95.7. The mean (\pm SD) importance weighted quality score for RS-TTO studies was 64.7 (\pm 17.9). An evaluation of Pearson's product-moment correlations indicated that quality score was not significantly correlated with

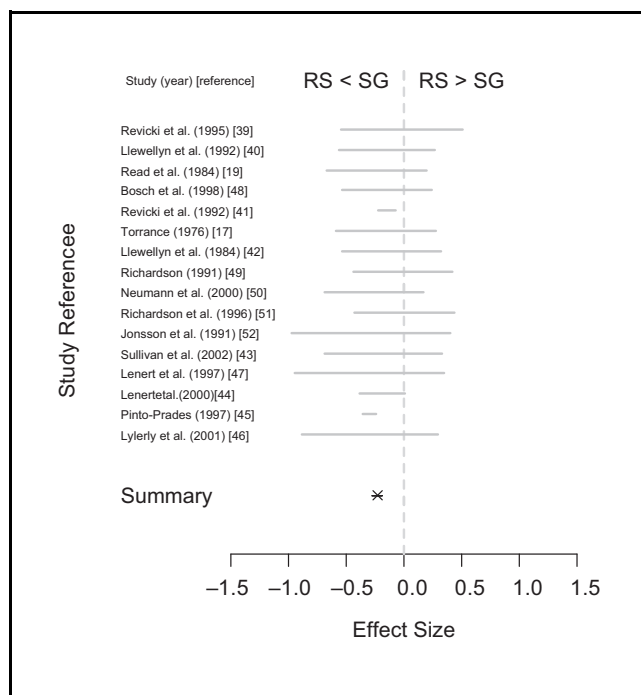


Figure 3 Plot of rating scale (RS) and standard gamble (SG) difference score effect sizes and confidence intervals for 16 studies.

effect size ($r = 0.23$, $P = \text{ns}$), standard error ($r = -0.28$, $P = \text{ns}$), or year of publication ($r = 0.0$, $P = \text{ns}$). Adding quality weights did not significantly influence the meta-analytic results in that the confidence interval for RS-TTO effect size still crossed zero.

In contrast, the meta-analysis on RS and SG values indicated that RS scores were significantly lower than SG scores: effect size (95% CI) = -0.23 ($-0.28, -0.19$). These results were robust over the range of reported correlations between RS and SG values. Figure 3 shows the plot of confidence intervals centered on effect size estimates (x-axis) for each of the 16 studies included in the analysis.

Again, the “ \times ” indicates an overall effect; the line through it is the confidence interval. The effect is sizable and the confidence interval around the estimate does not cross zero.

Quality scores for RS-SG studies also ranged between 21.0 and 95.7. The mean (\pm SD) importance weighted quality score for RS-SG studies was 59.4 (± 19.3). An evaluation of Pearson's product-moment correlations indicated that quality score was not significantly correlated with effect size ($r = 0.22$, $P = \text{ns}$), standard error ($r = 0.20$, $P = \text{ns}$), or year of publication ($r = -0.20$, $P = \text{ns}$). Adding quality weights did not significantly influence meta-analytic results in that

the confidence interval for RS-SG effect size did not overlap with 0.0 and registered SG scores as consistently higher than RS scores.

The meta-analyses on corrected SG scores revealed that the probability weighting correction was effective in reducing SG and RS difference but left a very small measurable difference between SG and RS scores: effect size (95% CI) = -0.09 ($-0.13, -0.05$). The correction adjusting for loss aversion, framing, and probability weighting²⁴ (see Table 3) eliminated differences altogether: effect size (95% CI) = 0.01 ($-0.03, 0.05$).

DISCUSSION

An early influential review of the health utility field suggested that TTO scores were higher than RS scores.¹ This assertion was based on the best available data at the time and has remained largely unchallenged. However, 15 years later we find that contrary to this notion that RS scores are lower than TTO scores, RS and TTO scores are about equal when data are examined systematically over many within-subject studies. This may indicate that when RS context bias diminishes, value measurement becomes consistent and TTO and RS values agree. Another interpretation of this result is that competing systematic TTO biases may cancel out. Hence, TTO scores may be relatively unbiased within a study. In either case, the discrepancy between our result that TTO and RS agree and the previous result that TTO scores exceed RS scores is likely due to diminishing RS context bias unique to the meta-analytic approach we used. In contrast, and as expected, SG biases, which are generally upward, result in higher scores than when the same individuals rate the same health states using the RS method. The disparity between SG and RS disappears when SG scores are corrected for probability weighting, framing, and loss aversion.

There are a few caveats to our results that deserve discussion. First, it is important to realize that our results do not suggest that RS and TTO scores are comparable or interchangeable within a study. Hence, our study should not be interpreted as offering support for the use of the RS in economic evaluations of health care. RS scores vary substantially within a study due to context effects unique to the study. Our findings show that when evaluated systematically across many studies, TTO scores do not appear to be higher than RS scores. We are inclined to interpret this as evidence that the systematic

biases in the TTO tend to cancel. Second, although no systematic RS biases are known, it is possible that 1 or more do exist,⁵³ which could threaten the interpretation that TTO scores overall do not exhibit a directional bias. However, given our current state of knowledge, we can be confident that TTO directional bias is not large in comparison to the directional bias exhibited by the SG method. Third, with respect to our analysis of SG corrections, the fundamental data element in our study is mean score for health state; it is not guaranteed that a transformed mean score will equal a mean of transformed scores. However, transformed mean scores will approach mean transformed scores as standard errors approach zero. In most cases, standard errors were low in the studies we evaluated. Fourth, other features of elicitation methodologies such as reliability, validity, and responsiveness to change are important but beyond the scope of this study.

A large body of literature assumes that because the SG is rooted in the axioms of expected utility theory and is the only scaling method that includes an element of risk inherent in most medical decisions, the SG represents the reference standard and that other methods (e.g., the RS) should be adjusted to match SG scores.⁵³ We do not agree with this point of view. There is much evidence to suggest that expected utility is not the correct descriptive model (i.e., it may not characterize observed preference behavior very well).⁵⁴ When decision makers deviate from expected utility, the SG method will generally yield biased utilities.²⁴ For this reason, our method of adjusting scores does not entrust the SG method with preeminence over other methods and does not relate RS or TTO scores via mapping them to SG as is commonly done.

A basic assumption of this study is that different methods should produce the same utilities. A practical rationale for this assumption is that if differences occur then the outcome of an economic evaluation will depend on the method used. In the absence of a gold standard for health utility measurement, this is undesirable. Such an assumption is not universally held. One theory that became popular in the 1970s and 1980s contends that risky utility (e.g., SG) and riskless value (e.g., TTO and RS) may differ by an increasing nonlinear transformation when risk aversion is considered.⁵⁵ In present day, this theory has become less popular for 2 reasons. First, it does not permit violations of expected utility theory, which are widely observed.⁵⁶ Second, it leads to serious problems in reconciling attitudes toward risk of small and large stakes gambles.⁵⁷

For these reasons risk behavior is now primarily modeled, at its source, as attitude toward chance (via nonlinear transformation of probabilities) and through the acknowledgment that decision makers are averse to losses.²³ For an excellent discussion of how this modern approach moves toward a unified notion of utility, one that has meaning prior to risk and not vice versa, see Wakker.⁵⁸ Empirical studies have shown that when attitude toward chance and loss aversion are considered, differences between riskless and risky utility tend not to prevail.^{59–61}

The findings of this study have implications for cost-effectiveness analysis. In cost-effectiveness analysis, health utility assessment is carried out so that quality weights can be assigned to health states in the analysis. As demonstrated here and elsewhere, methods and procedures applied to the same health state often result in values that are inconsistent with respect to each other. Inconsistencies mean that more than 1 quality weight can be assigned to any particular health state. However, the valid application of cost-effectiveness analyses requires that 1 and only 1 quality weight be assigned to any particular health state.

The present study is part of a growing number of studies suggesting that biases that lead to differences between measures can be reduced or eliminated. Biases appear to distort preferences in lawful and thus correctable ways, with corrections yielding greater consistency across methods. The findings of this study suggest that standard gambles may need to be corrected for probability weighting bias. Loss aversion and framing effects may also be of concern with the SG. In contrast, the findings of this study do not support a net directional systematic TTO bias and give further support to the use of raw TTO values in cost-effectiveness analysis. Finally, although RS contextual bias may diminish over many studies, unless contextual bias is manipulated and neutralized within an experiment, it is likely to adversely influence ratings in individual studies.

REFERENCES

1. Froberg DG, Kane RL. Methodology for measuring health-state preferences. II. Scaling methods. *J Clin Epidemiol.* 1989;42:459–71.
2. Drummond MF, O'Brien B, Stoddart G, Torrance GW. *Methods for the Economic Evaluation of Health Care Programmes.* 2nd ed. Oxford (UK): Oxford University Press; 1997.
3. Elkin EB, Cowen ME, Cahill D, Steffel M, Kattan MW. Preference assessment method affects decision-analytic recommendations: a prostate cancer treatment example. *Med Decis Making.* 2004;24:504–10.

4. Bleichrodt H. A new explanation for the difference between SG and TTO utilities. *Health Econ.* 2002;11:447–56.
5. Bleichrodt H, Johannesson M. An experimental test of a theoretical foundation for rating scale valuations. *Med Decis Making.* 1997;17:208–16.
6. Robinson A, Loomes, G, Jones-Lee, M. Visual analog scales, standard gambles, and relative risk aversion. *Med Decis Making.* 2001;21:17–27.
7. Parducci A. Category judgment: a range-frequency model. *Psychol Rev.* 1965;75:407–18.
8. Schwartz A. Rating scales in context. *Med Decis Making.* 1998;18:236.
9. Pollack I. Iterative techniques for unbiased rating scales. *Q J Exp Psychol.* 1965;17:139–48.
10. Pollack I. Neutralization of stimulus bias in the rating of grays. *J Exp Psychol.* 1965;69:564–78.
11. Birnbaum MH. How to show that $9 > 221$: collect judgments in a between-subjects design. *Psychol Methods.* 1999;4:243–9.
12. Birnbaum MH. Using contextual effects to derive psychophysical scales. *Percept Psychophys.* 1974;15:89–96.
13. McNeil BJ, Weichselbaum R, Pauker S. Speech and survival: tradeoffs between quality and quantity of life in laryngeal cancer. *N Engl J Med.* 1981;305:982–7.
14. Ihaka R, Gentleman R. R—a language for data analysis and graphics. *J Comp Graph Stat.* 1996;5:299–314.
15. Morris SB, DeShon RP. Combining effect size estimates in meta-analysis with repeated measures and independent-groups designs. *Psychol Methods.* 2002;7:105–25.
16. Howell DC. *Fundamental Statistics for the Behavioral Sciences.* 2nd ed. Boston: PWS-Kent; 1989.
17. Torrance GW. Social preferences for health states: an empirical evaluation of three measurement techniques. *Socio-Econ Plan Sci.* 1976;10:129–36.
18. Wolfson AD, Sinclair AJ, Bombardier C, McGeer A. Preference measurements for functional status in stroke patients: interrater and intertechnique comparisons. In: Kane RL, Kane RA, eds. *Values and Long Term Care.* Lexington (MA): Lexington Books; 1982. p 191–214
19. Read JL, Quinn RJ, Berwick DM, Fineberg HV, Weinstein MC. Preferences for health outcomes: comparison of assessment methods. *Med Decis Making.* 1984;4:315–29.
20. Nickerson CE. Assessing convergent validity of health-state utilities obtained using different scaling methods. *Med Decis Making.* 1999;19:487–96.
21. Kartman B, Gatz G, Johannesson M. Health state utilities in gastroesophageal reflux disease patients with heartburn: a study in Germany and Sweden. *Med Decis Making.* 2004;24:40–52.
22. Krabbe PFM, Essink-Bot M, Bonsel, GJ. The comparability and reliability of five health-state valuation methods. *Soc Sci Med.* 1997;45:1641–52.
23. Tversky A, Kahneman D. Advances in prospect theory: cumulative representation of uncertainty. *Journal of Risk and Uncertainty.* 1992;5:297–323.
24. Bleichrodt, H, Pinto JL, Wakker PP. Making descriptive use of prospect theory to improve the prescriptive use of expected utility. *Manage Sci.* 2001;47:1498–514.
25. van Osch SMC, Wakker PP, van den Hout WB, Stiggelbout AM. Correcting biases in standard gamble and time trade-off utilities. *Med Decis Making.* 2004;24:511–7.
26. Prelec D. The probability weighting function. *Econometrica.* 1998;66:497–527.
27. Wakker P, Stiggelbout A. Explaining distortions in utility elicitation through the rank-dependent model for risky choices. *Med Decis Making.* 1995;15:180–6.
28. Stalmeier PF, Goldstein MK, Holmes AM, et al. What should be reported in a methods section on utility assessment? *Med Decis Making.* 2001;21:200–7.
29. Tsevat J, Solzan JG, Kuntz KM, et al. Health values of patients infected with human immunodeficiency virus: relationship to mental health and physical functioning. *Med Care.* 1996;34:44–57.
30. Merlino LA, Bagchi I, Taylor TN, et al. Preference for fractures and other glucocorticoid-associated adverse effects among rheumatoid arthritis patients. *Med Decis Making.* 2001;21:122–32.
31. Mackeigan LD, O'Brien BJ, Oh PI. Holistic versus composite preferences for lifetime treatment sequences for type 2 diabetes. *Med Decis Making.* 1999;19:113–21.
32. Sculpher M, Michaels J, McKenna M, Minor J. A cost-utility analysis of laser-assisted angioplasty for peripheral arterial occlusions. *Int J Technol Assess Health Care.* 1996;12:104–25.
33. Sanderson K, Andrews G, Corry J, Lapsley H. Using the effect size to model change in preference values from descriptive health states. *Qual Life Res.* 2004;13:1255–64.
34. Bosch JL, Hunink MGM. Comparison of the health utilities index mark 3 (HUI3) and the EuroQol EQ-5D in patients treated for intermittent claudication. *Qual Life Res.* 2000;9:591–601.
35. Ackerman SJ, Beusterien KM, Mafilios MS, Wood MR. Measuring preferences for living in U.S. states: a comparison of the rating scale, time trade-off, and standard gamble. *Acad Radiol.* 1998;5 Suppl 2:S291–6.
36. Zug KA, Littenberg B, Baughman RD, et al. Assessing the preferences of patients with psoriasis—a quantitative, utility approach. *Arch Dermatol.* 1995;131:561–8.
37. Daly E, Gray A, Barlow D, McPherson K, Roche M, Vessey M. Measuring the impact of menopausal symptoms on quality of life. *BMJ.* 1993;307:836–40.
38. Schwarzing M, Stouthard MEA, Burstrom K, Nord E; European Disability Weights Group. Cross-national agreement on disability weights: the European Disability Weights Project. *Population Health Metrics.* 2003;1:9–19.
39. Revicki DA, Wu AW, Murray MI. Change in clinical status, health status, and health utility outcomes in HIV-infected patients. *Med Care.* 1995;33:AS173–82.
40. Llewellyn-Thomas JA, Thiel EC, McGreal MJ. Cancer patients' evaluations of their current health states: the influence of expectations, comparisons, actual health status, and mood. *Med Decis Making.* 1992;12:115–22.
41. Revicki DA. Relationship between health utility and psychometric health status measure. *Med Care.* 1992;30:MS274–82.
42. Llewellyn-Thomas H, Sutherland HJ, Tibshirani R, Ciampi A, Till JE, Boyd NF. Describing health states. Methodologic issues in obtaining values for health states. *Med Care.* 1984;22:543–52.

43. Sullivan SD, Lew DP, Devine EB, et al. Health state preference assessment in diabetic peripheral neuropathy. *Pharmacoeconomics*. 2002;20:1079–89.
44. Lenert LA, Ziegler J, Lee T, Sommi R, Mahmoud R. Differences in health values among patients, family members, and providers for outcomes in schizophrenia. *Med Care*. 2000;38:1011–21.
45. Pinto-Prades J. Is the person trade-off a valid method for allocating health care resources? *Health Econ*. 1997;6:71–81.
46. Lyerly AD, Myers ER, Faden RR. The ethics of aggregation and hormone replacement therapy. *Health Care Anal*. 2001;9:187–221.
47. Lenert LA, Soetikno RM. Automated computer interviews to elicit utilities: potential applications in the treatment of deep venous thrombosis. *JAMIA*. 1997;49–56.
48. Bosch JL, Tetteroo E, Mali WP, Hunink MG. Iliac arterial occlusive disease: cost-effectiveness analysis of stent placement versus percutaneous transluminal angioplasty: Dutch Iliac Stent Trial Study Group. *Radiology*. 1998;208:641–8.
49. Richardson J. Economic assessment of health care: theory and practice. *Australian Economics Review*. 1991;93:4–21.
50. Neumann PJ, Blumenschein K, Zillich A, et al. Relationship between FEV1% predicted and utilities in adult asthma. Presented at Society for Medical Decision Making, 22nd Annual Meeting, Cincinnati, OH, 2000.
51. Richardson J, Hall J, Salkeld G. The measurement of utility in multiphase health states. *Int J Technol Assess Health Care*. 1996;12:151–62.
52. Jonsson B, Horisberger B, Bruguera M, Matter L. Cost-benefit analysis of hepatitis-B vaccination. *Int J Technol Assess Health Care*. 1991;7:379–402.
53. Torrance GW, Feeny D, Furlong, W. Visual analog scales: do they have a role in the measurement of preferences for health states? *Med Decis Making*. 2001;21:329–34.
54. Tversky A, Kahneman D. The framing of decisions and the psychology of choice. *Science*. 1981;211:453–8.
55. Dyer JS, Sarin RK. Relative risk aversion. *Manage Sci*. 1982;28:875–86.
56. Starmer C. Developments in non-expected utility theory: the hunt for a descriptive theory of choice under risk. *J Econ Lit*. 2000;38:332–382.
57. Rabin M. Risk aversion and expected-utility theory: a calibration theorem. *Econometrica*. 2000;68(5):1281–92.
58. Wakker P. Separating marginal utility and probabilistic risk aversion. *Theory and Decision*. 1994;36:1–44.
59. Stalmeier PFM, Bezembinder TGG. The discrepancy between risky and riskless utilities: a matter of framing? *Med Decis Making*. 1999;19:435–47.
60. Abdellaoui M, Barrios C, Wakker PP. Reconciling introspective utility with revealed preference: experimental arguments based on prospect theory. *J Econ*. 2007;138:356–78.
61. Attema AE, Bleichrodt H, Wakker PP. Measuring the utility of life duration in a risk-free versus risky situation. Working paper, Erasmus University; 2006.