



Published in final edited form as:

Med Decis Making. 2012 ; 32(2): 273–286. doi:10.1177/0272989X11418671.

Agreement About Identifying Patients Who Change Over Time: Cautionary Results in Cataract and Heart Failure Patients

David Feeny, PhD,

The Center for Health Research, Kaiser Permanente Northwest and Health Utilities Incorporated

Karen Spritzer, BA,

Department of Medicine, University of California, Los Angeles

Ron D Hays, PhD,

Department of Medicine, University of California, Los Angeles

Honghu Liu, PhD,

School of Dentistry, University of California, Los Angeles

Theodore G. Ganiats, MD,

Department of Family and Preventive Medicine, University of California, San Diego

Robert M. Kaplan, PhD,

Department of Health Services Research, University of California, Los Angeles

Mari Palta, PhD, and

Department of Population Health Sciences, University of Wisconsin-Madison

Dennis G. Fryback, PhD

Department of Population Health Sciences, University of Wisconsin-Madison

Abstract

Background—Preference-based measures of health-related quality of life all use the same dead = 0.00 to perfect health = 1.00 scale, but there are substantial differences among measures.

Objective—The objective is to examine agreement in classifying patients as better, stable, or worse.

Design—The EQ-5D, Health Utilities Index Mark 2 and Mark 3, Quality of Well-Being – Self-Administered, Short-Form 36 (Short-Form 6D), and disease-targeted measures were administered prospectively in two clinical cohorts.

Setting—The study was conducted at academic medical centers: University of California, Los Angeles; University of California, San Diego; University of Wisconsin-Madison; and University of Southern California.

Patients—Patients undergoing cataract extraction surgery with lens replacement completed the 25-item National Eye Institute Visual Function Questionnaire (NEI-VFQ-25). Patients newly

Contact Information for Corresponding Author and Reprint Requests, David Feeny, The Center for Health Research, Kaiser Permanente Northwest, 3800 North Interstate Avenue, Portland, OR 97227-1110 USA, Telephone: (503) 528-3937, FAX: (503) 335-2428, david.feeny@kpchr.org.

Conflict of Interest. David Feeny has a proprietary interest in Health Utilities Incorporated, Dundas, Ontario, Canada. HUInc. distributes copyrighted Health Utilities Index (HUI) materials and provides methodological advice on the use of HUI. None of the other authors declare a conflict of interest.

An earlier version of the paper was presented at the 2010 meeting of Health Technology Assessment International, Dublin, June 6–9, 2010 and at the 17th Annual Meeting of the International Society for Quality of Life Research, London, October 27–30, 2010.

referred to congestive heart failure specialty clinics completed the Minnesota Living with Heart Failure Questionnaire (MLHF).

Measurements—In both cohorts subjects completed surveys at baseline, one and six months. The NEI-VFQ-25 and MLHF were used as gold standards to assign patients to categories of change. Agreement was assessed using kappa.

Results—376 cataract patients were recruited. Complete data for baseline and the one-month follow-up were available on all measures for 210 cases. Using criteria specified by Altman, agreement was poor for six of nine pairs of comparisons and fair for three pairs. 160 heart failure patients were recruited. Complete data for baseline and the six-month follow-up were available for 86 cases. Agreement was negligible for five pairs and fair for one.

Limitations—The study was conducted on selected patients at a few academic medical centers.

Conclusions—The results underscore the lack of interchangeability among different preference-based measures.

Introduction

Preference-based measures of health-related quality of life (HRQL) are needed for monitoring population health and for program evaluation for comparative effectiveness research. Most importantly, these measures are required for estimating quality-adjusted life years (QALYs). A number of widely used generic preference-based measures are available such as the EQ-5D (1), Health Utilities Index Mark 2 (HUI2) and Mark 3 (HUI3)(2), the Quality of Well-Being - Self-Administered scale (QWB-SA) (3), and Short-Form 6D (SF-6D) (4;5). Although these measures share a common core (6;7) and all include items on mobility, mental health, and pain, there are also important differences with respect to which attributes (dimensions or domains of health status) are included. HUI and the QWB-SA include vision, hearing, speech, and dexterity; the EQ-5D and SF-6D do not. The QWB-SA is unique in that it includes 58 symptoms or health problems, only some of which are included in the other measures. These measures also differ in the range of function or symptom severity covered in each attribute. The QWB-SA asks respondents if they have or do not have a problem such as pain and stiffness; in contrast HUI and SF-6D have gradients such as the categories mild, moderate, and severe pain.

These measures also differ with respect to the methods that were used to elicit preference scores with which to estimate their respective multi-attribute scoring functions, the methods for estimating those functions, and their functional forms (8). For instance, the QWB-SA scoring function is based on valuations using the visual analog scale (VAS) and a linear additive scoring function. SF-6D is based on the standard gamble (SG) and an *ad hoc* modified linear additive functional form. EQ-5D is based on the time trade-off (TTO) and an *ad hoc* modified linear additive functional form. HUI is based on transformed VAS and SG scores and a multiplicative functional form.

It is therefore not surprising that several investigators that have used two or more measures have concluded that the scores from these measures are not interchangeable (9–14). Further there is evidence from prospective studies that the estimates of absolute and/or relative change (responsiveness, including effect size (ES) and the standardized response mean (SRM)) (15) often do not agree (12;16–19).

The objective of this paper is to examine agreement among the above measures in classifying patients into the same categories of change: We want to know if the measures agree on which patients get better, remain stable, or get worse. Data from two prospective cohort studies that employed all five of the above measures as well as disease-targeted

measures are used to assess agreement among these measures: one study of patients undergoing cataract surgery, the other of patients referred for treatment for congestive heart failure by a specialty clinic.

This paper builds on an earlier paper (19) based on the data from the same study. That paper provided cohort-level estimates of responsiveness (SRM) for each of the five preference-based measures in each of the two cohorts. Responsiveness varied among measures and across cohorts. Results from that paper underscore the lack of interchangeability of scores among these measures.

This paper asks an important follow-up question. Even if overall responsiveness differs among measures, do they agree on who gets better, who gets worse, and who was stable?

Methods

Patients

Subjects for both components of the study had to be at least 35 years of age, able to give informed consent, able to hear and understand instructions in English, and have sufficient vision and ability in reading and writing English to complete questionnaires (19). *Cataract Surgery*. Patients were undergoing cataract extraction surgery with lens replacement. Patients were excluded if undergoing simultaneous glaucoma, corneal, or vitreoretinal procedures, or if they were unable to read large print versions of questionnaires. *Heart Failure*. Patients were newly referred to congestive heart failure clinics. Inclusion criteria included evidence of the presence of heart failure for at least three months defined as a left ventricular ejection fraction less than 40%. Patients classified as Class IV in the New York Heart Association system, those with a recent (< six months) myocardial infarction, unstable angina, recent (< three months) coronary artery bypass graft surgery, those on the heart transplant list, or those with recent (< three months) ventricular tachycardia were excluded.

Participants were recruited from four academic medical centers: The University of California, Los Angeles (UCLA), the University of California, San Diego (UCSD), the University of Wisconsin, and the University of Southern California (cataract patients). The study was approved by the Institutional Review Boards at each of these institutions (UCLA IRB #G05-06-096-11; UCSD Project #070435; Wisconsin M-2005-1171; USC #HS-06-00493).

Procedures

At enrollment patients were given a packet of self-administered questionnaires to complete and mail back to the UCSD Health Services Research Center (HSRC) within seven days. The HSRC mailed out the same packet for the one- and six-month follow up surveys.

Measures

The study included five of the most commonly used preference-based measures (8). There is substantial evidence on the reliability, cross-sectional construct validity, and responsiveness (longitudinal construct validity) of each of these measures in a wide variety of applications. The study also used a widely used disease-targeted measure for vision (25-Item National Eye Institute Visual Function Questionnaire) (NEI-VFQ-25)(20–22) and a prominent disease-targeted measure for heart failure (Minnesota Living with Heart Failure Questionnaire) (MLHF) (23–26).

EQ-5D-3L (hereafter: EQ-5D)—The health-status classification system of EQ-5D includes five attributes (mobility, self-care, usual activity, pain/discomfort, and anxiety/

depression) with three levels (no problem, some problem, extreme problem) per attribute (1). The EQ-5D also includes a visual analog scale (VAS) on which respondents provide a rating of their current overall health; the analyses reported here do not include the VAS scores. Health status at a point in time for a subject is described as a five-element vector, one level for each attribute. Preference-based scores for EQ-5D health states were derived using a scoring function based on TTO preferences elicited from a random sample of community dwelling residents of the United State and estimated with an *ad hoc* modified linear additive utility function (27). Scores are defined on the conventional scale in which dead = 0.00 and perfect health = 1.00; EQ-5D scores range from -0.11 (states worse than dead) to 1.00.

Health Utilities Index Mark 2 (HUI2) and Mark 3 (HUI3) (28;29)—HUI2 includes seven attributes: sensation [vision, hearing, speech], mobility, emotion, cognition, self-care, pain, and fertility. (The item on fertility was not administered in this study; fertility was assumed to be normal, level 1.) There are four or five levels per attribute in HUI2. The multiplicative HUI2 scoring function is based on preference elicitation using the VAS and SG from a random sample of community-dwelling subjects in Canada (28). Single-attribute utility scores are on a scale in which 0.00 is the score of the most disabled level in that attribute and 1.00 is the score for level 1, no problem or disability in that attribute. Overall HUI2 scores vary from -0.03 to 1.00. In addition to the overall HUI2 score, the single-attribute HUI2 sensation score was included in the analyses of data from the cataract cohort because of its relevance as a specific measure of visual function.

HUI3—HUI3 includes eight attributes (vision, hearing, speech, ambulation, dexterity, emotion, cognition, and pain and discomfort) with five or six levels per attribute. The multiplicative HUI3 scoring function is based on preference elicitation using the VAS and SG from a random sample of community dwelling subjects in Canada (29). Overall HUI3 scores vary from -0.36 to 1.00. In addition the overall HUI3 score, the single-attribute HUI3 vision score was included in the analyses of data from the cataract cohort because of its relevance.

Self-Administered Quality of Well-being Scale (QWB-SA)—The QWB-SA assesses self-reported functioning using a series of questions designed to record limitations in the previous three days, within three separate domains (mobility, physical activity, and social activity). In addition, QWB-SA includes a series of questions that ask about the presence or absence of different symptom/problem complexes. The four domain scores are combined into a total score that provides a numerical point-in-time expression of well-being that ranges from zero (0.00) for dead to one (1.00) for asymptomatic optimum functioning. The original QWB obtained preference ratings of 856 people from the general population (30). The QWB-SA used convenience samples to model preference for case descriptions and the models were shown to be highly correlated with the population ratings in the original QWB general population preferences elicitation survey. Scores range from 0.00 to 1.00; 0.09 is the minimum for a living health state. The self-administered QWB-SA has been shown to be highly correlated with the interviewer-administered QWB and to retain the psychometric properties. Extensive evaluation of reliability and validity have been published (3;3;30–32).

Self-Rated Health (SRH)—The self-rated health item (33), “In general, would you say that your health is excellent, very good, good, fair, or poor,” is a widely used measure of overall health and was therefore included in the analyses.

Short-Form 6D (SF-6D)—SF-6D is a preference-based measure based on a subset of items from the SF-36 (or SF-12)(4;5;34). SF-6D includes six attributes (physical functioning, role limitations, social functioning, pain, mental health, and vitality) with four

to six levels per attribute. The scoring function is based on SG preferences elicited from a random sample of community-dwelling subjects in the United Kingdom and estimated using an *ad hoc* linear additive functional form (4).

25-Item National Eye Institute Visual Function Questionnaire (NEI-VFQ-25)—

The NEI-VFQ-25 was designed to capture the influence of vision on a number of dimensions of HRQL including emotional well-being and social functioning (20–22). The NEI-VFQ-25 includes 25 items covering general health, general vision, near vision, distance vision, driving, peripheral vision, color vision, ocular pain, role limitations, dependency, social function, mental health, and expectations. The total score ranges from 0 to 100 with higher scores signifying better (less impaired) vision.

Visual Function Questionnaire - Utility (VFQ-UI)—Recently a preference-based index scoring system has been developed for the NEI-VFQ-25 (35) (Kowalski et al. submitted; Rentz et al. submitted), the VFQ-UI. The VFQ-UI includes a single item representing each of six domains of the NEI-VFQ: near vision (see well up close), distance vision (going out for films, sports events), role function (limited work time due to vision), mental health (worry about doing things that may embarrass because of vision), vision dependency (stay at home because of vision) and social function (see people’s reaction to things I say). The items were selected to cover a range of vision-related functioning using Rasch analyses on samples of patients with central vision loss or peripheral vision loss. The VFQ-UI defines eight vision-related health states ranging from no difficulty to stopped doing work scored on a 0.00 (dead) to 1.00 perfect health range using time-tradeoff derived preference scores.

Minnesota Living with Heart Failure Questionnaire (MLHF)—The MLHF includes 21 items covering symptoms, mental health, social life, fatigue, appetite, mobility, sleep, sexual activity, work and recreational activities, and side-effects of treatment (23–26). Overall scores range from 0 to 105 with higher scores signifying greater impairment (lower HRQL).

Criteria for clinically important change

It is important to assess a measured change with respect both to its statistical significance and its clinical importance or magnitude. Guyatt et al. (p 377) (36) provide a definition of a clinically important difference: “The MID [minimum important difference] is the smallest difference in score in the domain of interest that patients perceive as important, either beneficial or harmful, and which would lead the clinician to consider a change in the patient’s management.” There are two major methods for determining the clinical importance of a given magnitude of change: anchor-based and distribution-based approaches (36–43). In the anchor-based approach, the change in HRQL score is related to a known anchor. The anchor itself must be an independent measure and be readily interpretable such as the categories of the New York Heart Association functional classification system or ability to climb a flight of stairs. Further, there must be an appreciable association between the anchor and the target measure (36). In the distribution-based approach the magnitude of change is compared to some measure of the variability of scores. Cohen’s guidance on classifying effect sizes is an example: 0.20 small; 0.50 medium; 0.80 large (44). The anchor-based approach provides an estimate of clinically important change while the distribution-based approach provides a basis for translating raw score change into standardized units that can be used for comparisons with estimates from prior studies or existing rules of thumb(40).

For this study, a change of 0.03 or more in the overall preference score for each of the preference-based measures is interpreted as a clinically important change (2;8;11;37;45–56). Empirical estimates of clinically important change (differences) for the five preference-based measures vary from 0.01 to 0.08 with 0.03 being well represented in estimates for each of these measures.

For the single-attribute utility scores for HUI2 sensation (which includes vision) and HUI3 vision the guideline for a clinically important difference is 0.05(2). For the NEI-VFQ-25, a change of 5.0 or more in the composite score on a 0 to 100 scale is regarded as clinically important (57). For the MLHF instrument, a change of 5.0 or more in the total score (0 to 105) is regarded as clinically important (23–25). For self-rated health (SRH: excellent, very good, good, fair, poor), a movement of one or more categories is considered clinically important.

Statistical Analyses

Previous work (19) indicated that patients undergoing cataract surgery changed substantially between baseline and the one-month follow-up survey (after surgery) and were typically then stable in the period between the one- and six-month follow-ups. Analyses for the cataract cohort therefore focus on change between the baseline to one-month follow-up. Improvement was more gradual in the heart failure cohort (19). Analyses focus on the change between baseline and the six-month follow-up.

Measures of Agreement

Relative agreement in direction and size among change scores for the 10 measures used in the cataract cohort and seven measures used in the heart failure cohort was assessed using an intra-class correlation coefficient (ICC) based on a two-way mixed analysis of variance model (measures fixed, patients as random). Agreement between the disease-targeted measure (NEI-VFQ-25 for cataracts, MLHF for congestive heart failure) and each of the five (EQ-5D, HUI2, HUI3, QWB-SA, and SF-6D) preference-based measures and SRH as to whether patients had improved, were stable, or got worse was assessed using a number of measures including the per cent agreement, kappa (unweighted and weighted), and the Delta statistic, a measure of agreement that is less sensitive than kappa to the marginal distributions (58). The degree of agreement (kappa) were interpreted according to the criteria suggested by Altman (59): <0.20 poor; 0.21 – 0.40 fair; 0.41 – 0.60 moderate; 0.61 – 0.80 good; 0.81 – 1.00 very good. In addition, regarding the two disease-specific measures as gold standards, the sensitivity of each of the six generic measures to change on the disease-targeted instruments was estimated using receiver operating characteristic curves (ROC) analyses. The ROC analyses determine if the results are sensitive to the choice of the threshold for clinically important change (0.03) on the preference-based measures.

Primary analyses were conducted on a sub-set of subjects for whom there is complete data at baseline and the one month follow-up (cataract cohort) and baseline and the six-month follow-up (heart failure cohort) for all of the measured included in the analyses. Thus, any differences in agreement across measures will not be the result of differences in the subjects excluded due to missing data. Secondary data analyses were conducted for the larger sample size for which data at baseline and the designated follow-up (all available pairs with complete data) and the sample size vary by pair of measures.

Results

A total of 376 cataract patients and 160 heart failure patients were recruited to the study. The majority of patients were white, cataract patients tended to be female, heart failure patients

tended to be male, most cataract patients were 65+, and the heart failure patients tended to be younger with the majority in the 45–64 age group (Table 1).

For the cataract cohort data for baseline and one-month follow up assessments were available for 315 of the 376 cases. Complete data for all pairs for all measures were available for 210 cases. The distribution of demographic variables for those with and without complete data was similar and there were no statistically significant differences between the two groups.

For the heart failure cohort data for baseline and the six-month follow-up assessments were available for 110 of the 160 cases. Complete data for all pairs for all measures were available for 86 cases. Those with missing data were older than those without missing data and the difference between the two groups was statistically significant.

The overwhelming majority of respondents, 93%, reported that no one helped them to complete the questionnaires; 7% reported receiving help. Among those who received any help, 90% reported that someone read the questions to them; 55% reported that someone wrote the answers on the questionnaire for them; 6% reported that someone answered the questions for them; 4% reported that someone translated the questions into their language for them; and 9% reported some other kind of help. Therefore the overwhelming majority of responses were based on self-completion and self-assessment.

Scores for each of the measures at baseline and one month and the change scores for the cataract cohort are displayed in Table 2. Note that the mean change in the total score for the NEI-VFQ-25 (VFQ_t) of 9.96 exceeds the guideline for a clinically important difference of 5.00. Similarly, the mean change in overall HUI3 scores exceeds the 0.03 clinically important difference guideline. The mean change in HUI3 vision score and HUI2 sensation scores exceed the 0.05 clinically important difference guideline. The mean changes in scores for EQ-5D, QWB-SA, SF-6D, and SRH are less than the guidelines for a clinically important difference. The distribution of change scores for the VFQ_t is displayed in Figure 1. Using the change of 5 or more in VFQ_t score as the criterion, 43% of patients improved, 52% were stable, and 4% got worse.

Scores for each of the measures at baseline and six months and the mean change scores for the heart failure cohort are displayed in Table 3. Note the mean change of 8.72 in score for the MLHF exceeds the 5.00 guideline for a clinically important difference. The mean change in QWB-SA, SF-6D, and HUI3 scores exceed the guideline while the mean changes in scores for the EQ-5D, HUI2, and SRH do not. The distribution of change scores for the MLHF is displayed in Figure 2. Using the change of 5 or more in MLHF score as the criterion, 47% of patients improved, 35% were stable, and 19% got worse.

Agreement among change scores

The ICC among the 10 measures of change in the cataract cohort was 0.16 (95% confidence interval: 0.02 – 0.29). The ICC among the seven measures of change in the heart failure cohort was 0.07 (95% confidence interval: –0.14 – 0.28).

Agreement Cataract cohort

The per cent agreement varies between 33% and 57% and is displayed in Table 4 along with simple and weighted kappa statistics for agreement between pairs of measures in classifying patients as improved, stable, or worse. The simple kappa statistics for six of the nine pairs are poor and the kappa statistics for three of the pairs are fair. Fair agreement was obtained between the NEI-VFQ-25 total scores and the vision-targeted measures: HUI2 sensation, HUI3 vision, and the VFQ-UI. The results for weighted kappa are very similar to the results

for the simple kappa. Results for the delta statistics are also very similar, ranging from -0.05 (VFQ_t and SRH) to 0.31 (VFQ_t and HUI3 vision) to 0.40 (VFQ_t and VFQ-UI). Area under the curve results for the ROC analyses range from 0.44 (SRH) to 0.67 (HUI3 vision) to 0.72 (VFQ-UI). In many cases in the ROC analyses the area under the curve is less than 0.60 , indicating agreement little better than one would expect by chance. These results indicate that the lack of agreement is not sensitive to the choice of cut points for clinically important differences. Finally, results from secondary analyses for $n = 315$ (subjects for whom observations on any measure was available at baseline and at the one-month follow-up) were very similar to the results reported in Table 4 (data not shown).

Agreement heart failure cohort

The per cent agreement varies between 19% and 49% and is displayed along with simple and weighted kappa statistics for agreement between pairs of measures in classifying patients as improved, stable, or worse are for the heart failure cohort in Table 5. The simple kappa statistics are negative for five pairs, indicating agreement less than that which would occur by chance. Agreement between the MLFH and SRH is fair. Results for the weighted kappa are very similar. The results for the delta statistics also indicate little agreement, ranging from -0.33 (QWB-SA) to 0.26 (SRH). Area under the curve results from the ROC analyses range from 0.31 (QWB-SA) to 0.73 (SRH) and indicate that the results are not sensitive to the choice of cut points. Finally, results from secondary analyses for $n = 110$ (subjects for whom observations on any measure was available at baseline and at the six-month follow-up) were very similar to the results reported in Table 5 (data now shown).

Agreement among measures on classification of patients as worse, stable, or improved

Results on the extent of agreement among measures in classifying patients as improved, stable, or deteriorated for the cataract cohort are found in Table 6. Analogous results for the heart failure cohort are found in Table 7. The lack of agreement among measures evident in the ICC results reported above is evident in Tables 6 and 7. Nonetheless, many observations are aligned “on the diagonal”, indicating that there is some agreement between the disease-specific measures, VFQ and MLHF, and each of the five preference-based measures, on which patients changed and which did not.

Discussion

There is very little pair-wise agreement between the disease-targeted measures and the five preference-based measures about which patients improved, were stable, or deteriorated. In general, agreement for the cataract cohort was poor and for the heart failure cohort negligible. For the cataract cohort, the agreement between the relevant HUI single-attribute (“disease-targeted”) scores and the NEI-VFQ-25, those for HUI2 sensation and HUI3 vision, were the exceptions; agreement was fair. Agreement was also fair between the utility scored and conventional versions of the NEI-VFQ-25. Given that both of these measures are based on the same questionnaire, it is perhaps surprising that the agreement is only fair and is not clearly higher than agreement between the NEI-VFQ-25 and HUI2 sensation and HUI3 vision. In the heart failure cohort, fair agreement was observed only for the SRH. On the basis of the ROC analyses, the results reported here appear to be robust to the choice of cut points for a clinically important change.

The agreement analyses treat the disease-targeted measure as the gold standard. Yet even though there is evidence of cross-sectional and longitudinal construct validity for the two disease-targeted measures, neither can be regarded as a true gold standard. Furthermore, that vision-related or heart-related HRQL improved does not necessarily imply that overall HRQL improved. It is possible that the side effects of interventions could more than offset

the gains and therefore overall HRQL might not improve. It is also possible that even though vision- or heart-related HRQL improved, overall HRQL did not due to the burdens associated with comorbidities. The NEI-VFQ-25 asks subjects about a wide variety of difficulties that they might experience due to limited vision, including reading, hobbies, navigating, driving, going up and down stairs, interacting with others, dressing, and the amount of assistance the subject needs from others. Similarly, the MLHF asks about limitations in/problems with mobility, sexual activity, interacting with others, fatigue, hobbies, worry, concentration, memory, and depression that the subject experiences due to the subject's heart condition. Although the breadth of coverage of these disease-targeted measures probably reduces the scope for a discrepancy between trends in vision- or heart-related HRQL and overall HRQL, it does not eliminate the possibility for such discrepancies.

The results on overall change in measures underscore that scores from these five preference-based measures are not interchangeable (Table 2). In the cataract cohort, using published guidelines (2;57) on clinically important differences/changes, clearly clinically important change was detected by the NEI-VFQ-25 and HUI3. Given that vision is included in HUI3 and that the NEI-VFQ-25 is a vision-targeted measure, this results is not surprising. But vision is included in the QWB-SA and the overall score did not reflect the gain in HRQL that was measured by the NEI-VFQ-25 and HUI3. Of course, it should be noted that only the "worst" symptom for that subject in the QWB-SA is used to compute the overall score and further that in a relatively elderly cohort, it is likely that many subjects were experiencing symptoms that are more burdensome than impaired vision and thus the vision item frequently did not affect the calculation of the QWB-SA score. Others have noted the lack of responsiveness of the EQ-5D to the effects of cataract surgery(60).

For the heart failure cohort, using published guidelines (23–25;46;47;51) on clinically important differences/changes, the MLHF, HUI3, QWB-SA, and SF-6D recorded clinically important change. Fatigue and shortness of breath symptoms on the QWB-SA as well as the mobility, physical, and social activity scales may have captured some of the effects of heart failure on HRQL. Similarly, the physical functioning, vitality, and role attributes on SF-6D may have registered some of the effects. HUI3 ambulation may have performed similarly(61)

It should be noted that reliability is less than perfect for each of the measures used in the study (62), so disagreement between change scores is also influenced by measurement error and short-term fluctuations in health that are unrelated to the conditions of primary interest. Change over time is measured with even less precision than absolute scores at a point in time. Jones and Feeny (63) and Pickard et al. (64) found lower levels of agreement between proxy and self-report for change scores than was evident for cross-sectional comparisons of baseline and follow-up scores. Other investigators have pointed out that due to the size of measurement error typically found in HRQL measures, change must be often be quite substantial for measures to agree (62). Our results were not sensitive to the magnitude of change that was considered clinically important, but the magnitude of true underlying changes does influence the agreement that can be expected. Hence, some measures of change may have better agreement than found in our studies, when reflecting interventions with larger overall effects on HRQL.

Some limitations of the study should be noted. Because the analyses are based on subjects for whom both baseline and the designated follow-up assessments were available, the results are not necessarily representative of the experience of the entire inception cohorts. In the cataract cohort those with and without complete data were similar. In the heart failure cohort, those with missing data were older than those without missing data. If older patients

experienced less improvement in HRQL than younger patients, it is possible that the estimate of change based on subjects for whom we had complete data is biased upwards. As noted in the Results, 7% of subjects had help in completing questionnaires so responses could have been influenced by others. Another limitation is that while the scoring functions for the QWB-SA and EQ-5D are based on preference scores from random samples of community-dwelling adults in the US, the scoring functions for HUI2 and HUI3 are based on preferences from random samples of the Canadian population and the function for the SF-6D is based on UK preferences. There is evidence of the generalizability of the QWB scoring function, (30;65;66) the HUI2 scoring function, (28;67) and the HUI3 scoring function.(29;68–70) In contrast there is considerable variability across “national” EQ-5D scoring functions. Nonetheless having not relied exclusively on US-based scoring functions is unlikely to be an important factor influencing the results. Finally, we classified cataract and heart failure patients as changed if the absolute value of their change score was ≥ 5.0 whether or not the difference was statistically significant. Hays et al. (71) note that changes that are statistically significant at the level of an individual subject will typically exceed the guideline for a clinically important difference(72).

Conclusions

The results underscore the lack of interchangeability of scores among these five widely used preference-based measures. Not only are the absolute scores not necessarily interchangeable; in these results the change scores were also not interchangeable (12). The results also point to a lack of precision in estimating the magnitude of change in HRQL.

In making choices about which preference-based measure(s) to use in a study, investigators need to consider carefully the coverage of the health-status classification systems and the relevance of those systems to their clinical or population health application, evidence on the cross-sectional construct validity of the measures in that application, and evidence of the responsiveness (longitudinal construct validity) of the measures in that context. Further, users of the results of studies that have employed preference-based measures to assess HRQL need to interpret those results carefully.

Acknowledgments

Supported by Grant P01AG020679 from the National Institute on Aging. Drs. Kaplan and Hays were also provided support by NIH grants 1 P01 AG020679-01A2, UCLA Claude D. Pepper, Older Americans Independence Center, NIH/NIA 5P30AG028748, and CDC Grant U48 DP000056-04. Dr. Hays also received support from the UCLA Resource Center for Minority Aging Research/Center for Health Improvement in Minority Elderly (P30AG021684) and the UCLA/DREW Project EXPORT (P20MD000148 and P20MD000182). The funding agreement ensured the independence of the authors in the design, conduct, interpretation, data, writing, and publishing of the paper. The granting agencies have neither read nor approved of the contents of the paper. The authors acknowledge the contributions of Steven Tally, UCSD, to the work reported here. The authors also appreciate the help of Barbara Brody MPH and Denise Herman, MD from UCSD, Nancy Sweitzer, MD, PhD, and Neal Barney, MD, from UW, Greg Fonerow, MD, and John Bartlett, MD, from UCLA for their collaboration on subject acquisition. An earlier version of the paper was presented at the 2010 meeting of Health Technology Assessment International, Dublin, June 6–9, 2010.

Grant Support. Supported by Grant P01AG020679 from the National Institute on Aging. Drs. Kaplan and Hays were also provided support by NIH grants 1 P01 AG020679-01A2, UCLA Claude D. Pepper, Older Americans Independence Center, NIH/NIA 5P30AG028748, and CDC Grant U48 DP000056-04. Dr. Hays also received support from the UCLA Resource Center for Minority Aging Research/Center for Health Improvement in Minority Elderly (P30AG021684) and the UCLA/DREW Project EXPORT (P20MD000148 and P20MD000182).

Reference List

1. Rabin R, de Charro F. EQ-5D: A measure of health status from the EuroQol Group. *Annals of Medicine*. 2001 Jul; 33(5):337–343. [PubMed: 11491192]

2. Horsman J, Furlong W, Feeny D, Torrance G. The Health Utilities Index (HUI®): concepts, measurement properties and applications. *Health Qual Life Outcomes*. 2003 Oct 16.1(1):54. [PubMed: 14613568]
3. Kaplan, RM.; Anderson, JP. The General Health Policy Model: An Integrated Approach. In: Spilker, B., editor. *Quality of Life and Pharmacoeconomics in Clinical Trials*. Second ed.. Philadelphia: Lippincott-Raven Press; 1996. p. 309-322.
4. Brazier J, Roberts J, Deverill M. The estimation of a preference-based measure of health from the SF-36. *J Health Econ*. 2002 Mar; 21(2):271–292. [PubMed: 11939242]
5. Brazier JE, Roberts J. The estimation of a preference-based measure of health from the SF-12. *Med Care*. 2004 Sep; 42(9):851–859. [PubMed: 15319610]
6. Fryback DG, Palta M, Cherepanov D, Bolt D, Kim JS. Comparison of 5 Health-Related Quality-of-Life Indexes Using Item Response Theory Analysis. *Med Decis Making*. 2009 Oct 20.
7. Cherepanov D, Palta M, Fryback DG. Underlying Dimensions of the Five Health-Related Quality-of-Life Measures Used in Utility Assessment: Evidence From the National Health Measurement Study. *Med Care*. 2010 Aug; 48(8):718–725. [PubMed: 20613664]
8. Feeny, DH. Preference-based measures: Utility and quality-adjusted life years. In: Fayers, P.; Hays, R., editors. *Assessing quality of life in clinical trials*. Second ed.. Oxford: Oxford University Press; 2005. p. 405-429.
9. Marra CA, Esdaile JM, Guh D, Kopec JA, Brazier JE, Koehler BE, et al. A comparison of four indirect methods of assessing utility values in rheumatoid arthritis. *Med Care*. 2004 Nov; 42(11):1125–1131. [PubMed: 15586840]
10. Marra CA, Marion SA, Guh DP, Najafzadeh M, Wolfe F, Esdaile JM, et al. Not all “quality-adjusted life years” are equal. *J Clin Epidemiol*. 2007 Jun; 60(6):616–624. [PubMed: 17493521]
11. Marra CA, Woolcott JC, Kopec JA, Shojania K, Offer R, Brazier JE, et al. A comparison of generic, indirect utility measures (the HUI2, HUI3, SF-6D, and the EQ-5D) and disease-specific instruments (the RAQoL and the HAQ) in rheumatoid arthritis. *Soc Sci Med*. 2005 Apr; 60(7):1571–1582. [PubMed: 15652688]
12. Feeny DH, Wu L, Eng K. Comparing short form 6D, standard gamble, and Health Utilities Index Mark 2 and Mark 3 utility scores: results from total hip arthroplasty patients. *Quality of Life Research*. 2004 Dec; 13(10):1659–1670. [PubMed: 15651537]
13. Luo N, Johnson JA, Shaw JW, Feeny D, Coons SJ. Self-reported health status of the general adult U.S. population as assessed by the EQ-5D and Health Utilities Index. *Medical Care*. 2005 Nov; 43(11):1078–1086. [PubMed: 16224300]
14. Fryback DG, Dunham NC, Palta M, Hanmer J, Buechner J, Cherepanov D, et al. U.S. Norms for six generic health-related Quality of Life indexes from the National Health Measurement study. *Medical Care*. 2007 Dec; 45(12):1162–1170. [PubMed: 18007166]
15. Terwee CB, Dekker FW, Wiersinga WM, Prummel MF, Bossuyt PM. On assessing responsiveness of health-related quality of life instruments: guidelines for instrument evaluation. *Quality of Life Research*. 2003 Jun; 12(4):349–362. [PubMed: 12797708]
16. Marra CA, Rashidi AA, Guh D, Kopec JA, Abrahamowicz M, Esdaile JM, et al. Are indirect utility measures reliable and responsive in rheumatoid arthritis patients? *Qual Life Res*. 2005 Jun; 14(5):1333–1344. [PubMed: 16047508]
17. Hatoum HT, Brazier JE, Akhras KS. Comparison of the HUI3 with the SF-36 preference based SF-6D in a clinical trial setting. *Value Health*. 2004 Sep; 7(5):602–609. [PubMed: 15367255]
18. McDonough CM, Tosteson TD, Tosteson AN, Jette AM, Grove MR, Weinstein JN. A longitudinal comparison of 5 preference-weighted health state classification systems in persons with intervertebral disk herniation. *Med Decis Making*. 2011 Mar; 31(2):270–280. [PubMed: 21098419]
19. Kaplan RM, Tally S, Hays RD, Feeny D, Ganiats TG, Palta M, et al. Five Preference-Based Indexes in Cataract and Heart-Failure Patients Were Not Equally Responsive to Change. *Journal of Clinical Epidemiology*. 2011 May; 64(5):497–506. [PubMed: 20685077]
20. Mangione CM, Lee PP, Gutierrez PR, Spritzer K, Berry S, Hays RD. Development of the 25-item National Eye Institute Visual Function Questionnaire. *Arch Ophthalmol*. 2001 Jul; 119(7):1050–1058. [PubMed: 11448327]

21. Varma R, Wu J, Chong K, Azen SP, Hays RD. Impact of severity and bilaterality of visual impairment on health-related quality of life. *Ophthalmology*. 2006 Oct; 113(10):1846–1853. [PubMed: 16889831]
22. McDonnell PJ, Mangione C, Lee P, Lindblad AS, Spritzer KL, Berry S, et al. Responsiveness of the National Eye Institute Refractive Error Quality of Life instrument to surgical correction of refractive error. *Ophthalmology*. 2003 Dec; 110(12):2302–2309. [PubMed: 14644711]
23. Rector TS. A conceptual model of quality of life in relation to heart failure. *J Card Fail*. 2005 Apr; 11(3):173–176. [PubMed: 15812743]
24. Rector TS, Kubo SH, Cohn JN. Patients' self-assessment of their Congestive Heart Failure Part II. *Heart Failure*. 1987 Oct-Nov;:198–209.
25. Rector TS, Francis GS, Cohn JN. Patients' self-assessment of their Congestive Heart Failure. Part 1. Patient perceived dysfunction and its poor correlation with maximal exercise tests. *Heart Failure*. 1987 Oct-Nov;:192–196.
26. Garin O, Ferrer M, Pont A, Rue M, Kotzeva A, Wiklund I, et al. Disease-specific health-related quality of life questionnaires for heart failure: a systematic review with meta-analyses. *Quality of Life Research*. 2009 Feb; 18(1):71–85. [PubMed: 19052916]
27. Shaw JW, Johnson JA, Coons SJ. US valuation of the EQ-5D health states: development and testing of the D1 valuation model. *Med Care*. 2005 Mar; 43(3):203–220. [PubMed: 15725977]
28. Torrance GW, Feeny DH, Furlong WJ, Barr RD, Zhang Y, Wang Q. Multi-attribute utility function for a comprehensive health status classification system. *Health Utilities Index Mark 2*. *Med Care*. 1996 Jul; 34(7):702–722. [PubMed: 8676608]
29. Feeny DH, Furlong W, Torrance GW, Goldsmith CH, Zhu Z, DePauw S, et al. Multi-attribute and single-attribute utility functions for the health utilities index mark 3 system. *Medical Care*. 2002 Feb; 40(2):113–128. [PubMed: 11802084]
30. Kaplan RM, Bush JW, Berry CC. Health status: types of validity and the index of well-being. *Health Serv Res*. 1976; 11(4):478–507. [PubMed: 1030700]
31. Kaplan RM, Frosch DL. Decision making in medicine and health care. *Annu Rev Clin Psychol*. 2005; 1:525–556. [PubMed: 17716098]
32. Kaplan RM, Anderson JP, Patterson TL, McCutchan JA, Weinrich JD, Heaton RK, et al. Validity of the Quality of Well-Being Scale for persons with human immunodeficiency virus infection. *HNRC Group. HIV Neurobehavioral Research Center. Psychosom Med*. 1995 Mar; 57(2):138–147. [PubMed: 7792372]
33. Idler EL, Benyamini Y. Self-rated health and mortality: A review of twenty-seven community studies. *Journal of Health and Social Behavior*. 1997; 38(1):21–37. [PubMed: 9097506]
34. Ware J Jr, Kosinski M, Keller SD. A 12-Item Short-Form Health Survey: construction of scales and preliminary tests of reliability and validity. *Med Care*. 1996 Mar; 34(3):220–233. [PubMed: 8628042]
35. Revicki D, Rentz AM, Kowalski JW, Chen WH. Assessment of unidimensionality for the visual function questionnaire-utility index (VFQ-UI) items in patients with central vision loss. *Quality of Life Research*. 2010; 19(Supplement 1):49–50. – Ref Type: Abstract.
36. Guyatt GH, Osoba D, Wu AW, Wyrwich KW, Norman GR. Methods to explain the clinical significance of health status measures. *Mayo Clin Proc*. 2002 Apr; 77(4):371–383. [PubMed: 11936935]
37. Harrison MJ, Davies LM, Bansback NJ, McCoy MJ, Verstappen SM, Watson K, et al. The comparative responsiveness of the EQ-5D and SF-6D to change in patients with inflammatory arthritis. *Qual Life Res*. 2009 Nov; 18(9):1195–1205. [PubMed: 19777373]
38. Revicki DA, Osoba D, Fairclough D, Barofsky I, Berzon R, Leidy NK, et al. Recommendations on health-related quality of life research to support labeling and promotional claims in the United States. *Quality of Life Research*. 2000; 9(8):887–900. [PubMed: 11284208]
39. Revicki D, Hays RD, Cella D, Sloan J. Recommended methods for determining responsiveness and minimally important differences for patient-reported outcomes. *Journal of Clinical Epidemiology*. 2008 Feb; 61(2):102–109. [PubMed: 18177782]

40. Hays RD, Farivar SS, Liu H. Approaches and recommendations for estimating minimally important differences for health-related quality of life measures. *COPD*. 2005 Mar; 2(1):63–67. [PubMed: 17136964]
41. Leidy NK, Wyrwich KW. Bridging the gap: using triangulation methodology to estimate minimal clinically important differences (MCIDs). *COPD*. 2005 Mar; 2(1):157–165. [PubMed: 17136977]
42. Terwee CB, Roorda LD, Dekker J, Bierma-Zeinstra SM, Peat G, Jordan KP, et al. Mind the MIC: large variation among populations and methods. *J Clin Epidemiol*. 2010 May; 63(5):524–534. [PubMed: 19926446]
43. Beaton DE, van ED, Smith P, van d, V Cullen K, Kennedy CA, et al. Minimal change is sensitive, less specific to recovery: a diagnostic testing approach to interpretability. *J Clin Epidemiol*. 2011 May; 64(5):487–496. [PubMed: 21109396]
44. Cohen, J. *Statistical power analysis for the behavioral sciences*. 2nd ed.. Hillsdale, New Jersey: Lawrence Erlbaum Associates; 1988.
45. Drummond M. Introducing economic and quality of life measurements into clinical studies. *Annals of Medicine*. 2001 Jul; 33(5):344–349. [PubMed: 11491193]
46. Walters SJ, Brazier JE. What is the relationship between the minimally important difference and health state utility values? The case of the SF-6D. *Health Qual Life Outcomes*. 2003; 1:4. [PubMed: 12737635]
47. Walters SJ, Brazier JE. Comparison of the minimally important difference for two health state utility measures: EQ-5D and SF-6D. *Quality of Life Research*. 2005 Aug; 14(6):1523–1532. [PubMed: 16110932]
48. Pickard AS, Neary MP, Cella D. Estimation of minimally important differences in EQ-5D utility and VAS scores in cancer. *Health Qual Life Outcomes*. 2007; 5:70. [PubMed: 18154669]
49. Majumdar SR, Johnson JA, Bowker SL, Booth GL, Dolovich L, Ghali W, et al. A Canadian consensus for the standardized evaluation of quality improvement interventions in Type2 diabetes. *Canadian Journal of Diabetes*. 2005; 29(3):220–229.
50. Samsa G, Edelman D, Rothman ML, Williams GR, Lipscomb J, Matchar D. Determining clinically important differences in health status measures: a general approach with illustration to the Health Utilities Index Mark II. *Pharmacoeconomics*. 1999 Feb; 15(2):141–155. [PubMed: 10351188]
51. Kaplan RM. The minimally clinically important difference in generic utility-based measures. *COPD*. 2005 Mar; 2(1):91–97. [PubMed: 17136968]
52. Grootendorst P, Feeny D, Furlong W. Health Utilities Index Mark 3: evidence of construct validity for stroke and arthritis in a population health survey. *Med Care*. 2000 Mar; 38(3):290–299. [PubMed: 10718354]
53. Groessl EJ, Kaplan RM, Barrett-Connor E, Ganiats TG. Body mass index and quality of well-being in a community of older adults. *Am J Prev Med*. 2004 Feb; 26(2):126–129. [PubMed: 14751323]
54. Kontodimopoulos N, Pappa E, Papadopoulos AA, Tountas Y, Niakas D. Comparing SF-6D and EQ-5D utilities across groups differing in health status. *Quality of Life Research*. 2009 Feb; 18(1): 87–97. [PubMed: 19051058]
55. Sullivan PW, Lawrence WF, Ghushchyan V. A national catalog of preference-based scores for chronic conditions in the United States. *Med Care*. 2005 Jul; 43(7):736–749. [PubMed: 15970790]
56. Khanna D, Furst DE, Wong WK, Tsevat J, Clements PJ, Park GS, et al. Reliability, validity, and minimally important differences of the SF-6D in systemic sclerosis. *Qual Life Res*. 2007 Aug; 16(6):1083–1092. [PubMed: 17404896]
57. Globe DR, Wu J, Azen SP, Varma R. The impact of visual impairment on self-reported visual functioning in Latinos: The Los Angeles Latino Eye Study. *Ophthalmology*. 2004 Jun; 111(6): 1141–1149. [PubMed: 15177964]
58. Andres AM, Marzo PF. Chance-corrected measures of reliability and validity in K x K tables. *Stat Methods Med Res*. 2005 Oct; 14(5):473–492. [PubMed: 16248349]
59. Altman, DG. *Practical statistics for medical research*. London: Chapman & Hall; 1991.
60. Browne JP, van der Meulen JH, Lewsey JD, Lamping DL, Black N. Mathematical coupling may account for the association between baseline severity and minimally important difference values. *J Clin Epidemiol*. 2010 Aug; 63(8):865–874. [PubMed: 20172689]

61. Pressler SJ, Eckert GJ, Morrison GC, Murray MD, Oldridge NB. Evaluation of the health utilities index mark-3 in heart failure. *J Card Fail.* 2011 Feb; 17(2):143–150. [PubMed: 21300304]
62. Palta M, Han-Yeng C, Kaplan RM, Feeny D, Cherepanov D, Fryback DG. Standard error of measurement of five health utility indexes across the range of health for use in estimating reliability and responsiveness. *Med Decis Making.* 2010
63. Jones CA, Feeny DH. Agreement between patient and proxy responses of health-related quality of life after hip fracture. *J Am Geriatr Soc.* 2005 Jul; 53(7):1227–1233. [PubMed: 16108944]
64. Pickard AS, Johnson JA, Feeny DH, Shuaib A, Carriere KC, Nasser AM. Agreement between patient and proxy assessments of health-related quality of life after stroke using the EQ-5D and Health Utilities Index. *Stroke.* 2004 Feb; 35(2):607–612. [PubMed: 14726549]
65. Balaban DJ, Sagi PC, Goldfarb NI, Nettler S. Weights for scoring the quality of well-being instrument among rheumatoid arthritics. A comparison to general population weights. *Med Care.* 1986 Nov; 24(11):973–980. [PubMed: 3773579]
66. Hector RD Sr, Anderson JP, Paul RC, Weiss RE, Hays RD, Kaplan RM. Health state preferences are equivalent in the United States and Trinidad and Tobago. *Qual Life Res.* 2010 Jun; 19(5):729–738. [PubMed: 20237958]
67. Wang Q, Furlong W, Feeny D, Torrance G, Barr R. How robust is the Health Utilities Index Mark 2 utility function? *Med Decis Making.* 2002 Jul; 22(4):350–358. [PubMed: 12150600]
68. Le Gales C, Buron C, Costet N, Rosman S, Slama PR. Development of a preference-weighted health status classification system in France: the Health Utilities Index 3. *Health Care Manag Sci.* 2002 Feb; 5(1):41–51. [PubMed: 11862978]
69. Raat H, Bonsel GJ, Hoogeveen WC, Essink-Bot ML. Feasibility and reliability of a mailed questionnaire to obtain visual analogue scale valuations for health states defined by the Health Utilities Index Mark 3. *Medical Care.* 2004 Jan; 42(1):13–18. [PubMed: 14713735]
70. Ruiz M, Rejas J, Soto J, Pardo A, Rebollo I. [Adaptation and validation of the Health Utilities Index Mark 3 into Spanish and correction norms for Spanish population]. *Med Clin (Barc).* 2003 Feb 1; 120(3):89–96. [PubMed: 12605729]
71. Hays RD, Brodsky M, Johnston MF, Spritzer KL, Hui KK. Evaluating the statistical significance of health-related quality-of-life change in individual patients. *Eval Health Prof.* 2005 Jun; 28(2): 160–171. [PubMed: 15851771]
72. McLeod LD, Coon CD, Martin SA, Fehnel SE, Hays RD. Interpreting patient-reported outcome results: US FDA guidance and emerging methods. *Expert Rev Pharmacoecon Outcomes Res.* 2011 Apr; 11(2):163–169. [PubMed: 21476818]

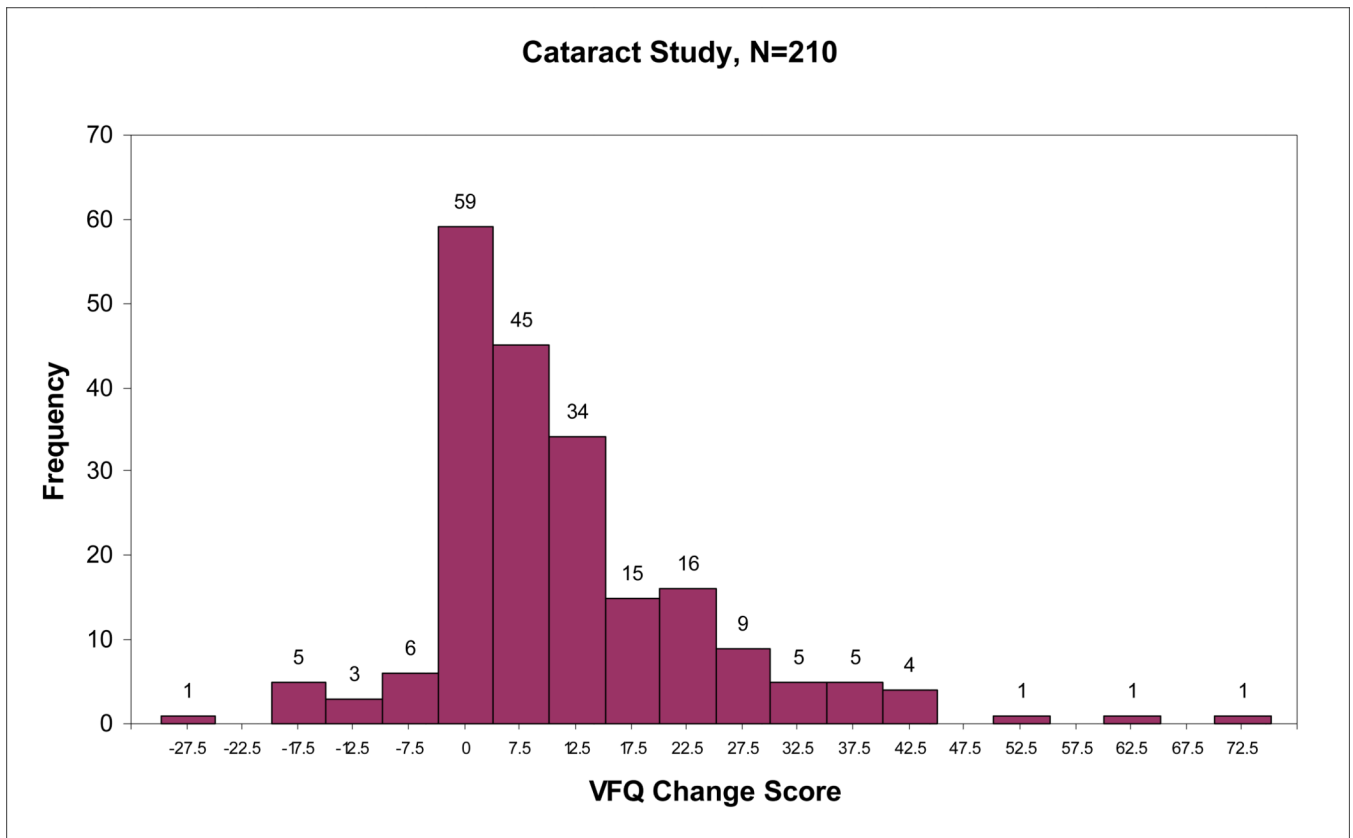


Figure 1. Distribution of Change in National Eye Institute Visual Function Questionnaire - 25 Total Scores

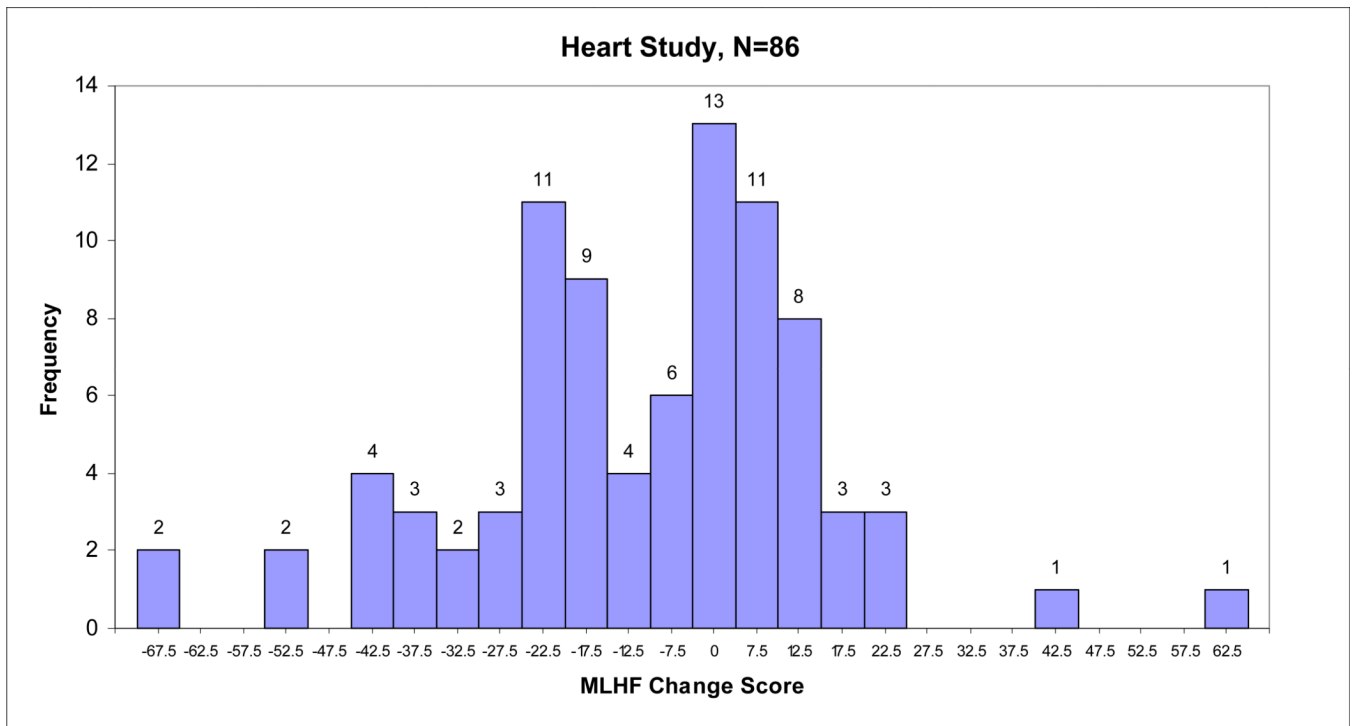


Figure 2.
Distribution of Change in Minnesota Living With Heart Failure Total Scores

Table 1

Demographic characteristics of the samples

	Cataract Patients			Heart Failure Patients			
	enrolled (n=376)	complete data – in sample (n=210)	incomplete data – not in sample (n=166)	enrolled (n=160)	Complete data – in sample (n=86)	incomplete data – not in sample (n=74)	in Sample vs not
Demographic							
Age:							
35–44	5 (1)	3 (1)	2 (1)	24 (15)	11 (13)	13 (18)	chi(2)=
45–64	115 (31)	71 (34)	44 (27)	101 (63)	63 (73)	38 (51)	8.96
65–91	256 (68)	136 (65)	120 (72)	35 (22)	12 (14)	23 (31)	Pr = 0.0113
Race:							
white	328 (87)	184 (88)	144 (87)	126 (79)	74 (86)	52 (70)	chi(3)=
black	12 (3)	7 (3)	5 (3)	19 (12)	8 (9)	11 (15)	6.01
Asian	19 (5)	13 (6)	6 (4)	5 (3)	1 (1)	4 (5)	Pr = 0.1083
other	4 (1)	2 (1)	2 (1)	2 (1)	2 (2)	0 (0)	
missing*	13 (3)	4 (2)	9 (5)	8 (5)	1 (1)	7 (9)	
Education:							
< HS	21 (6)	6 (3)	15 (9)	20 (13)	7 (8)	13 (18)	chi(6)=
HS graduate	60 (16)	32 (15)	28 (17)	45 (28)	25 (29)	20 (27)	10.55
some college	78 (21)	43 (20)	35 (21)	47 (29)	33 (38)	14 (19)	Pr = 0.1034
2-year assoc	27 (7)	14 (7)	13 (8)	12 (8)	7 (8)	5 (7)	
4-yr coll grad	90 (24)	56 (27)	34 (20)	16 (10)	7 (8)	9 (12)	
MA degree	57 (15)	38 (18)	19 (11)	9 (6)	3 (3)	6 (8)	
doctorate/professional	34 (9)	19 (9)	15 (44)	6 (4)	4 (5)	2 (3)	
missing*	9 (2)	2 (1)	7 (4)	5 (3)	0 (0)	5 (7)	
Female	222 (59)	124 (59)	98 (59)	52 (33)	26 (30)	26 (35)	chi(1)=
							0.44
							Pr = 0.5092

* missing not used in tests

Table 2

Baseline, One-month, and Change Scores, Cataract Cohort, n = 210

	Measure	Mean	Median	Std Dev	Minimum	Maximum
Baseline	VFQ _i	76.51	80.63	15.42	17.23	98.67
	EQ-5D	0.83	0.83	0.17	0.08	1.00
	HUI2	0.79	0.82	0.17	0.08	1.00
	HUI2 Sensation	0.76	0.76	0.14	0.00	1.00
	HUI3	0.66	0.69	0.27	-0.28	1.00
	HUI3 Vision	0.80	0.95	0.22	0.00	1.00
	QWB-SA	0.59	0.61	0.14	0.15	1.00
	SF-6D	0.74	0.75	0.12	0.33	1.00
	SRH	2.50	2.00	0.93	1.00	5.00
	VFQ _{ii}	0.86	0.92	0.12	0.41	0.97
One Month	VFQ _i	86.47	90.13	12.57	26.83	100.00
	EQ-5D	0.84	0.83	0.16	0.17	1.00
	HUI2	0.81	0.87	0.19	0.12	1.00
	HUI2 Sensation	0.84	0.87	0.17	0.00	1.00
	HUI3	0.72	0.80	0.28	-0.32	1.00
	HUI3 Vision	0.91	0.95	0.15	0.00	1.00
	QWB-SA	0.60	0.61	0.14	0.15	0.97
	SF-6D	0.73	0.74	0.12	0.39	1.00
	SRH	2.52	2.00	0.92	1.00	5.00
	VFQ _{ii}	0.90	0.94	0.09	0.46	0.97
Change	VFQ _i	9.96	7.86	13.44	-25.49	71.63
	EQ-5D	0.02	0.00	0.13	-0.51	0.54
	HUI2	0.02	0.03	0.14	-0.61	0.50
	HUI2 Sensation	0.08	0.00	0.19	-0.65	1.00
	HUI3	0.05	0.03	0.21	-0.82	0.77
	HUI3 Vision	0.12	0.00	0.22	-0.38	1.00

Measure	Mean	Median	Std Dev	Minimum	Maximum
QWB-SA	0.01	0.00	0.13	-0.63	0.39
SF-6D	-0.01	0.00	0.09	-0.29	0.22
SRH	0.02	0.00	0.69	-3.00	2.00
VFQ _{ii}	0.04	0.01	0.11	-0.24	0.53

Note: EQ-5D = 5 dimension Euro QOL measure; HUI2 = Health Utilities Index Mark 2; HUI3 = Health Utilities Index Mark 3; QWB-SA = Quality of Well Being – Self Administered Score; SF-6D = Short-Form 6D; SRH = self-rated health; VFQ_t = total score of National Eye Institute Visual Function Questionnaire (NEI-VFQ-25); VFQ_{ii} = Preference-based Score based on NEIVFQ-25.

Table 3

Baseline, Six-month, and Change Scores, Heart Failure Cohort, n = 86

	Measure	Mean	Median	Std Dev	Minimum	Maximum
Baseline	MLHF	48.26	48.50	25.40	0.00	101.00
	EQ-5D	0.77	0.79	0.18	0.18	1.00
	HUI2	0.76	0.85	0.22	0.14	1.00
	HUI3	0.62	0.73	0.32	-0.25	1.00
	QWB-SA	0.54	0.55	0.14	0.22	0.87
	SF-6D	0.63	0.63	0.11	0.39	0.93
	SRH	3.79	4.00	0.83	2.00	5.00
Six Months	MLHF	39.53	36.50	24.97	0.00	89.00
	EQ-5D	0.76	0.78	0.18	0.21	1.00
	HUI2	0.76	0.84	0.24	0.04	1.00
	HUI3	0.65	0.74	0.33	-0.34	1.00
	QWB-SA	0.58	0.59	0.15	0.21	1.00
	SF-6D	0.66	0.64	0.13	0.41	1.00
	SRH	3.49	4.00	0.98	1.00	5.00
Change	MLHF	-8.72	-7.50	22.12	-69.00	60.00
	EQ-5D	-0.01	0.00	0.18	-0.63	0.52
	HUI2	0.00	0.00	0.17	-0.68	0.43
	HUI3	0.03	0.00	0.23	-0.69	0.87
	QWB-SA	0.04	0.02	0.16	-0.45	0.46
	SF-6D	0.03	0.02	0.11	-0.17	0.42
	SRH	-0.30	0.00	0.90	-3.00	1.00

Note: EQ-5D = 5 dimension Euro QOL measure; HUI2 = Health Utilities Index Mark 2; HUI3 = Health Utilities Index Mark 3; MLHF = Minnesota Living with Heart Failure; QWB-SA = Quality of Well Being - Self Administered Scale; SF-6D = Short-Form 6D; SRH = self-rated health.

Table 4

Agreement Among 10 Measures, Cataract Cohort, n = 210

Pair	%Agreement	Kappa Statistic	95% Confidence Interval	Weighted Kappa Statistic	95% Confidence Interval
VFQ _i and EQ-5D	39	0.08	(-0.01, 0.16)	0.10	(0.01, 0.18)
VFQ _i and HUI2	48	0.11	(0.01, 0.21)	0.14	(0.04, 0.25)
VFQ _i and HUI2 Sensation	55	0.22	(0.11, 0.32)	0.22	(0.12, 0.32)
VFQ _i and HUI3	44	0.07	(-0.02, 0.17)	0.11	(0.01, 0.21)
VFQ _i and HUI3 Vision	57	0.25	(0.15, 0.35)	0.25	(0.15, 0.36)
VFQ _i and QWB-SA	40	0.06	(-0.02, 0.15)	0.12	(0.03, 0.21)
VFQ _i and SF-6D	39	0.09	(0.00, 0.17)	0.09	(0.00, 0.17)
VFQ _i and SRH	33	0.01	(-0.06, 0.08)	-0.05	(-0.12, 0.02)
VFQ _i and VFQ _{iii}	56	0.26	(0.17, 0.35)	0.33	(0.24, 0.42)

Note: EQ-5D = 5 dimension Euro QOL measure; HUI2 = Health Utilities Index Mark 2; HUI3 = Health Utilities Index Mark 3; QWB-SA = Quality of Well Being – Self Administered Scale; SF-6D = Short-Form 6D; VFQ_i = total score of National Eye Institute Visual Function Questionnaire (NEIVFQ-25); VFQ_{iii} = Preference-based Score based on NEIVFQ-25. Self-rated Health (SRH) was scored as 1 = poor; 2 = fair; 3 = good; 4 = very good; and 5 = excellent.

Table 5

Agreement Among Seven Measures, Heart Failure cohort, n = 86

Pair	% Agreement	Kappa Statistic	95% Confidence Interval	Weighted Kappa Statistic	95% Confidence Interval
MLHF and EQ-5D	19	-0.25	(-0.37, -0.13)	-0.30	(-0.45, -0.15)
MLHF and HUI2	29	-0.10	(-0.26, 0.05)	-0.19	(-0.36, -0.02)
MLHF and HUI3	26	-0.17	(-0.32, -0.02)	-0.23	(-0.40, -0.06)
MLHF and QWB-SA	22	-0.22	(-0.37, -0.07)	-0.30	(-0.46, -0.14)
MLHF and SF-6D	26	-0.11	(-0.25, 0.04)	-0.22	(-0.38, -0.06)
MLHF and SRH	49	0.25	(0.12, 0.39)	0.34	(0.19, 0.49)

Note: EQ-5D = 5 dimension Euro QOL measure; HUI2 = Health Utilities Index Mark 2; HUI3 = Health Utilities Index Mark 3; MLHF = Minnesota Living with Heart Failure; QWB-SA = Quality of Well Being – Self Administered Scale; SF-6D = Short-Form 6D. Self-rated Health (SRH) was scored as 1 = poor; 2 = fair; 3 = good; 4 = very good; and 5 = excellent.

Table 6

Comparisons of Change among Measures of Health-Related Quality of Life from Baseline to One Month in Cataract Surgery Cohort

	Got worse (n=15)	Stayed same (n=62)	Showed Improvement (n=133)	Row totals
EQ5D				
-	<u>4</u>	13	21	38
0	7	<u>39</u>	73	119
+	4	10	<u>39</u>	53
HUI2				
-	<u>8</u>	16	22	46
0	2	<u>20</u>	39	61
+	5	26	<u>72</u>	103
HUI3				
-	<u>8</u>	20	30	58
0	3	<u>14</u>	32	49
+	4	28	<u>71</u>	103
QWB-SA				
-	<u>11</u>	27	34	72
0	2	<u>12</u>	39	53
+	2	23	<u>60</u>	85
SF-6D				
-	<u>10</u>	20	43	73
0	3	<u>23</u>	42	68
+	2	19	<u>48</u>	69
SRH				
-	<u>1</u>	7	27	35
0	7	<u>46</u>	84	137
+	7	9	<u>22</u>	38

Note: - means Got Worse; 0 means Stayed Same; + means Showed Improvement.

Table 7

Comparisons of Change among Measures of Health-Related Quality of Life from Baseline to Six Months in Heart Failure Cohort

	Got Worse (n=46)	Stayed Same (n=13)	Showed Improvement (n=27)	Row totals
EQ5D				
-	<u>11</u>	6	16	33
0	16	<u>2</u>	8	26
+	19	5	<u>3</u>	27
HUI2				
-	<u>11</u>	6	13	30
0	8	<u>5</u>	5	18
+	27	2	<u>9</u>	38
HUI3				
-	<u>10</u>	8	14	32
0	10	<u>3</u>	4	17
+	26	2	<u>9</u>	37
QWB-SA				
-	<u>9</u>	6	15	30
0	7	<u>3</u>	5	15
+	30	4	<u>7</u>	41
SF-6D				
-	<u>9</u>	2	12	23
0	8	<u>7</u>	9	24
+	29	4	<u>6</u>	39
SRH				
-	<u>24</u>	3	4	31
0	20	<u>8</u>	13	41
+	2	2	<u>10</u>	14

Note: - means Got Worse; 0 means Stayed Same; + means Showed Improvement.