

Short Communication

An addendum to: a meta-analysis of hypothetical bias in stated preference valuation

GIANLUCA STEFANI^{1,*}, RICCARDO SCARPA², GINEVRA V. LOMBARDI¹

¹ Department of Economics and Management, University of Florence, Florence, Italy

² Economics Department, Waikato Management School, Hamilton, New Zealand

Abstract. A recent study published by Murphy *et al.* (2005) reported results of a meta-analysis of hypothetical bias using 28 valuation studies. The authors found a median ratio of hypothetical to actual values of 1.35 but they did not investigate the ratio of variances of the hypothetical and actual value distributions, which is of great relevance in joint stated and revealed preference analysis. We propose an addendum to Murphy *et al.* (2005) to provide some insights on the distribution of the scale factor across 23 studies for which relevant data is available. We distinguish three types of dispersion parameters reported in the literature. We find that the ratio of real to hypothetical standard deviations of marginal distributions of WTP is about 0.6.

Keywords. Contingent Valuation, experiments, scale identification, meta-analysis, stated preferences

JEL Codes. C9, H41, Q26, Q28

1. Introduction

Stated preference methods are widely used in nonmarket valuation of environmental goods. However, they have been criticised for a number of reasons revolving around the issues of credibility and reliability of hypothetical responses (Cummings *et al.* 1997, Diamond & Hausman 1993, Green *et al.* 1998, to name but a few in the context of contingent valuation). The difference between responses in hypothetical and real payment settings, known as hypothetical bias, is an issue that has given rise to a fierce debate among scholars and has motivated much research effort. A key question in the ensuing research agenda has been the estimation of a calibration factor (CF). This is the ratio between hypothetical (stated preference) and actual (or revealed preference) values. Using CF, values elicited with hypothetical choices may be corrected to obtain value estimates similar to those obtainable from revealed preference studies.

A methodological approach that has generated much attention in this area of applied research has been the use of meta-analysis. At least three literature reviews or meta-analysis studies have investigated the scope and extension of CF. Harrison and Rutstrom (2008)

¹ Corresponding author: gianluca.stefani@unifi.it.

using 35 observations report a CF ranging from 0.75 to 26. List and Gallet (2001) analyse 29 studies with a total of 174 observations of willingness to pay (WTP) and willingness to accept (WTA) estimates. According to their study hypothetical values are about three times larger than real ones, with CF being larger for WTA rather than WTP or when the values are elicited for public rather than private goods. Little and Berrens (2004) expanded the dataset of List and Gallet to include 17 additional observation. The CF from their dataset ranges from 2.93 to 3.34 with a median value of 3.13. Murphy *et al.* (2005), drawing on the study by List and Gallet, proposed a new meta-analysis focussing on WTP estimates and including only observations that employ the same mechanism to elicit hypothetical and real values. The authors selected 28 studies that yield a total of 83 observations for which the distribution of CF is skewed with a mean value of 2.60 and a median value of 1.35. The authors found mixed results about the determinants of CF. Students and group setting seem to widen CF, while discrete choice format and valuation of private goods would have the opposite effect. In a cautionary note, the authors warn that results are sensitive to model specification and that the choice of explanatory variables is affected by the lack of a theory explaining hypothetical bias.

All four studies purport the ratio between hypothetical and actual values as a key factor in criterion validity of stated preference estimates, under the assumption that values elicited from revealed preference data are closer to the truth. Our point of departure is the observation that distributions of ratios of value estimates are not completely characterised by location parameters alone (such as the mean or the median).

Dispersion parameters are also of crucial importance, especially in the context of joint preference estimation from merged revealed and stated preference data (e.g. Hensher, Louviere and Swait, 1999). In this context there are good theoretical reasons for the existence of a difference in error scale from different data sources, which has been corroborated by much empirical evidence (Louviere, 2001)². Similarly, Cameron *et al.* (2002) state that “What would be most valuable for predicting actual demand behavior from stated preference choice data would be some means of using common underlying systematic preference parameters, [...] mapping the dispersion parameter from the particular stated preference method into the likely corresponding dispersion parameter for a revealed preference choice context. This might allow prediction of the distribution of WTP for real market choices.”

For example, the results from one of the first papers addressing the impact of real vs. hypothetical treatments on values elicited with contingent valuation (CV) of public goods (Cummings *et al.*, 1997) were indeed questioned with respect to the assumption of equal variance across treatments two years later (Haab *et al.* 1999).

We define as inverse relative scale factor (IRSF) the ratio of standard deviations of real over hypothetical value distributions:

$$IRSF = \sqrt{\frac{\sigma_r^2}{\sigma_h^2}} = \frac{\sigma_r}{\sigma_h} \quad (1)$$

² According to Louviere (2001) “experimental manipulations, differences in contexts, actions taken by managers, and the like impact not only distributions of response means but also variances of these distributions”.

where σ_r and σ_h are standard deviations with subscripts referring to “real” and “hypothetical” distributions, respectively. We named the ratio inverse relative scale factor since the scale factor is usually defined as $\mu = 1/\sigma$ (Adamowicz, Louviere and Williams, 1994) and the relative scale factor as $SF = \mu_r/\mu_h$ whilst our index is given by $IRSF = \mu_h/\mu_r$.

The aim of this note is to provide a first estimate of the distribution of the IRSF from a subset of 23 studies out of the original 28 considered by Murphy *et al.* (2005), for which relevant data on scale is available. Our focus is on deriving estimates of the IRSF rather than exploring the determinants of hypothetical bias, therefore the note should be considered as an “addendum” rather than a “comment” to the original paper by Murphy.

The remainder of this note is set out as follows. Section 2 illustrates alternative measures of dispersion of WTP. Section 3 deals with data and estimation procedures. Sections 4 and 5 provides a summary of findings and regression results while section 6 concludes.

We provide an assessment of the empirical distribution of the SF across a sample of stated preference studies finding that differences of variances are mild, a result similar to that provided for CF by Murphy *et al.* (2005). We found that CF and SF are correlated and that factors that affect the former also tend to affect the latter.

2. Alternative measure of dispersion of WTP

Depending on the estimation framework adopted, different measures of dispersion are reported in the studies we reviewed. So, we provide a simple model that helps clarifying the differences among alternative measures.

Let us start with a simple linear-in-the-parameters random WTP model:

$$WTP_i = x_i'\beta + \varepsilon_i \tag{2}$$

A first important distinction we make is between the marginal or unconditional variance of WTP and the variance of the error term of the model:

$$VAR(WTP) = E_x[VAR(WTP|x)] + VAR_x[E(WTP|x)] \tag{3}$$

or

$$VAR(WTP) = VAR(\varepsilon) + VAR_x(x\beta) \tag{4}$$

Equation 4 decomposes the unconditional variance of WTP into two terms. The first terms is the variance of the error term and the second term is the variance of the conditional mean of WTP with respect to a vector of covariates x . It is clear that the ratio of unconditional variances will be always different from the ratio of the error term variances unless x is fixed in the hypothetical and real treatment or the ratio of the $VAR_x(x\beta)$ is the same of the ratio of the $VAR(\varepsilon)$:

$$\frac{\sigma_{wtpH}^2}{\sigma_{wtpR}^2} \neq \frac{\sigma_{\varepsilon H}^2}{\sigma_{\varepsilon R}^2} \tag{5}$$

Further measures of dispersion can arise as some researchers calculate fitted WTP for different representative persons. Then, drawing from the asymptotically joint normal distribution of the maximum likelihood parameter estimates, they build up a sampling distribution of fitted WTP estimates following the procedure originally set out by Krinsky and Robb (1986) to estimate confidence intervals for elasticities. Then the distribution reflects the estimation precision for all of the parameters in the model and not only the error precision, and it shows how estimation efficiency affects the range of plausible values for WTP for a representative subject. In the case of the linear model 2, which is estimated using OLS, it is well known that the asymptotic estimator of $\text{VAR}(b|X)$ is given by:

$$\widehat{\text{VAR}}(b|X) = \sigma_u^2 (X'X)^{-1} \quad (6)$$

Therefore the variance of estimated WTP for a representative subject (at the mean values of x , \bar{x}) is:

$$\text{VAR}(b\bar{x}|X) = \bar{x}' \left[\sigma_u^2 (X'X)^{-1} \right] \bar{x} \quad (7)$$

Which again is different from either $\text{VAR}(\varepsilon)$ and its estimator $\sigma_u^2 = \frac{\widehat{\varepsilon}'\widehat{\varepsilon}}{n-k}$ or from (4).

3. Data and estimation

We supplemented the dataset employed by Murphy *et al.* (2005)³ by recording measures of dispersion for the WTP distribution irrespective of the form in which the measures were provided by the authors. We were able to collect data from 23 out of the original 28 studies providing 67 observations⁴. In addition we retrieve 4 more observations from two studies surveyed by Little and Berrens (2004). Overall, our dataset includes 25 studies and 71 observations.

In the augmented dataset available, measures of dispersion can be classified according to both the type of measure of dispersions outlined in the previous section and the format the dispersion is provided with. We classified different formats for dispersion measure into 4 groups as follows.

- 1) *Standard deviation*. Studies based on experimental auctions and open-ended elicitation formats generally provide data on standard deviations of WTP values distributions.
- 2) *Confidence interval*. Most dichotomous choice and some open ended CV methods provide confidence intervals for the estimates of the mean WTP. As the sample size is similar for real and hypothetical treatments, the width of confidence intervals is proportional to the standard deviation of WTP. Therefore we maintain that the ratio of the sizes of confidence intervals is a close proxy for IRSF.
- 3) *Sigma*. Studies that employ dichotomous choice data often use probit or logit models to explain outcome probabilities. In such cases it is possible to recover the stand-

³ Both dataset, bibliography and description of variable have been made available by Murphy on its own webpage at: <http://faculty.cbpp.uaa.alaska.edu/jmurphy/meta/meta.html>

⁴ The five excluded studies are: Blumenschein, *et al.* (2001); Boyce, *et al.* (1989); Duffield and Patterson (1992); Murphy, *et al.* (2002) and Sinden (1988).

ard deviation σ of the underlying distribution from the inverse of the estimate of the parameter of the bid variable⁵ as described in Cameron and James (1987). Actually, in the case of logit models the inverse of the parameter of the bid variable gives $\kappa = \sigma\sqrt{3/\pi}$ that is the dispersion of the error term in the logistic regression. However, since κ is a linear function of σ the constants cancel out and the probit and logit ratios are equal:

$$\frac{\kappa_r}{\kappa_h} = \frac{\sigma_r}{\sigma_h} \tag{8}$$

- 4) *Scale factor.* Finally, a single study (Carlsson and Martinsson, 2001) carried out using multiple choices, reports directly the corresponding scale factor.

From each of the studies employed here a IRSF value is obtained by dividing the measure of dispersion of the real or actual subsample by the one estimated from the hypothetical subsample.

Table 1. Observations classified according to format and type of dispersion.

Format	Type of distribution			Total
	Error	Parameters	Marginal	
Confidence interval	1	5	8	14
Scale_fact	1	0	0	1
Sigma	14	0	0	20
Stand. Dev	0	8	34	36
Total	16	13	42	71

Most of the observations in the dataset are estimates of marginal WTP distributions, formatted either as standard deviations or confidence intervals. Only 16 observations provide an estimate of error distribution followed by the group of observations where the measure provided is a function of model parameters standard errors through Krinsky-Robb type procedures.

4. Results: summary statistics

Overall, the mean value of the IRSF in our sample is 0.67 with a standard deviation of 0.41. However, the IRSF distribution is quite different across the three distribution types confirming their different nature. All means and medians are smaller than 1, consistently with theoretical expectations. Haab *et al.* (1999) state that real experiments control more effectively for sources of variability, therefore the distribution of elicited values is likely to be less dispersed than in hypothetical settings.

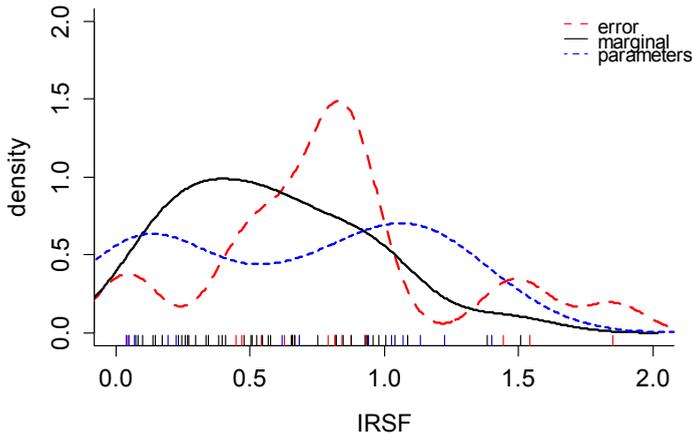
⁵ A single study that directly provides the scale factor for a multinomial logit model also belongs to this group.

Table 2. Summary statistics of IRSF across distribution types.

Form	Type of distribution			Total
	Error	Parameters	Marginal	
min	0.45	0.04	0.07	0.04
max	1.85	1.40	1.51	1.85
mean	0.90	0.67	0.58	0.67
med	0.82	0.69	0.54	0.66
st.dev	0.39	0.50	0.35	0.41

It is worth noticing that with a similar number of observations, measures of IRSF based on all model parameters distributions are more dispersed than those based on model error distribution. Across all types of distribution there are observations with IRSF values higher than 1, a result that mimics what was found by Murphy for the CF. Kernel density estimates of the three distributions of IRSF are plotted in figure 1⁶.

Figure 1. Density estimates of IRSF.



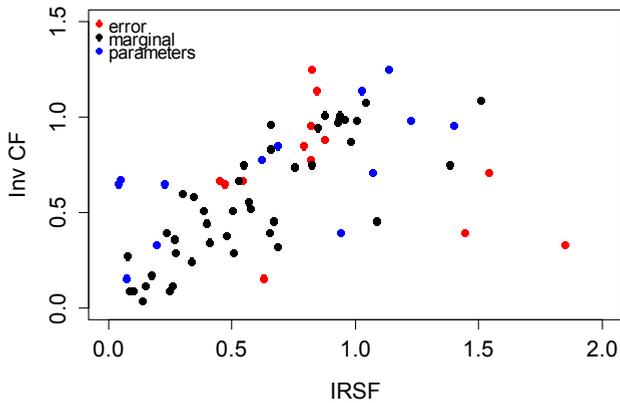
IRSF shows an overall weak and positive correlation ($r=0.58$) with the inverse of CF (ICF), which is the ratio of actual vs hypothetical mean. It is the IRSF obtained from marginal distributions of WTP that shows the highest correlation with ICF. Interestingly, IRSF from error distributions does not seem to be correlated with ICF. However, this results is likely to be affected by the presence of outliers as it can be seen from figure 2.

⁶ A normal kernel density estimate was employed. The bandwidth parameter was selected with Sheater and Jones (1991) formula.

Table 3. Correlation between ICF and IRSF.

	ρ	CI
error	-0.13	(-0.56 , 0.36)
parameters	0.66	(0.18 , 0.89)
marginal	0.81	(0.67 , 0.90)
all	0.58	(0.41 , 0.72)

Figure 2. ICF vs IRSF.



5. Results: regression analysis

To try and explain hypothetical bias we also regress IRSF on ICF and on the explanatory variables used by Murphy *et al.* (2005) (tab. 4). This allows us to see if there is any further marginal and significant effect besides ICF.

Table 4. Regression results: marginal distribution type only.

Estimate	Estimate	Std. Error	t value
(Intercept)	0.15	0.06	2.33
ICF	0.81	0.09	8.61
Choice	-0.08	0.09	-0.82
Private	-0.12	0.09	-1.28
Student	0.02	0.06	0.28
Within	0.19	0.09	1.95
Calibrate	0.19	0.08	2.46

Multiple R-squared: 0.77.
Adjusted R-squared: 0.73.

We include in the regression only observations derived from marginal distributions of WTP due to limitations on the degrees of freedom for the error and parameter groups both with less than 20 observations each. . An R-squared of 0.77 is obtained and the only statistically significant coefficients are those for *ICF*, *within* sample (at the 10% level) and *calibrate*. *Within* might give IRSF closer to 1 either because of carry over effects⁷ or simply because the hypothetical and the real treatment groups are identical. This finding is consistent with evidence of a larger variance in responses in between subject experimental designs (Louviere, 2001). The *calibrate* variable refers to either ex-ante calibration techniques such as budget reminder or cheap talk scripts or ex-post calibration such as using lab experiments to calibrate field data or other uncertainty adjustments (Murphy *et al.*, 2005). The calibration techniques is likely to mitigate the erratic behaviour observed in hypothetical treatments. However, as in the case of the CF, for the SF we also lack a comprehensive theory that explains hypothetical bias. So, the causality of significant parameters should be interpreted with caution.

6. Conclusions

Our point of departure is the observation that most meta-analyses on discrete choice contingent valuation studies comparing real and hypothetical choice settings ignore the role of scale factor. Yet, the issue of estimation efficiency (bias and mean square error) is likely to be as important as the bias question in comparing stated and revealed preferences.

Building on Murphy *et al.* (2005) our study provides some insights on the distribution of the inverse relative scale factor across 25 stated preference studies. The results show that, on average, the IRSF is about 0.6-0.7 and is correlated with the ratio between real and hypothetical average WTPs. However, there are important differences in the distribution of the IRSF depending on which type of WTP distribution is considered: marginal WTP distribution, WTP model error distribution and WTP estimate distribution considered as non linear function of model parameters distribution.

7. Acknowledgments

The authors would like to thank, with the usual disclaimer, James Murphy for the comments provided on various versions of this paper.

6. References

Adamowicz, W., Louviere, J. and Williams, M. (1994) Combining revealed and stated preference methods for valuing environmental amenities. *Journal of Environmental Economics and Management* 26(3): 271-292.

⁷ A carryover effect is an effect that “carries over” from one experimental condition to another. Whenever subjects perform in more than one condition (as in within-subject designs) behaviour in the first condition tends to persist in the subsequent ones.

- Blumenschein, K., Johannesson, M., Yokoyama, K. and Freemand, P. (2001). Hypothetical versus real willingness to pay in the health care sector: results from a field experiment. *Journal of Health Economics* 20: 441-457.
- Boyce, R., Brown, T., McClelland, G., Peterson, G. and Schulze, W. (1989). Experimental Evidence of Existence Values in Payment and Compensation Contexts. In: Proceedings of the Joint Meetings of the Western Committee on Benefits and Costs of Natural Resource Planning (W-133) , Western Regional Science Association, pp. 305-336.
- Cameron, T. (2002). Alternative Non-market Value-Elicitation Methods: Are the Underlying Preferences the Same?. *Journal of Environmental Economics and Management* 44: 391-425.
- Cameron, T. and James, D.(1987). Estimating Willingness to Pay from Survey Data - an Alternative Pre-Test-Market Evaluation Procedure. *Journal of Marketing Research* 24 (4): 389-395.
- Carlsson, F. and Martinsson, P. (2001). Do hypothetical and actual marginal willingness to pay differ in choice experiments? Application to the valuation of the environment. *Journal of Environmental Economics and Management* 41: 179-192.
- Cummings, R., Elliott, S., Harrison, G.W. and Murphy, J. (1997). Are Hypothetical Referenda Incentive Compatible?. *Journal of Political Economy* 105(3): 609-621.
- Diamond, P. and Hausman, J. (1993). On contingent valuation measurement of nonuse values. In: Hausman, J.A., (ed) Contingent valuation: A Critical Assessment. North-Holland, Amsterdam, The Netherlands, pp 3-38.
- Duffield, J. and Patterson, D. (1992). Field Testing Existence Values: Comparison of Hypothetical and Cash Transaction Values. Presented at Joint Western Regional Science Association/W-133. South lake Tahoe, Nevada.
- Green, D., Jacowitz, K., Kahneman, D. and McFadden, D. (1998). Referendum contingent valuation, anchoring, and willingness to pay for public goods. *Resource and Energy Economics* 20(2): 85-116.
- Haab, T., Huang, J. and Whitehead, J. (1999). Are Hypothetical Referenda Incentive Compatible? A Comment. *Journal of Political Economy* 107(1): 186-196.
- Harrison, G. (2006). Experimental Evidence on Alternative Environmental Valuation Methods. *Environmental, Resource Economics* 34: 125-162.
- Hensher, D., Louviere, J. and Swait, J. (1999). Combining sources of preference data. *Journal of Econometrics* 89: 197-221.
- List, J. and Gallett, C. (2001). What Experimental Protocol Influence Disparities Between Actual and Hypothetical Stated Values? Evidence from a Meta-Analysis. *Environmental and Resource Economics* 20: 241-254.
- Little, J. and Berrens R. (2004). Explaining Disparities between Actual and Hypothetical Stated Values:Further Investigation Using Meta-Analysis. *Economics Bulletin*, 3(6): 1-13.
- Louviere, J. (2001). What If Consumer Experiments Impact Variances as well as Means? Response Variability as a Behavioral Phenomenon, *Journal of Consumer Research* 28(3): 506-511.
- Louviere, J., Hensher, D., Swait, J., Adamowicz, W. (2000). Stated Choice Methods: Analysis and Applications. Cambridge University Press, Cambridge (UK)
- Murphy, J., Allen, G., Stevens, T. and Weatherhead, D. (2005). A Meta-Analysis of Hypo-

- thetical Bias in Stated Preference Valuation. *Environmental and Resource Economics* 30: 313-325.
- Murphy, J. and Weatherhead, D. (2002). An Empirical Study of Hypothetical Bias in Voluntary Contribution Contingent Valuation: Does Cheap Talk Matter? Paper prepared for the World Congress of Environmental and Resource Economists Monterey, CA.
- Sheater, S. and Jones, M. (1991). A Reliable Data-based Bandwidth Selection Method for Kernel Density Estimation. *Journal of the Royal Statistical Society* 53: 683-690.
- Sinden, J.A. (1988). Empirical Tests of Hypothetical Bias in Consumers' Surplus Survey. *Australian Journal of Agricultural Economics* 32: 98-112.