

IZA DP No. 8760

Survey Design and the Determinants of Subjective Wellbeing: An Experimental Analysis

Angus Holford
Stephen Pudney

January 2015

Survey Design and the Determinants of Subjective Wellbeing: An Experimental Analysis

Angus Holford

*ISER, University of Essex
and IZA*

Stephen Pudney

ISER, University of Essex

Discussion Paper No. 8760
January 2015

IZA

P.O. Box 7240
53072 Bonn
Germany

Phone: +49-228-3894-0

Fax: +49-228-3894-180

E-mail: iza@iza.org

Any opinions expressed here are those of the author(s) and not those of IZA. Research published in this series may include views on policy, but the institute itself takes no institutional policy positions. The IZA research network is committed to the IZA Guiding Principles of Research Integrity.

The Institute for the Study of Labor (IZA) in Bonn is a local and virtual international research center and a place of communication between science, politics and business. IZA is an independent nonprofit organization supported by Deutsche Post Foundation. The center is associated with the University of Bonn and offers a stimulating research environment through its international network, workshops and conferences, data service, project support, research visits and doctoral program. IZA engages in (i) original and internationally competitive research in all fields of labor economics, (ii) development of policy concepts, and (iii) dissemination of research results and concepts to the interested public.

IZA Discussion Papers often represent preliminary work and are circulated to encourage discussion. Citation of such a paper should account for its provisional character. A revised version may be available directly from the author.

ABSTRACT

Survey Design and the Determinants of Subjective Wellbeing: An Experimental Analysis^{*}

We analyse the results of experiments on questionnaire design and interview mode in the first four waves (2008-11) of the UK *Understanding Society* Innovation Panel survey. The randomised experiments relate to job, health, income, leisure and overall life-satisfaction questions and vary the labeling of response scales, mode of interviewing, and location of questions within the interview. We find significant evidence of an influence of interview mode and question design on the distribution of reported satisfaction measures, particularly for women. Results from the sort of conditional modeling used to address real research questions appear less vulnerable to design influences.

JEL Classification: C23, C25, C81, J28

Keywords: survey design, wellbeing, satisfaction, response bias, Understanding Society

Corresponding author:

Angus Holford
Institute for Social and Economic Research
University of Essex
Wivenhoe Park
Colchester, CO4 3SQ
United Kingdom
E-mail: ajholf@essex.ac.uk

^{*} This work was supported by the UK Longitudinal Studies Centre (award no. RES- 586-47-0002) and the European Research Council (project no. 269874 [DEVHEALTH]). We are grateful to Jon Burton, Anette Jäckle and Noah Uhrig for their help with the IP data. Thanks also to participants at the March 2012 LSE seminar and January 2014 AEA conference session on wellbeing measurement.

NON-TECHNICAL SUMMARY

In the effort to reduce the cost of collecting data and increase the number of people who respond, surveys employ an increasing variety of interview methods. These include face-to-face and telephone interviews, and paper, computer or web-based self-completion surveys. The interview mode constrains the format of questions that can be asked. It is not practical to force a respondent to listen to (for example) 7 possible answers to a question on the telephone, while it is straightforward to ask them to tick one box next to a range of options written in front of them.

The way you ask a question, and the format in which you require an answer, may have a big influence on the answer that you get. This is particularly the case for potentially sensitive questions on health and wellbeing.

Using data from randomised experiments in the Understanding Society Innovation Panel, we first describe the effects of survey design and interview mode on respondents' reported wellbeing. We show that even apparently minor differences in design, such as a vertical rather than horizontal list of possible responses in a private (self-completion) mode, generate large differences in reported wellbeing.

There is currently significant interest among policymakers in the determinants of subjective wellbeing. They should be concerned if policy prescriptions from research into this question are also sensitive to survey design and interview mode. Secondly therefore we investigate the impact of survey design and interview mode on the conclusions drawn from analysis of two specific research questions:

- (i) Do income and wages matter differently in determining the life, job and health satisfaction of men and women?
- (ii) How much additional income does a person need to offset the effects of a persistent health condition on their reported wellbeing?

We show some evidence that switching from a private mode (telephone or self-completion), to a public one (face-to-face), causes women to downplay the importance of income in determining their income and health satisfaction. We find no evidence for any effect of interview mode on the relative importance of health conditions and income in determining life satisfaction.

1 Introduction

Subjective assessments of wellbeing play an increasingly important role in applied economics. They underpin some approaches to economic evaluation, for example in assessment of health technologies and interventions (Ferrer-i-Carbonell and van Praag 2002); they have been proposed for legal assessment of compensation claims in the courts (Oswald and Powdthavee 2008), and they have been used as indicators of developmental outcomes and some aspects of the underlying non-cognitive skills emphasised by Heckman et al (2006). In political circles, wellbeing is also in the air. In France, President Nicolas Sarkozy's *Commission sur la Mesure de la Performance Économique et du Progrès Social* of 2008-9 (Stiglitz et al 2009) proposed measures of economic performance with wider scope than traditional indicators like GDP. In November 2010, the British Prime Minister announced plans for the Office for National Statistics to develop and publish official measures of wellbeing, observing that "prosperity alone can't deliver a better life" (Cameron 2010). Other national governments and international organisations including the OECD, Eurostat and the UN have made similar moves to extend the range of welfare indicators they produce. Much of the discussion about wellbeing measurement in the economics literature has focused on conceptual issues relating to distinctions between satisfaction, happiness, positive and negative affect, etc, distinctions between different domains of wellbeing (Stiglitz and Fitoussi, 2013), and validating international comparisons of the level and determinants of wellbeing (Kapteyn et al 2013). More practical questions about survey design for subjective questions have long been discussed in the survey methods literature, but have not featured much in the economic debate over wellbeing measurement.

Everyone knows that the way you ask a question may have a big influence on the answer that you get, and subjective questions on health and wellbeing are no exception to this. Perhaps more important than the question itself is the form in which you require the answer – and evidence of difficulty in interpreting response scales has been found in qualitative work

by Ralph et al (2011). Although there has been interest in reliability issues for subjective wellbeing data (see, for example, Kristensen and Westergaard-Nielsen 2007 and Krueger and Schkade 2008), the economics literature on happiness and wellbeing has devoted little attention to the influence of survey design and context and its possible implications for data analysis. However, Conti and Pudney (2011) analysed quasi-experimental evidence arising from variations over time in the design of job satisfaction questions in the British Household Panel Survey (BHPS), finding evidence of a substantial influence of the design of questions and response scales, the mode of interview and the interview context on the distribution of reported satisfaction. One of the most striking aspects of the BHPS evidence is the significantly different impact that survey design features have on the response behaviour of men and women and the consequent large distortions that may be induced in research findings on gender differences in the determinants of job satisfaction. The lack of robustness of gender differences was also found by Studer (2012), who analysed a randomised experiment comparing continuous and discrete rating scales in a Dutch web-based panel survey.

Although choice of question design and interview mode have been examined by survey methodologists, their conclusions are typically limited to simple indicators of data quality and impacts on summary statistics like means and sample proportions. But, in practice, measures of wellbeing are used in much more sophisticated ways, for conditional modeling and comparison over time and across groups within society. The political drivers of the move to measurement are quite clear about this: “If anyone was trying to reduce the whole spectrum of human happiness into one snapshot statistic, I would be the first to roll my eyes [...]. But that’s not what this is all about” (Cameron 2010).

In the economics literature, claims for the validity of subjective wellbeing indicators are often based on the loose prediction argument that these variables are correlated in the way one would expect with observable variables and with subsequent behavioural outcomes. The

predictive criterion is an important one, but not very powerful. It is easy to construct examples where an indicator has measurement error serious enough to cause catastrophic biases for the sort of analysis economists are interested in, but is sufficiently highly correlated with the ‘true’ variable to satisfy the requirement for predictive correlation. Our aim in this paper is to add to the evidence on measurement reliability, by using a set of randomised experiments to investigate the impact of various dimensions of survey design in the context of a major UK survey that is used for wellbeing analysis: the UK Household Longitudinal Survey (UKHLS), also known as *Understanding Society*, which is the successor to the BHPS. This paper extends Conti and Pudney’s (2011) BHPS analysis of job satisfaction by considering also self-reported health and several other domains of life satisfaction, and by using experimental control and a wider range of variations in survey design.

2 The UKHLS Innovation Panel

The UKHLS is a very large-scale household panel survey which has absorbed the long-established BHPS sample. Its design differs in many ways from that of the BHPS and one of its innovative features is a sub-panel, known as the *Innovation Panel* (IP), reserved exclusively for experimental work. The IP experiments constitute one of the very few examples of experimentation in survey design sustained over a substantial number of waves of a panel study. There is an annual competition open to all researchers to propose new experiments.¹ The IP sample of 1500 households was drawn in Spring 2008 and has been re-interviewed annually since then. There is a relatively high attrition rate (McFall et al 2013), so a refreshment sample of 500 households was added at wave 4 in Spring 2011. The core content of the IP interview is similar to the UKHLS main survey, but there are considerable variations in content from wave to wave. Each wave carries several experiments: during waves 1-4, an annual average of three procedural experiments (interview mode, incentives, etc)

¹The German Socio-Economic Panel (GSOEP) has also now developed an IP sub-panel of the GSOEP.

and five measurement experiments (questionnaire format, wording, etc). As a consequence, the observed outcomes are affected by multiple interventions and the complicated nature of the sample that results is a disadvantage of the IP concept. However, each experiment is randomised, so cross-contamination between experiments can be expected to contribute variance rather than bias. The major advantage of the IP is that experiments take place within the context of a large, continuing panel survey, so the IP is believed to be superior in terms of external validity to the small special-purpose experiments which are typical of much of the survey methods literature.

Experimental design: treatment groups

Fieldwork for waves 1-4 of the IP were conducted in April-June of each of the four years 2008-2011, and included experiments to investigate the impact of question design, interview mode and positioning of questions within the interview, using random assignment of households to treatment groups. All individuals within each household received the same experimental treatment. In this study, we focus on questions covering five satisfaction domains: a 4-item module covering satisfaction with health, income, leisure and life overall, and a separate job satisfaction question. Each measure involves a 7-point response scale, as used in the BHPS from 1991 to 2008. This excludes all wave 1 health/income/leisure/life satisfaction data and a random half of the wave 1 job satisfaction responses, which used an 11-point scale.²

Wave 1 A random half of the sample received a job satisfaction question with a 7-point scale. (Question wordings are given in the next sub-section). The question was delivered orally by an interviewer following a computer-generated script without use of a showcard. Verbal equivalents were given only for the two polar points on the scale (“completely dissatisfied” and “completely satisfied”).

²See Burton et al (2008) for a comparison of IP responses using 7-point and 11-point scales.

Wave 2 The wave 2 design was a composite of two separate randomisations. First, households were assigned in the ratio 2:1 to Computer-Assisted Telephone Interviewing (CATI) or face-to-face interviewing (F2F) during an interviewer visit to the home. During F2F interviews, most questions were delivered by Computer-Assisted Personal Interviewing (CAPI), but Computer-Assisted Self-Interviewing (CASI) was used for the satisfaction module in a randomly-assigned subgroup. There were also independent assignments to treatment groups formed by varying question design and position of the question within the interview. As part of the design, for assignments that would have resulted in a requirement to administer by CATI a question that was in fact infeasible by telephone (because it required a show-card or reading of a long list of allowable responses), the closest feasible approximation to the allocated treatment was substituted. In some cases telephone contact was unsuccessful, in which case some or all members of the household were instead interviewed F2F, if that proved possible. There were no variations in question position within the interview for job satisfaction, so there were 8 treatment groups at wave 2 for job satisfaction, rather than the 14 for other domains.

Wave 3 Wave 3 was conducted entirely F2F, so there were no CATI groups. Apart from that, the set of experimental treatments was identical to those used at wave 2, but a group rotation was used to generate temporal variation in treatments. There was also a further – unintended – experiment at wave 3, since an error in the CAPI script led to some questions being repeated in a different format within the same interview. The impact of that inadvertent repetition is examined in the final part of section 3.

Wave 4 For the 4-item satisfaction module, the wave 4 experiment was a simple comparison of two private modes: a paper self-completion (PSC) questionnaire and a CASI version. The job satisfaction question was administered to all employed respondents in standard CAPI mode.

TABLE 1
Experimental treatment groups: satisfaction variables with 7-point response scale

	Treatment group	Question placement	Sample size								
			Wave 1 <i>ITT*</i>	Wave 1 <i>Actual†</i>	Wave 2 <i>ITT*</i>	Wave 2 <i>Actual†</i>	Wave 3 <i>ITT*</i>	Wave 3 <i>Actual†</i>	Wave 4 <i>ITT*</i>	Wave 4 <i>Actual†</i>	
<i>Satisfaction with Health, Income, Leisure and Life overall</i>											
1	Visit: CASI full labels	Late	-	-	126	167	135	123	1,208	1,060	-
2	Visit: CASI polar labels	Late	-	-	110	136	306	286	-	-	-
3	Visit: showcard full labels	Late	-	-	47	55	291	270	-	-	-
4	Visit: oral 2-stage	Late	-	-	59	80	286	260	-	-	-
5	Visit: showcard polar labels	Late	-	-	48	69	424	391	-	-	-
6	Visit: oral polar labels	Late	-	-	63	76	314	287	-	-	-
7	Phone: CATI 2-stage	Late	-	-	404	297	-	-	-	-	-
8	Phone: CATI polar labels	Late	-	-	402	308	-	-	-	-	-
9	Visit: showcard full labels	Early	-	-	58	65	-	-	-	-	-
10	Visit: oral 2-stage	Early	-	-	56	69	-	-	-	-	-
11	Visit: showcard polar labels	Early	-	-	52	56	-	-	-	-	-
12	Visit: oral polar labels	Early	-	-	59	74	-	-	-	-	-
13	Phone: CATI 2-stage	Early	-	-	188	140	-	-	-	-	-
14	Phone: CATI polar labels	Early	-	-	198	161	-	-	-	-	-
15	Visit: PSC full labels	Separate	-	-	-	-	-	-	1,177	855	-
<i>Job satisfaction</i>											
1	Visit: CASI full labels	Mid	-	-	76	101	159	140	-	-	-
2	Visit: CASI polar labels	Mid	-	-	61	78	170	161	-	-	-
3	Visit: showcard full labels	Mid	-	-	60	61	169	155	1,380	1,247	-
4	Visit: oral 2-stage	Mid	-	-	64	87	159	149	-	-	-
5	Visit: showcard polar labels	Mid	-	-	51	67	155	141	-	-	-
6	Visit: oral polar	Mid	731	672	64	79	158	140	-	-	-
7	Phone: CATI 2-stage	Mid	-	-	331	246	-	-	-	-	-
8	Phone: CATI polar labels	Mid	-	-	322	256	-	-	-	-	-

* Intention-to-treat: assigned treatment † Treatment received by responders

The longitudinal pattern of treatments affecting the IP satisfaction modules is detailed in Table 1, together with potential sample numbers on an intention-to-treat (ITT) and an actual response basis. To avoid difficult selection issues, all the results presented here are on an ITT basis, although there is very little difference in the findings when the analysis is repeated using actual treatments, owing to the high compliance rate in most groups. The remaining experiments carried by the IP in waves 1-4 were irrelevant to our objectives.

Experimental design: question wording and response scales

Questions were asked sequentially for three aspects of life satisfaction, using (for all except groups 4, 7, 10 and 13) the following question stem:

*How dissatisfied or satisfied are you with the following aspects of your situation:
(a) your health; (b) the income of your household; (c) the amount of leisure time
you have.*

These three domain-specific questions were followed by an overall assessment:

Using the same scale, how dissatisfied or satisfied are you with your life overall?

For groups 3 and 9, a fully-labeled showcard specified response options in a vertical list ordered from top to bottom as: *7 Completely satisfied; 6 Mostly satisfied; 5 Somewhat satisfied; 4 Neither satisfied nor dissatisfied; 3 Somewhat dissatisfied; 2 Mostly dissatisfied; 1 Completely dissatisfied.* For groups 1 and 2, the questions were administered by the more private Computer-Assisted Self-Interviewing (CASI) method, and the seven alternatives were displayed horizontally across the screen of a laptop computer for selection directly by the respondent. Polar-point labeled variants of the question (groups 5, 6, 8, 11, 12 and 14) omitted the textual labels from options 2 to 6. If the polar-labeled response scale was communicated orally (groups 6, 8, 12 and 14), explanations of the two extreme points were read out by the interviewer.

Groups 4, 7, 10 and 13 received a question with a fully-labeled response scale, designed to be deliverable by telephone, when use of a showcard or reading of a full list of responses would have been impractical. The question has a two-stage structure:

(i) *How dissatisfied or satisfied are you with [...]? Would you say you are...
(1 Dissatisfied; 2 Neither Dissatisfied nor Satisfied; 3 Satisfied)*

(ii) *[If dissatisfied or satisfied...] And are you Somewhat, Mostly or Completely
[Satisfied / Dissatisfied] with [...]? (1 Somewhat; 2 Mostly; 3 Completely)*

At wave 2, questions on satisfaction with health, family income, leisure and life overall were either asked early (about 25% of the way through the interview, following a block of questions on transport mode choices) or late (about 95% of the way through the interview, following questions on political affiliation and values). At waves 3 and 4, these questions were always asked late, except for group 15 at wave 4, where the questions were contained in a paper self-completion questionnaire completed during the interviewer's visit.

People in employment or self-employment were also asked about their job satisfaction: shortly after mid-interview, following a section dealing with employment or self-employment details, including occupation, hours and earnings. The question stem is:

*All things considered, which number best describes how dissatisfied or satisfied
you are with your job overall?*

The same 1-7 response scale and labeling options were used as for the single-stage life satisfaction questions, and groups 4 and 7 received a 2-stage variant.

3 The impact on data distributions

We are mainly interested in the effect of survey design on the answers to substantial empirical research questions but we first look for evidence that the distribution of responses to questions

on subjective wellbeing are influenced by aspects of design. Table 2 gives the mean responses for each treatment group, separately for men and women, together with wave-specific chi-square tests comparing the distribution of responses from each treatment group with the distribution in the pooled sample.³ There is some evidence of an impact of the set of experimental variations, but the small group sizes mean that tests for individual treatment groups have low power. Table 3 gives results from overall tests of joint significance for the whole set of experimental variations, by survey wave and domain of satisfaction. In view of earlier findings on gender differences, we carry out these tests separately for men and women. Let Y_d be the satisfaction score for the d th domain. We use Monte Carlo permutation versions of two tests; an ANOVA F -statistic for Y_d , and a chi-square test for the equality of the vector of proportions in each response category with the pooled sample proportions. See Good (2006) for a review of permutation tests and Heckman et al (2010) for an application to experimental evaluation.

The test results of Table 3 show that design effects are frequently significant, although the pattern of effects is unexpected. Experimental variation at wave 2 produces large impacts on the response distributions: for women, they are significant in all five domains at the 5% level using an ANOVA permutation F test and one using the chi-square test for equality of response proportions. For men, we find a significant ANOVA test in three domains and a p-value of 0.063 or lower in all cases, besides some evidence at the 10% level for chi-square test. But at wave 3, where group sizes are larger and we would expect better power, there is less evidence of an effect. This is especially the case for women, with a solitary rejection at the 5% level using the chi-square test and only two at the 10% level in the ANOVA. For men there are two rejections at the 5% level using the chi-square test, but only a single 5% rejection (satisfaction with health) in the ANOVA test, with no rejections at all by the H -statistic.

³Here, the test statistic $\chi^2 = (\hat{p} - \bar{p})' V^{-1} (\hat{p} - \bar{p})$ where \hat{p} is the 7×1 vector of observed response probabilities, \bar{p} is the empirical distribution of response probabilities across all groups, and V^{-1} is the generalised inverse of an estimated covariance matrix. A multi-group generalisation of this statistic is used for Table 3.

TABLE 2
Mean satisfaction scores and permutation tests for equality of response distribution to pooled sample

Treatment group	Place- ment	Women				Men					
		Health	Income	Leisure	Life	Job	Health	Income	Leisure	Life	Job
<i>Wave 2</i>											
CASI full labels	Late	5.05***	4.81	4.94	5.15**	5.11	5.53***	4.58**	4.56**	5.37	4.97
CASI polar labels	Late	5.15**	4.61	4.67	5.20**	4.71**	5.05	5.07**	5.05	5.40	5.48
Showcard full labels	Late	5.33***	4.56	5.26	6.07	5.74	5.88	5.35	5.41	6.35	5.35
Oral 2-stage	Late	6.00***	5.22	5.42	6.07	5.46	5.13	4.56	5.48	5.48	5.35
Showcard polar labels	Late	4.52	5.30	5.04	5.43*	5.38	5.32	4.64	4.14	5.45	4.65
Oral polar labels	Late	4.52	4.77	4.92	5.26	4.90	5.24	4.44	4.44	5.36	4.64
CATI 2-stage	Late	5.28	4.92*	5.32*	5.75*	5.64***	5.19*	4.54*	5.13***	5.54	5.45***
CATI polar labels	Late	5.27	4.73***	4.89**	5.70	5.37**	5.40	4.67**	5.16**	5.45	5.27*
Showcard full labels	Early	5.52**	5.33	5.41	6.07	.	5.14	4.87	5.30	5.57	.
Oral 2-stage	Early	4.96	5.08	5.20	5.72	.	5.76	5.14	5.19	6.24	.
Showcard polar labels	Early	5.79	4.25	4.86	5.54	.	4.94	5.00	5.06	5.50	.
Oral polar labels	Early	4.91**	4.27**	4.41***	5.18	.	5.00	4.13	5.23	5.45	.
CATI: 2-stage	Early	5.10	4.66	5.19*	5.82**	.	5.56	5.00	4.99	5.54**	.
CATI: polar labels	Early	5.15***	4.90**	5.14**	5.42***	.	5.54	4.96	5.10*	5.59**	.
Overall mean		5.21	4.78	5.06	5.51	5.40	5.35	4.74	5.05	5.53	5.26
<i>Wave 3</i>											
CASI full labels	Late	5.30	4.46**	4.97**	5.50	5.06***	4.92	4.72***	4.53**	5.40*	4.81
CASI polar labels	Late	5.15	4.96	4.91	5.49*	5.16	5.38**	4.94	4.97	5.41	4.87
Showcard full labels	Late	5.29***	4.88*	5.03**	5.75**	5.51**	5.59**	4.90**	5.16	5.66	5.07
Oral 2-stage	Late	5.08***	4.93*	5.20	5.69	5.49	5.21**	4.74**	5.08	5.69	5.32*
Showcard polar labels	Late	5.35	4.88	4.97	5.61	5.24	5.36*	4.93	5.14	5.67	5.27
Oral polar	Late	5.52	4.96*	5.21	5.73	5.35***	5.45	4.95	5.21*	5.64	5.13
Overall mean		5.29	4.88	5.05	5.63	5.30	5.36	4.88	5.07	5.60	5.07
<i>Wave 4</i>											
CASI full labels	Separate	4.67***	4.44	4.67	5.10**	.	4.70***	4.56*	4.74	5.12	.
Paper self-completion	Separate	4.95***	4.64	4.86	5.28**	.	5.12***	4.71*	4.85	5.24	.
Overall mean		4.80	4.53	4.76	5.21	.	4.88	4.62	4.78	5.17	.

Statistical significance stars from chi-square test for equality of vector of proportions in each response category with pooled sample proportions, Monte Carlo permutation with 10000 replications: *** = 1%, ** = 5%, * = 10%

TABLE 3
Permutation test *P*-values for joint hypothesis of no treatment effects

Satisfaction domain	Women			Men		
	<i>Wave 2</i>	<i>Wave 3</i>	<i>Wave 4</i>	<i>Wave 2</i>	<i>Wave 3</i>	<i>Wave 4</i>
Health	0.086 <i>0.000</i>	0.049 <i>0.089</i>	0.000 <i>0.329</i>	0.843 <i>0.047</i>	0.046 <i>0.027</i>	0.003 <i>0.008</i>
Income	0.131 <i>0.037</i>	0.457 <i>0.182</i>	0.232 <i>0.310</i>	0.474 <i>0.029</i>	0.063 <i>0.889</i>	0.089 <i>0.884</i>
Leisure	0.860 <i>0.024</i>	0.200 <i>0.269</i>	0.193 <i>0.207</i>	0.778 <i>0.022</i>	0.398 <i>0.105</i>	0.525 <i>0.716</i>
Life overall	0.027 <i>0.000</i>	0.245 <i>0.139</i>	0.016 <i>0.291</i>	0.063 <i>0.063</i>	0.355 <i>0.190</i>	0.448 <i>0.677</i>
Job	0.617 <i>0.000</i>	0.150 <i>0.052</i>	.	0.489 <i>0.063</i>	0.016 <i>0.373</i>	.

All p-values from Monte Carlo permutation with 10,000 replications. **Bold:** p-value for chi-square test statistic for equality of vector of response proportions with pooled sample proportions; *Italic:* p-value for ANOVA *F*-statistic.

A possible interpretation of the weaker effect at wave 3 is linked to the rotation of treatment groups between waves 2 and 3. Since almost every wave 3 respondent had responded via a different mode or question design a year earlier, the recollection of that response may have nullified the effect of treatment at wave 3 – which would be consistent with Pudney’s (2008, 2011) findings of dynamic contamination of responses to a different subjective wellbeing question in the BHPS. If that explanation is accepted, then it casts doubt on the validity of observed measures of change in wellbeing in panel data.

Specific design aspects

The experimental treatment groups differ in a number of dimensions, and tests of the impact of specific aspects of design (rather than combinations of aspects) are more informative. Table 4 reports the results of extending simple ANOVA comparisons to the health, income, leisure and life satisfaction domains, using a seemingly unrelated regressions approach allowing unrestricted correlation between the four satisfaction scores. The analysis is restricted to wave 2 (the analogous estimates for wave 3 show little impact) and focuses on two interview mode contrasts (CASI v. F2F and CATI v. F2F) and two question design contrasts (polar-point v. full labeling of the response scale and 2-stage v. 1-stage question design). The analysis is applied within subsamples which have approximately the same composition

in terms of all other experimental aspects for the two groups being compared, so that there should be negligible compositional bias in the comparisons reported. For comparison, we also include analogous single-equation results for the smaller group of employed respondents who are asked a separate job satisfaction question. For each panel of Table 4, the first four rows show the mean effects on domain-specific satisfaction scores; the fifth row gives a joint P -value for the joint hypothesis that all four mean effects are zero. The clearest evidence from these joint tests is for CASI rather than F2F interviewing and for 2-stage rather than 1-stage question design, but both results apply only to female respondents. Evidence on job satisfaction shows the same pattern.

TABLE 4
IP wave 2: Impacts on mean responses of specific design aspects

Satisfaction domain	<i>Interview mode</i>		<i>Response scale design</i>	
	<i>CASI</i> [‡]	<i>CATI</i> [†]	<i>Polar-point</i> [∇]	<i>2-Stage</i> [★]
<i>Women</i>				
Health	-0.225 (0.222)	0.109 (0.170)	-0.044 (0.222)	0.132 (0.122)
Income	0.069 (0.221)	0.008 (0.174)	-0.438** (0.219)	0.152 (0.125)
Leisure	-0.340 (0.223)	0.173 (0.185)	-0.328 (0.223)	0.374*** (0.133)
Overall	-0.627*** (0.164)	0.123 (0.135)	-0.253 (0.168)	0.255*** (0.097)
Joint P -value [◇]	0.0003	0.8020	0.2386	0.0244
n	227	727	227	727
Job [★]	-0.654*** (0.267)	0.334 (0.210)	-0.385 (0.272)	0.319** (0.149)
<i>Men</i>				
Health	0.008 (0.219)	0.112 (0.176)	-0.401* (0.216)	-0.036 (0.126)
Income	-0.158 (0.237)	0.174 (0.188)	0.147 (0.236)	0.031 (0.134)
Leisure	-0.191 (0.268)	0.047 (0.205)	-0.093 (0.267)	0.031 (0.147)
Overall	-0.334* (0.187)	-0.088 (0.153)	-0.172 (0.188)	0.124 (0.109)
Joint P -value [◇]	0.4009	0.5873	0.0696	0.7078
n	177	603	177	603
Job [♠]	0.211 (0.318)	0.401* (0.224)	-0.005 (0.313)	0.261* (0.157)

Standard errors in parentheses. Significance: *** = 1%; ** = 5%; * = 10%.

[‡] Comparison with F2F interview + showcard: based on treatment groups 1-3, 5, 9, 11.

[†] Comparison with F2F oral (no showcard): based on treatment groups 4, 6-8, 10, 12-14.

[∇] Comparison with fully-labeled scale: based on treatment groups 1-3, 5, 9, 11.

[★] Comparison with 1-stage question design: based on treatment groups 4, 6-8, 10, 12-14.

[◇] SURE generalisation of the ANOVA test allowing for responses correlated across domains.

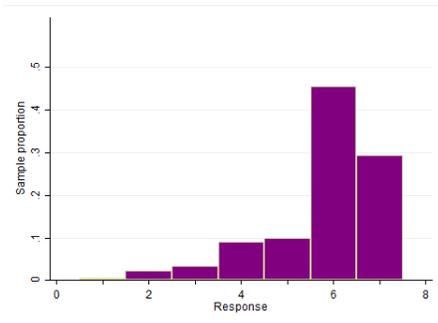
[♠] Single-equation ANOVA test for subset of employed/self-employed individuals.

We find no evidence that early or late positioning of questions within the questionnaire causes any significant shifts in the response distribution. This is in contrast to the large questionnaire context effects found in some other survey applications (Schuman and Presser 1981, Tourangeau 1999) and the evidence of respondent fatigue which may affect responses late in the interview (Herzog and Bachman 1981, Helgeson and Ursic 1994). Note that we do not investigate the ordering of individual questions within the satisfaction module – something that has been found to influence respondents’ interpretation of satisfaction questions (Schwarz et al 1991, Tourangeau et al 1991). We now expand on the effects of response scale and interview mode on response distributions with reference to visual evidence.

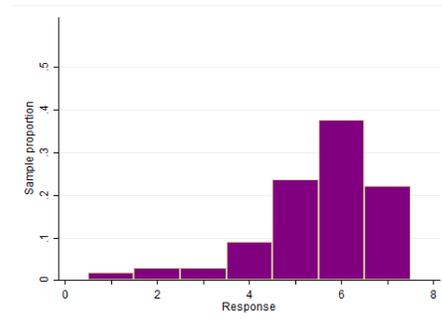
2-stage versus 1-stage questions

There has been some debate about the use of two-stage branching (or unfolding) question structures. Some authors find these designs yield better reliability (Krosnick and Berent 1993),⁴ while others found that some respondents have difficulty interpreting the question appropriately without access to the full range of allowed responses (Hunter 2005, p.10-11). Comparing the 2-stage design with 1-stage alternatives in Table 2, we find higher mean scores for the 2-stage design in 16 out of 18 cases for women and 12 out of 18 for men. Table 3 shows that these differences are statistically significant for women (leisure and life overall) but not men. Figure 1 shows the empirical response distributions and suggest that the main effect of the 2-stage design is to move responses from the $Y = 5$ category to $Y = 6, 7$, thus raising the mean score. There is little evidence of any difference between the 1-stage and 2-stage designs at wave 3.

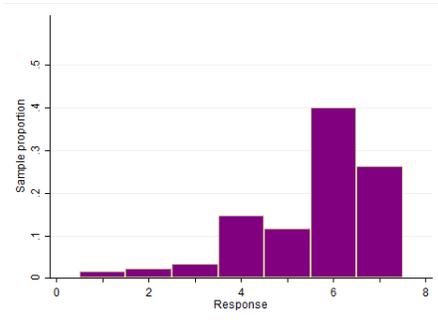
⁴Note that differences in question structure were confounded with labeling differences in the Krosnick-Berent study of test-retest reliability. We would also argue that test-retest reliability should be seen as a measure of consistency over time rather than ‘reliability’.



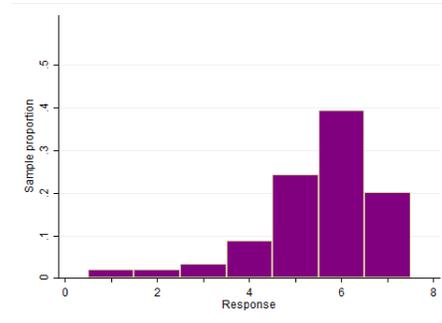
(a) 2-stage design: women ($n = 379$)



(b) 1-stage designs: women ($n = 641$)



(c) 2-stage design: men ($n = 328$)



(d) 1-stage designs: men ($n = 522$)

Figure 1 Wave 2 sample distributions for life satisfaction: 2-stage vs. 1-stage question designs

Polar-point versus full labels

Unlike Weng (2004) and Conti and Pudney (2011), there is only weak evidence of an impact of polar point rather than full labeling of the response scale (Table 3). Its impact on mean scores is negative in most cases (Table 4), resulting from a shift from responses at $Y = 6$ to $Y = 5$ (Figure 2). This effect is surprising, given our expectation that exclusive labeling of extreme points would attract responses to those extremes.

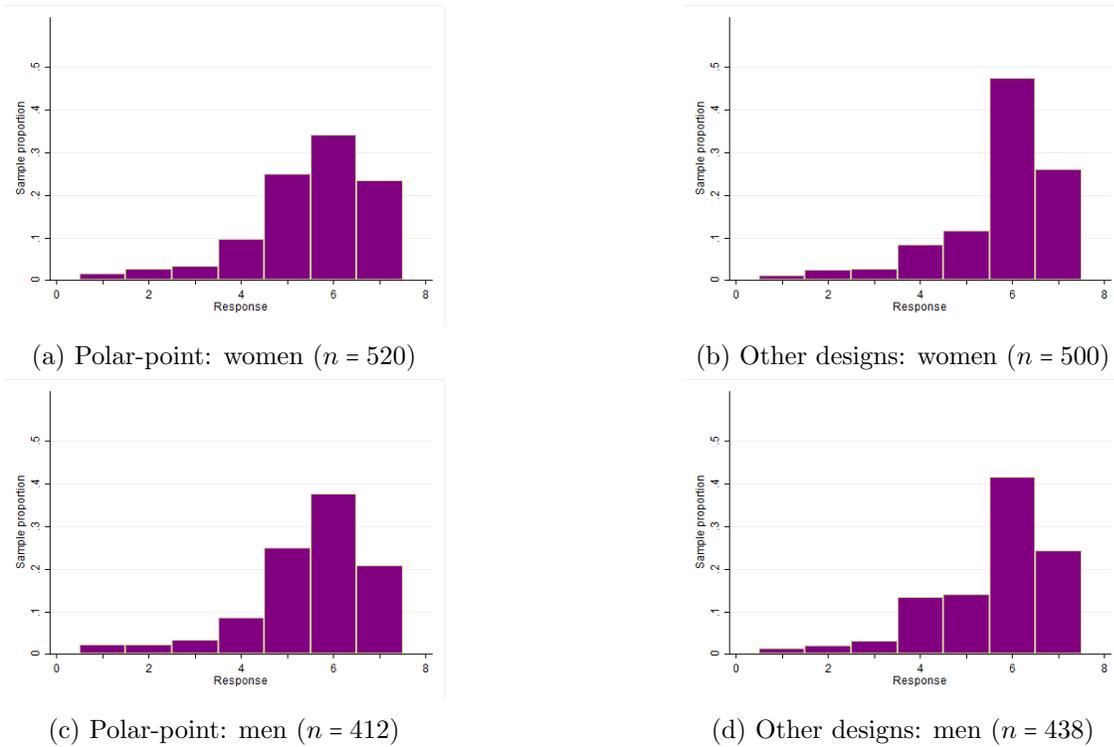


Figure 2 Wave 2 sample distributions for life satisfaction: Polar-point vs. other question designs

CASI versus *F2F*

At waves 2 and 3, Table 2 suggested a definite pattern for CASI compared to other more public interview modes: looking across all five satisfaction domains and both forms of CASI, 18 of the 20 mean scores are below the overall average for women and 14 of 20 are below average for men. Figure 3 compares the distributions for CASI responses to the life satisfaction question with other one-stage F2F designs at wave 2. The distributions are dominated by a mode at $Y = 6$, which is a general feature of categorical responses to satisfaction questions, possibly reflecting an aversion to extremes, as suggested by Studer (2012). The comparison of CASI with other designs suggests a shift of mass from $Y = 6$ and 7 to $Y \leq 4$: overall, CASI increases the sample proportion of $Y \leq 4$ from 16% to 23% and reduces the sample proportion of $Y \geq 6$ from 61% to 52%.

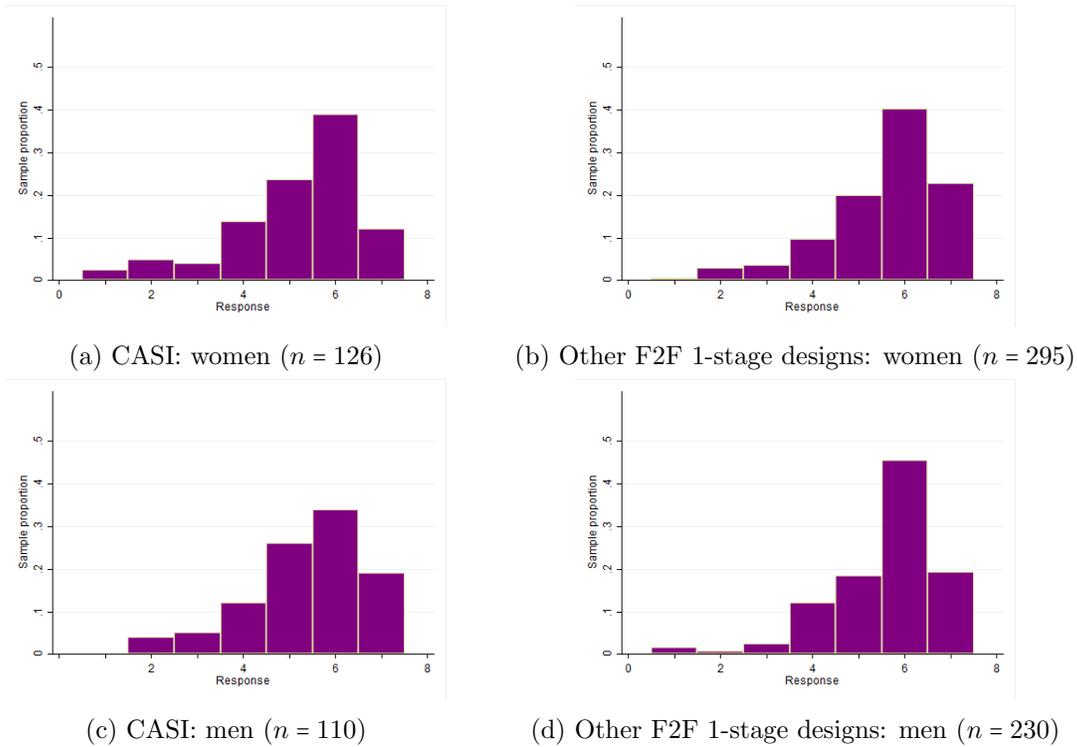


Figure 3 Wave 2 sample distributions for life satisfaction: CASI vs. 1-stage F2F

It is not a simple matter to interpret these mode effects, since they involve differences in several dimensions, including the format of visual display of the response scale (Jenkins and Dillman 1997), the degree of respondent privacy and presence of an outsider (the interviewer). Privacy and the social desirability of alternative responses are especially important for sensitive issues (Hochstim 1967, De Leeuw 1992, Aquilino 1997) and a further important factor may be a desire by some individuals to maintain a bargaining position within the family, rendering some satisfaction questions sensitive in oral interviews where other family members may be within earshot (Conti and Pudney 2011).⁵

⁵For example, Appendix Table A1, detailing context-specific permutation tests for the effects of certain design aspects, shows strong evidence of an impact of CASI against the F2F mode for the health and leisure domains (men) and income (women) in both the full and polar-labeled contexts in wave 3. We argue this represents a privacy effect, with men reluctant to express public dissatisfaction with their health or leisure and women reluctant to voice concerns about income. In all of these cases, CASI delivers a significantly lower mean satisfaction score.

Comparing private modes: CASI *versus* Paper self-completion

At wave 4, there is a significant effect for CASI rather than PSC, especially for satisfaction with health among male respondents, for whom CASI produces a much smaller mean response (4.70) than PSC (5.12). Figure 4 shows the wave 4 response distributions for satisfaction with health, by gender and interview mode. Compared with PSC, CASI has the effect of transferring probability mass to categories $Y = 1$ and 2 , from $Y = 6$ in particular. This reduces the mean score, but also changes the mass of the lower tail, which has implications for the common applied practice of using binary indicators of low satisfaction. The impact on the response distribution is surprising because CASI and PSC are both private modes designed to do essentially the same thing: shield the respondent from social pressures during interview. Assuming they both achieve that aim, the remaining difference between them must presumably relate to the way in which the response scale is conveyed on the computer screen or paper questionnaire and then interpreted by the respondent. However, both use the same fully-labeled response scale. In CASI they are displayed vertically from $1 =$ completely dissatisfied at the top of the screen to $7 =$ completely satisfied at the bottom, whereas the paper questionnaire displays them horizontally from 1 at the left to 7 at the right. The significant differences we find are consistent with the warning from Christian et al (2009) that the visual design of response scales can have a significant influence on responses. It is likely to become a particular issue in future multi-mode surveys which have difficulty in avoiding endogenous selection from the set of interview modes, each of which has a distinct ‘look’.

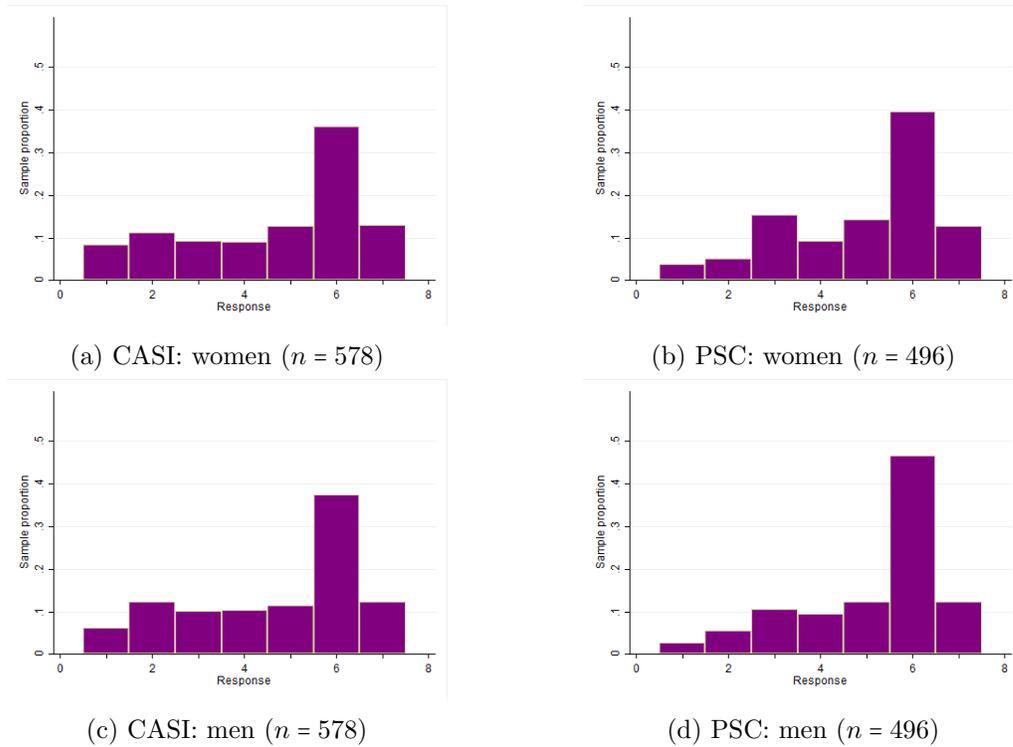


Figure 4 Wave 4 sample distributions for satisfaction with health: CASI vs. PSC

Repeated measures within wave 3

Some respondents at wave 3 received the health, income, leisure and life satisfaction questions in two different forms within the same interview. This was an error in programming the CAPI script, rather than a deliberate experiment, but it offers a direct opportunity to assess the effects of different treatments on the same set of people. This allows us to compare responses to different designs more efficiently than through random assignment to single treatments. However, if the fact of repetition changes behaviour directly, or if there are significant question order or respondent fatigue effects, the results will be confounded to some extent.

Four groups received repeated questions. Group I received the 2-stage question, delivered orally early in the interview and then the single-stage question with fully-labeled showcard about 20 minutes later on average. Group II received the single-stage question orally, with verbal descriptions of the two extreme points, then later the same question using a polar-labeled showcard. Groups III and IV had the same treatments as I and II in reverse order.

The top panel of Table 5 gives correlations between early and late scores and estimates of the mean differences. The test-retest correlations are in the range 0.57-0.86, which is rather higher than the range of correlations for life satisfaction quoted by Andrews and Whithey (1976), Kammann and Flett (1983) and Krueger and Schkade (2008), who used longer retest intervals but unchanged question design. If we make classical measurement error assumptions, the correlation between early and late measures gives the usual measure of test-retest reliability as the share of measurement error in total variance: implying a range of values of 0.16-0.75 for the noise/signal ratio (of the measurement error variance to the variance of the ‘true’ variable).

We investigate differences in the early and late mean scores and in the proportion of high scores ($Y \geq 6$ or $Y = 7$). For respondent i , the satisfaction score at time $t = 0$ (early) or 1 (late) is Y_{igt} , where $g = \text{group I, II, III or IV}$. At time t , members of groups I and II receive treatment sequences a, b and b, a respectively, while members of groups III and IV receive sequences c, d and d, c , where a, b, c, d denote the oral 2-stage question, fully-labeled showcard, oral polar-labeled question and polar-labeled showcard respectively. Assume additive effects:

$$Y_{igt} = \mu_0 + (\mu_a - \mu_b)\xi_{igt}^a + (\mu_c - \mu_d)\xi_{igt}^c + \mu_R(1 - t) + \varepsilon_{igt} \quad (1)$$

where ξ_{igt}^a, ξ_{igt}^c are indicators of receiving treatments a and c respectively, and the ε_{it}^g are mutually independent zero mean measurement errors. The coefficient $(\mu_a - \mu_b)$ is the effect of using a 2-stage question structure rather than a showcard, $(\mu_c - \mu_d)$ is the effect of delivering the polar-labeled response question orally rather than by showcard, and μ_R is the effect of repetition. We estimate the coefficients by least squares random effects regression; the results are presented in the last panel of Table 5. We see significant effects for health satisfaction only, where use of the oral 2-stage question raises reported satisfaction relative to fully-labeled showcards, and question repetition has a positive effect of a similar magnitude.

TABLE 5
Repeated measures in IP wave 3

<i>Treatment sequence</i>	Satisfaction domain			
	<i>Health</i>	<i>Income</i>	<i>Leisure</i>	<i>Life</i>
<i>Test-retest Pearson correlation coefficients</i>				
Oral 2-stage question → Fully labeled showcard [‡]	0.665	0.707	0.672	0.571
Fully labeled showcard → Oral 2-stage question [†]	0.770	0.737	0.745	0.591
Polar labeled oral → Polar labeled showcard [†]	0.743	0.786	0.686	0.723
Polar labeled showcard → Polar labeled oral [•]	0.708	0.860	0.817	0.749
<i>Mean scores: random effects regression</i>				
Oral 2-stage question v. Fully labeled showcard: $(\mu_a - \mu_b)$	0.147* (0.076)	-0.003 (0.077)	0.118 (-0.085)	0.031 (0.064)
Polar labeled oral v. Polar labeled showcard: $(\mu_c - \mu_d)$	0.082 (0.103)	0.056 (0.096)	-0.120 (0.108)	-0.049 (0.083)
Repetition effect: μ_R	0.143* (0.077)	0.056 (0.078)	-0.006 (0.087)	0.012 (0.064)
Sample size $n =$	512	503	511	512

[‡] $n = 124$; [†] $n = 117$; [•] $n = 153$. Test statistics based on robust standard errors.

4 Survey design and satisfaction models

The demand for data is a derived demand – we are interested in data only because of the research results that can be produced from them. Much of the survey methods literature ignores this fundamental point and restricts consideration of the impact of design features to the statistical reliability of relatively simple summary measures computed from the data. Instead, most applied researchers are interested in the statistical relationships between variables, using models which represent complex conditional distributions in the data. In the research literature on wellbeing, this type of modeling takes the form of relationships between satisfaction as a dependent variable and a set of covariates describing the individual’s characteristics and circumstances in some detail (see Van Praag and Ferrer-i-Carbonell 2004, and Clark et al 2008 for surveys). Typical analysis methods include fixed-effects regression and random-effects ordered probit. We apply these modeling approaches and investigate the impact of experimental variations in survey design on the estimates.

It is no simple matter to assess the impact of a set of experimental design variations on these complex analyses. With 15 treatment groups and models involving 20 or more coefficients for both genders over five satisfaction domains, there are at least 3,000 experimental

effects to be estimated in the most general approach. We resolve this ‘curse of dimensionality’ by focusing on the answers to specific research questions rather than model parameters. In this section, we consider two issues: first, the possible gender difference in pecuniary influences on wellbeing; and, second, the magnitude of the compensating income variation which would be required to offset the wellbeing effects of a persistent health condition. In both cases, we investigate the effect of using F2F interviewing rather than other more private modes.

Two single-equation model specifications are used, both based on the following latent regression:

$$Y_{it}^* = \mathbf{x}_{it}\boldsymbol{\beta} + \mathbf{x}_{it}^+\zeta_{it}\boldsymbol{\gamma} + u_i + \varepsilon_{it} \quad (2)$$

Here Y_{it}^* is the (latent) satisfaction score, \mathbf{x}_i is the full vector of covariates, \mathbf{x}_i^+ is the subset of covariates of particular interest for a particular research question and ζ_i is a dummy indicating cases featuring a specific design aspect of interest. u_i and ε_{it} are unobservables. Let \mathbf{T}_{it} be a vector of indicators describing the design treatment experienced by individual i at time t ; the observed score Y_{it} is then related to Y_{it}^* and \mathbf{T}_{it} in alternative ways by the two models:

(i) *Fixed-effects (FE) regression*: $Y_{it} = Y_{it}^* + \mathbf{T}_{it}\boldsymbol{\alpha}$ and u_i is eliminated by removing within-group means.

(ii) *Generalised random-effects ordered (GREO) probit*: $Y_{it} = r$ iff $A_{it}^{r+1} \geq Y_{it}^* > A_{it}^r$ where the threshold parameters are linear functions of design aspects: $A_{it}^r = \mathbf{T}_{it}\boldsymbol{\alpha}^r$.

Gender and the income-wellbeing relation

A common finding in the literature on job satisfaction is that the pecuniary aspects of a job are less important to women than to men (see, for example, Booth and van Ours 2008). This was called into question by Conti and Pudney (2011), whose results suggested that responses from women interviewed F2F were subject to bias and that the gender difference

largely disappeared when data from a more private PSC questionnaire were used instead. Table 6 explores this for satisfaction with income (waves 2-4), job (waves 1-4) and health (also waves 2-4) satisfaction.

TABLE 6
Gender-income-design interactions in three satisfaction models

Coefficient	Satisfaction with income		Job satisfaction		Satisfaction with health	
	<i>GREO</i>	<i>FE</i>	<i>GREO</i>	<i>FE</i>	<i>GREO</i>	<i>FE</i>
	<i>probit</i>	<i>regression</i>	<i>probit</i>	<i>regression</i>	<i>probit</i>	<i>regression</i>
	<i>Coefficients[‡]</i>					
Female	-0.283 (0.250)	- -	0.418 (0.372)	- -	-0.614*** (0.231)	- -
Income	0.629*** (0.062)	0.162 (0.115)	0.198* (0.119)	0.070 (0.185)	0.003 (0.058)	-0.092 (0.122)
Female × Income	0.112 (0.083)	0.161 (0.157)	-0.117 (0.157)	0.004 (0.232)	0.188** (0.076)	0.148 (0.165)
Female × F2F	0.859** (0.380)	0.699 (0.461)	0.129 (0.417)	-0.135 (0.519)	1.105*** (0.364)	1.018** (0.485)
Income × F2F	0.079 (0.096)	0.030 (0.115)	-0.076 (0.127)	-0.241 (0.154)	0.175*** (0.092)	0.093 (0.121)
Female × Income × F2F	-0.259** (0.126)	-0.181 (0.152)	-0.038 (0.176)	0.064 (0.218)	-0.343*** (0.122)	0.309* (0.160)
	<i>Joint tests of design effects: P-values</i>					
Additive design effects [†]	0.0000	0.0437	0.0000	0.0000	0.0000	0.0000
F2F interactions	0.0687	0.1304	0.6545	0.2679	0.0194	0.1086

Standard errors in parentheses. Significance: * = 10%; ** = 5%; *** = 1%. ‡ *Income* is log equivalised gross household income for the income and health satisfaction equations and log hourly earnings for job satisfaction. Other covariates included in the model are: age, age², single/widowed/divorced, no. children, non-white, wave dummies. Health satisfaction models only: Non-disabling and disabling health conditions.

[†] Design aspects in T_{it} are: CASI, CATI, Polar-labeled, 2-stage design, F2F.

In both the GREO and FE models, for all three satisfaction measures, the additive design variables are jointly significant at the 5% level. The FE regressions show no further design impacts and, indeed, no significant income effect or gender-income interaction at all. For the GREO models however, there is some evidence of a design interaction which could affect the empirical picture of gender differences in relation to money as a contributor to wellbeing, but only for income and health satisfaction.

In the GREO probit model for income satisfaction, the use of F2F rather than private interview modes seems to have two gender-specific effects: a large general increase in the

levels of satisfaction reported by women; and a significant reduction in the female \times income coefficient from 0.112 to -0.147. In other words, switching from private CASI to public F2F modes causes women significantly, on average, to downplay the importance of income in determining their income satisfaction. Both effects are individually significant at the 5% level, although jointly, the whole set of F2F-interactions are only significant at the 7% level. The same interpretation can be made from the health satisfaction model, where F2F mode reduces the female \times income coefficient from 0.188 to -0.155, with the whole set of F2F-interactions this time significant at the 2% level. These results are consistent with Conti and Pudney’s (2011) findings for BHPS job satisfaction data, although the smaller sample sizes here reduce the statistical clarity somewhat.

Compensating variations for health conditions

Statistical models of wellbeing have often been used to estimate the income variation equivalent to events or resources like marriage, divorce, childbirth, unemployment and social capital (for example, Blanchflower and Oswald 2004, Di Tella and MacCulloch 2008 and Groot et al 2007). In health, the same approach has been used by Ferrer-i-Carbonell and van Praag (2002), Groot and Maassen van den Brink (2004, 2006), Mentzakis (2011), Zaidi and Burchardt (2005) and Morciano et al (2013) to estimate the personal costs of disease and disability. We have argued elsewhere (Hancock et al 2013) that this indirect method of constructing an estimate of the compensating variation (CV) as a by-product of a parametric model of wellbeing, is particularly sensitive to even minor misspecifications, often giving huge overestimates. Hancock et al (2013) argue for a more stable direct nonparametric approach, but indirect parametric estimation of the CV remains standard practice and so we examine the impact of survey design on it. We consider linear and quadratic models of overall life satisfaction, based on the latent regression (2), with the leading terms of the linear index specified as $\mathbf{x}_{it}\boldsymbol{\beta} = \beta_1 H_{1it} + \beta_2 H_{2it} + \phi(M_{it}) + \dots$, where: H_{1it} is a binary indicator of the existence of a “long-standing health condition” that is not reported to

cause any disability; H_{2it} indicates such a condition with associated disability; M_{it} is annual gross household income (in £'000) per equivalent adult; and $\phi(M_{it}) = \beta_3 M_{it}$ or $\beta_3 M_{it} + \beta_4 M_{it}^2$. In these two cases, the CV for health state $H_j (j = 1, 2)$ is $-\beta_j/\beta_3$ (linear model) or $-(B + \sqrt{B^2 - 4\beta_j\beta_4})/2\beta_4$ (quadratic model), where $B = \beta_3 + 2\beta_4 M_{it}$.⁶

TABLE 7
Compensating income variations in two satisfaction models

	Linear in income		Quadratic in income	
	<i>Coefficients (standard errors)[‡]</i>			
Income (£'000 p.a. per equivalent adult)	0.0078***	(0.0016)	0.0182***	(0.0031)
Income ²	.	.	-0.0001***	(0.00003)
Non-disabling health condition	-0.219***	(0.068)	-0.222***	(0.068)
Disabling health condition	-0.461**	(0.058)	-0.453***	(0.058)
Income × F2F	-0.0011	(0.0027)	-0.0027	(0.0064)
Income ² × F2F	.	.	0.00002	(0.00008)
Non-disabling condition × F2F	-0.085	(0.117)	-0.085	(0.117)
Disabling condition × F2F	-0.110	(0.092)	-0.119	(0.092)
<i>Joint tests of design effects: P-values</i>				
Additive design effects [†]	0.0000		0.0000	
F2F interactions	0.5785		0.7451	
<i>Estimated compensating variations (standard errors), £'000 p.a. per equivalent adult[▲]</i>				
Non-disabling condition (not F2F)	27.95***	(10.58)	120.33***	(27.60)
Non-disabling condition (F2F)	34.19**	(14.95)	142.1	(141.9)
P- value for difference	0.718		0.878	
Disabling condition (not F2F)	58.80***	(14.96)	85.66*	(51.2)
Disabling condition (F2F)	64.22***	(21.73)	105.4	(188.9)
P- value for difference	0.821		0.917	

[‡] Other covariates included in the model are: age, age², single/widowed/divorced, no. children, non-white, retired, wave dummies. [†] Design aspects in ξ_{it} are: CASI, CATI, Polar-labeled, 2-stage design, F2F. [▲] CV estimates at mean income for the quadratic model.

Table 7 reports GREO probit estimates of the disability and income coefficients, and their interactions with the F2F interview mode. Again, additive design effects are highly significant, but here we are unable to detect any interaction between interview mode and health or income. Consistent with Hancock et al's (2013) findings, the implied CV estimates are extremely large, even for a non-disabling health condition: almost £28,000 for the linear

⁶Log income is often used in applied work, giving a CV of the form $M_{it} \exp\{-\beta_j/\beta_3\}$. This tends to produce even less robust CV estimates than the linear or quadratic income models and we do not report the results here.

model and – quite implausibly – £120,000 at mean income for the better-fitting quadratic model. The F2F interaction raises these large values still further, but the increase is not statistically significant.

5 Conclusions

There are three reasonably clear conclusions from our analysis of the wave 1-4 experiments in the UKHLS Innovation Panel, a couple of puzzling results, and some implications for the design of multi-wave experiments in large longitudinal surveys.

First, there is strong overall evidence that the choice of interview mode and question/response scale design has a detectable influence on the distribution of responses to questions on subjective health and wellbeing. This is particularly true for computer-assisted self-interviewing (CASI) relative to other interview modes and there is some, weaker, evidence of an influence for the way the response scale is designed.

Second, the evidence for an influence of design features – especially interview mode – is stronger for female respondents than for males. This is consistent with evidence from other sources, and suggests a greater degree of sensitivity to the social context of the interview for women than men on average.

Our third conclusion is more important for the purposes of econometric analysis. We have taken two research questions as examples to assess the practical importance of these design effects: (*i*) Is there a gender difference in the impact of pecuniary factors on expressed wellbeing? (*ii*) What income variation is equivalent in wellbeing terms to a persistent health condition? We find that the answer to question (*i*) is influenced by the use of face-to-face (F2F) rather than more private modes of interview, with (after controlling for a wide range of other characteristics) women tending to give higher and less strongly income-related assessments of satisfaction with income only when F2F interviewing is used. For research

question (ii), we found no evidence for any effect of interview mode on the tradeoff between income and health, and therefore no impact on compensating income differentials. Despite the significant effects that we have found, on this evidence it seems fair to say that, with the possible exception of gender effects, the sort of conditional modelling used in economics seems more robust with respect to design differences than are simpler unconditional summary statistics.

But there are some puzzles accompanying these conclusions. At wave 3, which involved a more powerful comparison between fewer treatment groups, the evidence for design effects was actually weaker than at wave 2 – a finding which could possibly be explained in part by the ‘contamination’ of current responses by recalled past responses, as found by Pudney (2008, 2011). A second puzzle is that, at wave 4 where the comparison was between two relatively private interview modes (CASI and paper self-completion questionnaire), there was a large significant mean difference between responses, with CASI producing lower ratings of wellbeing. Given the similarity of the degree of privacy of those two modes, visual differences in response scale (e.g. vertical rather than horizontal presentation) may be involved in the impact that CASI appears to have.

Finally, resources like the UKHLS Innovation Panel are (arguably) a good way of ensuring that experiments are relevant to the reality of large-scale surveys but there is a risk that the resulting multiplicity of experiments within a moderately-sized sample may reduce power and complicate the interpretation of experimental effects, unless the complex of experiments can be designed in an integrated way. The problem of designing multiple experiments spanning multiple waves of a panel survey has not been studied systematically and it is not clear that the UKHLS Innovation Panel used in this paper has yet found a good way of managing the process of experimental design. Although randomised, the multi-treatment experiments considered here were confined to three or four waves and are arguably less effective in revealing framing and mode effects than the longer-term (and unplanned) BHPS experiment

exploited by Conti and Pudney (2011), which involved sustained question repetition with different interview modes.

References

- [1] Andrews F.M. and Withey, S.B. (1976). *Social Indicators of Wellbeing: Americans' Perceptions of Life Quality*. New York: Plenum Press.
- [2] Aquilino, W. S. (1997). Privacy effects on self-reported drug use: interactions with survey mode and respondent characteristics. In Harrison L. and Hughes A. (eds.) *The Validity of Self-Reported Drug Use: Improving the Accuracy of Survey Estimates*, 383-415, Rockville: National Institute on Drug Abuse, NIDA Research Monograph 167.
- [3] Blanchflower, D. G. and Oswald, A.J. (2004). Well-being over time in Britain and the USA, *Journal of Public Economics* **88**, 1359-1386.
- [4] Booth, A.L. and van Ours, J.C. (2008). Job satisfaction and family happiness: the part-time work puzzle, *Economic Journal* **118**, F77-F99.
- [5] Burton, J., Laurie, H. and Uhrig, S. C. N. (eds.) (2008). *Understanding Society* Some preliminary results from the wave 1 *Innovation Panel*. University of Essex: Understanding Society Working Paper no. 2008-03.
- [6] Cameron, D. (2010). Speech on wellbeing, London 25 November 2010 <http://www.number10.gov.uk/news/pm-speech-on-well-being/> (accessed 8 October 2013).
- [7] Clark, A., Frijters, P. and Shields, M. (2008). Relative income, happiness and utility: an explanation for the Easterlin paradox and other puzzles, *Journal of Economic Literature* **46**, 95-144.
- [8] Conti, G. and Pudney, S.E. (2011). Survey design and the analysis of satisfaction, *Review of Economics and Statistics* **93**, 1087-1093.
- [9] Christian, L.M., Parsons, N.L. and Dillman, D.A. (2009). Designing scalar questions for web surveys, *Sociological Methods and Research* **37**, 393-425.
- [10] De Leeuw, E. (1992). *Data quality in mail, telephone and face to face surveys*. Amsterdam: TT Publications.
- [11] Di Tella, R., and MacCulloch, R.J. (2008). Gross national happiness as an answer to the Easterlin paradox? *Journal of Economic Development* **86**, 22-42.
- [12] Ferrer-i-Carbonell, A. and van Praag, B.M.S. (2002). The subjective costs of health losses due to chronic diseases. An alternative model for monetary appraisal, *Health Economics* **11**, 709-722.
- [13] Good, P. I. (2006). *Resampling Methods: A Practical Guide to Data Analysis* (3rd edition). Basel: Birkhäuser.
- [14] Groot, W. and Maassen van den Brink, H. (2004). A direct method for estimating the compensating income variation for severe headache and migraine *Social Science and Medicine* **58**, 305-314.

- [15] Groot, W. and Maassen van den Brink., H. (2006). The compensating income variation of cardiovascular disease, *Health Economics* **15**, 1143-1148.
- [16] Groot, W., Maassen van den Brink., H. and van Praag, B.M.S. (2007). The compensating income variation of social capital. University of Munich: CESIFO Working Paper no. 1889. *Health Economics* **11**, 709-722.
- [17] Hancock, R.M., Morciano, M. and Pudney, S.E. (2013). Nonparametric estimation of a compensating variation: the cost of disability, University of Essex: ISER Working Paper 2013-26.
- [18] Heckman, J.J., Stixrud, J. and Urzua, S. (2006). The effects of cognitive and noncognitive abilities on labor market outcomes and social behavior. *Journal of Labor Economics* **24**, 411-482.
- [19] Heckman, J.J., Moon, S.H., Pinto, R., Savelyev, P. and Yavitz, A. (2010). Analyzing social experiments as implemented: a reexamination of the evidence from the HighScope Perry preschool program, *Quantitative Economics* **1**, 1-46.
- [20] Helgeson, J.G. and Ursic, M.L. (1994). The role of affective and cognitive decision-making processes during questionnaire completion, *Public Opinion Quarterly*, **11**, 493-510.
- [21] Herzog, A.R. and Bachman, J.G. (1981). Effects of questionnaire length on response quality, *Public Opinion Quarterly*, **45**, 549-559.
- [22] Hochstim, J. (1967). A critical comparison of three strategies of collecting data from households, *Journal of the American Statistical Association*, **62**, 976-989.
- [23] Holford, A.J. and Pudney, S.E. (2013). The *Understanding Society Innovation Panel*: Notes on the construction of a gross annual household income variable for waves 1-4. Mimeo, University of Essex.
- [24] Hunter, J. (2005). Cognitive Test of the 2006 NRFU: Round 1. Washington DC: US Bureau of the Census, Study Series Report (Survey Methodology no.2005-07).
- [25] Jenkins, C.R. and Dillman, D.A. (1997). Towards a theory of self-administered questionnaire design. In Lyberg, L., Biemer, P., Collins, M., de Leeuw, E., Dippo, K., Schwarz, N. and Trewin, D. (eds.) *Survey Measurement and Process Quality*. New York: Wiley.
- [26] Kammann, N.R. and Flett, R. (1983). Affectometer 2: A scale to measure current level of general happiness. *Australian Journal of Psychology*, **35**, 259-265.
- [27] Kapteyn, A., Smith, J.P. and Van Soest, A. (2013). Are Americans Really Less Happy with Their Incomes? *Review of Income and Wealth* **59**, 44-65.
- [28] Kristensen, N. and N. Westergaard-Nielsen (2007). Reliability of job satisfaction measures, *Journal of Happiness Studies* **8**, 273-292.
- [29] Krosnick, J.A., and Berent, M.K. (1993). Comparisons of party identification and policy preferences: the impact of survey question format, *American Journal of Political Science*, **37** (3), 941-964.

- [30] Krueger, A.B. and D.A. Schkade (2008). The reliability of subjective well-being measures, *Journal of Public Economics*, **92**, 1833-1845.
- [31] McFall, S., Burton, J., Jäckle, A., Lynn, P. and Uhrig, S.C.N. (2013). Understanding Society: The UK Household Longitudinal Study Innovation Panel, Waves 1-5, User Manual. University of Essex: Institute for Social and Economic Research (<https://www.understandingsociety.ac.uk/documentation/innovation-panel>, accessed 30 Sep 2013).
- [32] Mentzakis, E. (2011). Allowing for heterogeneity in monetary subjective wellbeing valuations, *Health Economics* **20**, 331-347
- [33] Morciano, M., Hancock, R.M. and Pudney, S.E. (2013). Disability costs and equivalence scales in the older population in Great Britain, *Review of Income and Wealth* forthcoming.
- [34] Oswald, A.J. and Powdthavee, N. (2008). Does happiness adapt? A longitudinal study of disability with implications for economists and judges, *Journal of Public Economics* **92**, 1061-1077.
- [35] Pudney, S.E. (2008). The dynamic consistency of responses to survey questions on wellbeing, *Journal of the Royal Statistical Society Series A* **171**, 21-40.
- [36] Pudney, S.E. (2011). Perception and retrospection: the dynamic consistency of responses to survey questions on wellbeing, *Journal of Public Economics* **95**, 300-310.
- [37] Ralph, K., Palmer, K. and Olney, J. (2011). Subjective well-being: a qualitative investigation of subjective well-being questions. London: Office for National Statistics, research report.
- [38] Schuman, H. and Presser, S. (1981). *Questions and Answers in Attitude Surveys: Experiments in Question Form, Wording and Context*. New York: Academic Press.
- [39] Schwarz, N., Strack, F. and Mai, H. (1991). Assimilation and contrast effects in part-whole question sequences: a conversational logic analysis, *Public Opinion Quarterly*, **55**, 3-23.
- [40] Stiglitz, J., Sen, A. and Fitoussi, J.-P. (2009). Report by the Commission on the Measurement of Economic Performance and Social Progress.
- [41] Stiglitz, J., and Fitoussi, J.-P. (2013). On the measurement of social progress and well-being: some further thoughts, *Global Policy* **4**, 290-293.
- [42] Studer, R. (2012). Does it matter how happiness is measured? Evidence from a randomised controlled experiment, *Journal of Economic and Social Measurement* **37**, 317-336.
- [43] Tourangeau, R. (1999). Context effects on answers to attitude questions, in Sirken, M.G., Herrmann, D.J., Schechter, S., Schwarz, N., Tanur, J.M. and Tourangeau, R. (eds.) *Cognition and Survey Research*. New York: Wiley. data. *Social Science and Medicine* **57**, 1621-1629.

- [44] Tourangeau, R., Rasinski, K.A., and Bradburn, N. (1991), Measuring happiness in surveys: a test of the subtraction hypothesis, *Public Opinion Quarterly*, **55**, 255-266.
- [45] Van Praag, B.M.S and Ferrer-i-Carbonell, A. (2004). *Happiness Quantified. A Satisfaction Calculus Approach*. Oxford: Oxford University Press.
- [46] Weng, L.-J. (2004). Impact of the number of response categories and anchor labels on coefficient alpha and test-retest reliability, *Educational and Psychological Measurement* **64**, 956-972.
- [47] Zaidi, A. and Burchardt, T. (2005). Comparing incomes when needs differ: equalization for the extra costs of disability in the UK. *Review of Income and Wealth* **51**, 89-114.

Appendix A: Additional tables

TABLE A1
P-values for permutation tests on specific design aspects

Context	Wave	Women				Men					
		Health	Income	Leisure	Life	Job	Health	Income	Leisure	Life	Job
Showcard full labels	2	0.075	0.292	0.712	0.752	0.256	0.613	0.836	0.078		
		<i>0.744</i>	<i>0.074</i>	<i>0.787</i>	<i>1.000</i>	<i>0.128</i>	<i>0.339</i>	<i>0.851</i>	<i>0.012</i>		
Oral 2-stage	2	0.001	0.685	0.933	0.603	0.650	0.653	0.804	0.700		
		<i>0.028</i>	<i>0.772</i>	<i>0.620</i>	<i>0.216</i>	<i>0.245</i>	<i>0.285</i>	<i>0.611</i>	<i>0.079</i>		
Showcard polar	2	0.156	0.664	0.848	0.169	0.937	0.621	0.554	0.959		
		<i>0.019</i>	<i>0.940</i>	<i>0.757</i>	<i>0.809</i>	<i>0.465</i>	<i>0.570</i>	<i>0.149</i>	<i>1.000</i>		
Oral polar	2	0.696	0.422	0.648	0.894	0.937	0.640	0.274	0.379		
		<i>0.416</i>	<i>0.384</i>	<i>0.375</i>	<i>0.870</i>	<i>0.587</i>	<i>0.557</i>	<i>0.170</i>	<i>0.809</i>		
CATI 2-stage	2	0.339	0.314	0.283	0.565	0.381	0.281	0.046	0.999		
		<i>0.410</i>	<i>0.234</i>	<i>0.575</i>	<i>0.701</i>	<i>0.086</i>	<i>0.046</i>	<i>0.573</i>	<i>1.000</i>		
CATI polar labels	2	0.122	0.007	0.150	0.067	0.257	0.658	0.074	0.448		
		<i>0.492</i>	<i>0.397</i>	<i>0.246</i>	<i>0.076</i>	<i>0.505</i>	<i>0.163</i>	<i>0.822</i>	<i>0.497</i>		
Full labels versus polar labels											
CASI	2	0.002	0.048	0.624	0.261	0.062	0.030	0.640	0.845	0.503	0.638
		<i>0.732</i>	<i>0.513</i>	<i>0.394</i>	<i>0.845</i>	<i>0.338</i>	<i>0.098</i>	<i>0.111</i>	<i>0.184</i>	<i>0.939</i>	<i>0.252</i>
CASI	3	0.524	0.062	0.432	0.361	0.032	0.105	0.015	0.218	0.306	0.465
		<i>0.532</i>	<i>0.038</i>	<i>0.800</i>	<i>1.000</i>	<i>0.726</i>	<i>0.075</i>	<i>0.427</i>	<i>0.131</i>	<i>1.000</i>	<i>0.840</i>
Showcard	2	0.126	0.558	0.023	0.139	0.214	0.158	0.009	0.309	0.107	0.212
		<i>0.550</i>	<i>0.047</i>	<i>0.247</i>	<i>0.002</i>	<i>0.293</i>	<i>0.411</i>	<i>0.440</i>	<i>0.046</i>	<i>0.123</i>	<i>0.176</i>
Showcard	3	0.013	0.446	0.122	0.035	0.015	0.018	0.840	0.261	0.531	0.471
		<i>0.727</i>	<i>1.000</i>	<i>0.749</i>	<i>0.310</i>	<i>0.241</i>	<i>0.183</i>	<i>0.906</i>	<i>0.946</i>	<i>0.958</i>	<i>0.450</i>
Two-stage versus polar-labeled questions											
Oral	2	0.000	0.230	0.485	0.370	0.119	0.211	0.813	0.544	0.587	0.705
		<i>0.022</i>	<i>0.065</i>	<i>0.079</i>	<i>0.007</i>	<i>0.226</i>	<i>0.370</i>	<i>0.134</i>	<i>0.171</i>	<i>0.132</i>	<i>0.162</i>
Oral	3	0.003	0.006	0.543	0.649	0.030	0.492	0.115	0.012	0.160	0.481
		<i>0.026</i>	<i>0.911</i>	<i>0.976</i>	<i>0.806</i>	<i>0.562</i>	<i>0.235</i>	<i>0.325</i>	<i>0.579</i>	<i>0.753</i>	<i>0.492</i>
CATI	2	0.099	0.025	0.030	0.055	0.001	0.129	0.073	0.002	0.033	0.013
		<i>0.979</i>	<i>0.741</i>	<i>0.032</i>	<i>0.121</i>	<i>0.001</i>	<i>0.340</i>	<i>0.611</i>	<i>0.730</i>	<i>0.714</i>	<i>0.296</i>

All p-values from Monte Carlo permutation with 10,000 replications. **Bold:** p-value for chi-square test statistic for equality of vector of response proportions with pooled sample proportions; *Italic:* p-value for ANOVA F-statistic

TABLE A1 (continued)
P-values for permutation tests on specific design aspects

Context	Women						Men					
	Wave	Health	Income	Leisure	Life	Job	Health	Income	Leisure	Life	Job	
Full labels	4	0.000	0.242	0.212	0.018	0.018	0.005	0.110	0.554	0.483	.	
		<i>0.011</i>	<i>0.057</i>	<i>0.071</i>	<i>0.184</i>	.	<i>0.000</i>	<i>0.190</i>	<i>0.374</i>	<i>0.251</i>	.	
CASI versus F2F												
Full labels	3	0.616	0.359	0.200	0.516	0.023	0.135	0.408	0.052	0.525	0.759	
		<i>1.000</i>	<i>0.083</i>	<i>0.816</i>	<i>0.146</i>	<i>0.072</i>	<i>0.006</i>	<i>0.490</i>	<i>0.027</i>	<i>0.234</i>	<i>0.344</i>	
Polar labels	3	0.104	0.037	0.036	0.369	0.000	0.132	0.008	0.009	0.035	0.140	
		<i>0.554</i>	<i>0.027</i>	<i>0.669</i>	<i>0.358</i>	<i>0.285</i>	<i>0.037</i>	<i>0.357</i>	<i>0.012</i>	<i>0.163</i>	<i>0.078</i>	
F2F with showcard versus CATI												
2-stage questions	2	0.003	0.502	0.945	0.731	0.798	0.145	0.562	0.492	0.081	0.316	
		<i>0.279</i>	<i>0.214</i>	<i>0.895</i>	<i>0.477</i>	<i>0.594</i>	<i>0.661</i>	<i>0.613</i>	<i>0.400</i>	<i>0.192</i>	<i>0.762</i>	
Polar labels	2	0.187	0.671	0.984	0.247	0.426	0.864	0.395	0.136	0.956	0.381	
		<i>0.125</i>	<i>0.038</i>	<i>0.328</i>	<i>0.074</i>	<i>0.287</i>	<i>0.108</i>	<i>0.214</i>	<i>0.041</i>	<i>0.738</i>	<i>0.009</i>	

All p-values from Monte Carlo permutation with 10,000 replications. **Bold:** p-value for chi-square test statistic for equality of vector of response proportions with pooled sample proportions; *Italic:* p-value for ANOVA F -statistic

TABLE A2
Covariate sample means

<i>Covariate</i>	<i>Mean</i>	<i>Covariate</i>	<i>Mean</i>
Age	49.2	Log equivalised household income (£'000 p.a.)*	2.907
Single/widowed/divorced	0.189	Equivalised household income (£'000 p.a.)*	22.04
No. of dependent children	0.534	Weekly hours of work [†]	37.3
Non-white	0.086	Log Hourly wage (£) [†]	2.25
Retired	0.254	Hourly wage (£) [†]	11.07
		Non-disabling health condition	0.132
		Disabling health condition	0.216

* See Holford and Pudney (2013) for explanation of the method of constructing IP2 income variables; [†] Mean computed from positive sample values. Values are pooled sample means for men and women and waves 1-4.