



Fondazione Eni Enrico Mattei

**Model Selection and Tests for Non  
Nested Contingent  
Valuation Models:  
An Assessment of Methods**

Margarita Genius\* and Elisabetta Strazzera\*\*

NOTA DI LAVORO 34.2001

**JUNE 2001**

ETA – Economic Theory and Applications

\*Department of Economics, University of Crete

\*\*DRES and CRENoS, University of Cagliari

This paper can be downloaded without charge at:

The Fondazione Eni Enrico Mattei Note di Lavoro Series Index:

[http://www.feem.it/web/attiv/\\_attiv.html](http://www.feem.it/web/attiv/_attiv.html)

Social Science Research Network Electronic Paper Collection:

[http://papers.ssrn.com/paper.taf?abstract\\_id](http://papers.ssrn.com/paper.taf?abstract_id)

Fondazione Eni Enrico Mattei  
Corso Magenta, 63, 20123 Milano, tel. +39/02/52036934 – fax +39/02/52036946  
E-mail: [letter@feem.it](mailto:letter@feem.it)  
C.F. 97080600154

**MODEL SELECTION AND TESTS FOR NON NESTED CONTINGENT  
VALUATION MODELS: AN ASSESSMENT OF METHODS**

**Margarita Genius**

Department of Economics, University of Crete

**Elisabetta Strazzera**

DRES and CRENoS, University of Cagliari

**Address for correspondence:**

Elisabetta Strazzera,

DRES and CRENoS,

University of Cagliari,

Via Fra Ignazio 78,

I-09123, Cagliari, Italy

tel. +39 070 675 3763; fax +39 070 675 3760

**e-mail: STRAZZERA@UNICA.IT**

## 1. Introduction

Survey data for contingent valuation analyses are often obtained through a dichotomous choice questioning framework: individuals are asked if they would be willing to pay some specified amount to insure access to some public good, and the answer may be Yes or No. In *single bound* models the elicitation procedure stops at this stage; while in *multiple bound* models further payment questions follow. Individual responses are then analyzed by means of statistical models to produce an estimate of the value that the public places on the good.

While non parametric or semi-parametric approaches are becoming more popular among contingent valuation practitioners, it is often necessary, for inference or prediction purposes, to uncover a functional relationship between the demand for the public good and individual socioeconomic characteristics. Since the dependent variable is discrete, estimates of the relevant parameters are generally obtained through a maximum likelihood procedure, and the value of the mean, or median, willingness to pay is calculated as a function of the estimated parameters. It is well known that maximum likelihood estimates are consistent if the model specification is correct, but that this does not hold in general for misspecified models: the risk of producing biased estimates of the benefits stemming from the public good is quite serious, and this may diminish the reliability of the analysis for public choice purposes.

Since distributional assumptions are so crucial in the estimation results, it would seem obvious that tests for model specification should play an important part in the statistical analysis of discrete data. In contingent valuation studies it can be observed that the application of tests for nested models is quite common, for example by means of likelihood ratio tests; the analysis, though, is much less accurate when the competing hypotheses are non nested.

The analysis of non nested models has followed two distinct approaches in literature: model selection criteria, and hypothesis testing (cfr. Gourieroux and Monfort (1995)). In the model selection approach, each competing model is evaluated by means of a numerical criterion: for a given sample observation, the procedure consists of selecting the model that optimizes the chosen criterion. A typical example in linear regression is the

(adjusted)  $R^2$  criterion, while in maximum likelihood estimation a commonly used criterion is the information criterion proposed by Akaike (1973), or one of its variants.

The problem of the model selection approach is that it produces a deterministic outcome, defined by the ranking of the values of the criterion, and it does not take into account the probabilistic nature of that result. Vuong (1989) points out that differences in the criterion values may not be statistically significant: yet the deterministic model selection approach would consider a model superior to another one, while in fact they may be considered as statistically equivalent. He then sets the information criterion in a testing framework, where the null hypothesis is that the two competing models are equally close to the true model.

The hypothesis testing approach takes a step further, extending the classical testing procedures to the case of non nested hypotheses: examples are the generalized Wald test, the generalized score test, and the Cox test, which is a generalized likelihood ratio test; or, in a different line, the tests based on artificial nesting: the Davidson-MacKinnon (1981) test, the Atkinson test and the Quandt test belong to the latter category (cfr. Gouriéroux and Monfort, cit.).

In contingent valuation analyses, non nested competing models are generally assessed by means of selection criteria, such as Akaike's (1973), while we are not aware of any testing approach in this field; and it might be added that such applications are very few in discrete data modeling in general. As put forth by Pesaran and Weeks (2000), the extra computational difficulties that the testing approach entails may explain why this path has been so neglected. However, there may be also a more theoretically founded justification for the choice of model selection criteria over hypothesis testing to test economic theories: as pointed out by Granger et al. (1995), the choice of the null hypothesis and the significance level is arbitrary, and this is even more so when testing is applied to non nested hypothesis. In their view, when the choice of the particular model is data dependent it is "better to use well-thought-out" model selection procedures rather than formal hypothesis testing.

The aim of this paper is to investigate on the performance of either approach in selecting among different contingent valuation models applied to simulated data. In

particular, we compare three methods that are based on the Kullback-Leibler Information Criterion (KLIC):

- the Akaike information criterion;
- the Vuong test;
- the Cox test, in the simulated approach of Pesaran and Pesaran (1993).

The structure of the paper is the following: section 2 gives a brief background about the KLIC and explains the 3 procedures above, section 3 describes the experimental setting of the simulation exercise, section 4 reports the results of the experiments and finally section 5 contains our conclusions.

## 2. Methods

In order to describe the different statistics or criteria we introduce some notation and terminology.

Consider a sequence  $(Y_i, X_i)$   $i=1,2,\dots$  of i.i.d. random vectors. The modeler is interested in the conditional probability distribution of the vector  $Y_i$  given  $X_i$ . Define the true conditional density as:

$$\ell_0(y|x) = \prod_{i=1}^n \varphi_0(y_i | x_i),$$

which is unknown. To evaluate its proximity to a specified parametric model, that we define as:

$$\ell(y|x;\theta) = \prod_{i=1}^n \varphi(y_i | x_i; \theta), \theta \in \Theta,$$

we make use of the notion of Kullback-Leibler Information Criterion (KLIC):

$$K_n(\ell(y|x;\theta)/\ell_0(y|x)) = \frac{1}{n} E_0 \left( \log \frac{\ell_0(Y|x)}{\ell(Y|x;\theta)} \right).$$

We will be interested in comparing pairs of competing parametric families of conditional densities of  $Y_i$  given  $X_i$  given by

$$H_f : \{f(y_i | x_i; \beta), \beta \in B \subset \mathfrak{R}^F \},$$

$$H_g : \{g(y_i | x_i; \gamma), \gamma \in \Gamma \subset \mathfrak{R}^G \},$$

where the models  $H_f$  and  $H_g$  are strictly non-nested.

It can be shown (cfr. Gourieroux and Monfort, cit.) that the asymptotic Kullback-Leibler proximity between the true probability distribution and a given parametric model is approximated by

$$\tilde{K} = \frac{1}{n} \sum_{i=1}^n \log \varphi_0(y_i | x_i) - \frac{1}{n} \sum_{i=1}^n \log \varphi(y_i | x_i; \hat{\theta}_n),$$

where  $\hat{\theta}_n$  is the maximum likelihood estimator of  $\theta$ .

Since  $\varphi_0$  is unknown  $\tilde{K}$  cannot be used; it can be noticed, though, that when two models are compared, the first term of  $\tilde{K}$  remains constant, so that minimization of the criterion only depends on the second term, i.e. on the maximum likelihood of the two competing models.

Denoting by  $\hat{\beta}_n$  and  $\hat{\gamma}_n$  the (quasi) maximum likelihood estimators of  $\beta$  and  $\gamma$  under  $H_f$  and  $H_g$  respectively, this amounts to calculating:

$$LR_n(\hat{\beta}_n, \hat{\gamma}_n) = \sum_{i=1}^n \log f(y_i | x_i; \hat{\beta}_n) - \sum_{i=1}^n \log g(y_i | x_i; \hat{\gamma}_n),$$

i.e. the likelihood ratio of the two models.

The drawback of using  $LR_n$  as such, is that it increases for more general models. In order to overcome this problem, Akaike (1973) proposed a correction of this criterion, that penalizes more complex models. The Akaike Information Criterion (AIC) penalizes the log-likelihood of each model by a quantity equal to the number of its parameters:

$$AIC = \left( \sum_{i=1}^n \log f(y_i | x_i; \hat{\beta}_n) - p \right),$$

where  $p$  is the number of parameters. The Akaike criterion for model selection (AICMS) simply consists in comparing the AIC values for the two models:

$$AICMS = \left( \sum_{i=1}^n \log f(y_i | x_i; \hat{\beta}_n) - p \right) - \left( \sum_{i=1}^n \log g(y_i | x_i; \hat{\gamma}_n) - q \right).$$

If the value is positive the first model is chosen, otherwise the second will be deemed best. Obviously, if the two models are characterized by the same number of parameters  $p$  and  $q$ , the Akaike criterion for model selection reduces to  $LR_n$ .

A criticism to the use of model selection criteria such as Akaike's is that they are deterministic: the model that satisfies the given criterion is selected. However, some authors point out that this result is just the outcome of a random draw from the sample space, and as such should be treated in probabilistic terms.

This issue is addressed by Vuong (1989), whose approach sets the model selection criterion in a hypothesis testing framework. More specifically, it tests whether the models under consideration are equally close to the true model, where closeness is measured by the KLIC.

The null hypothesis is given by:

$$H_0 : E_0 \left[ \log \frac{f(y_i/x_i; \beta^*)}{g(y_i/x_i; \gamma^*)} \right] = 0, \text{ (both models are equivalent)}$$

against

$$E_0 \left[ \log \frac{f(y_i/x_i; \beta^*)}{g(y_i/x_i; \gamma^*)} \right] > 0, \text{ (} H_f \text{ is better than } H_g \text{), or}$$

$$E_0 \left[ \log \frac{f(y_i/x_i; \beta^*)}{g(y_i/x_i; \gamma^*)} \right] < 0, \text{ (} H_g \text{ is better than } H_f \text{),}$$

where  $\beta^*$  and  $\gamma^*$  are the pseudo-true values of  $\beta$  and  $\gamma$  respectively. The tests statistics proposed by Vuong are the following:

-an unadjusted LR statistic given by

$$n^{-1/2} LR_n(\hat{\beta}_n, \hat{\gamma}_n) / \hat{\omega}_n,$$

where

$$\hat{\omega}_n = \sqrt{\frac{1}{n} \sum_{i=1}^n \left[ \log \frac{f(y_i/x_i; \hat{\beta})}{g(y_i/x_i; \hat{\gamma})} \right]^2 - \left[ \frac{1}{n} \sum_{i=1}^n \log \frac{f(y_i/x_i; \hat{\beta})}{g(y_i/x_i; \hat{\gamma})} \right]^2};$$

-an adjusted LR statistic given by

$n^{-1/2} L \tilde{R}_n(\hat{\beta}_n, \hat{\gamma}_n) / \hat{\omega}_n$  where

$L \tilde{R}_n(\hat{\beta}_n, \hat{\gamma}_n) \equiv LR_n(\hat{\beta}_n, \hat{\gamma}_n) - \xi_n$ , and  $\xi_n$  is a correction factor that penalizes each model for model complexity. Different correction factors, as well as a slightly different version of the denominator term, give rise to different variants of the Vuong's statistics, that in any case, for non nested models, is asymptotically standard normal under  $H_0$ .

While Vuong's approach is to test if the two models are statistically different, the Cox approach aims at testing if the true conditional probability distribution belongs to one of the competing models under examination. The null hypothesis may be that the true Data Generating Process (DGP) belongs to  $H_f$ ; but it also may be that the DGP belongs to  $H_g$ . Due to the special role of the null hypothesis in this context, it is not obvious which null hypothesis we should choose. Many (see Pesaran and Pesaran (1993), henceforth P&P; Weeks (2000)) advocate performing the non-nested test twice by reversing the role of the null and alternative hypothesis. This procedure could very well lead to a situation where both models are accepted or both are rejected.

Following P&P, the standardized Cox statistic is asymptotically normal and under the null  $H_f$  is given by

$$S_f(\hat{\beta}_n, \hat{\gamma}_n) = \sqrt{n} T_f / \hat{v}_f,$$

where  $\hat{v}_f$  is an estimate of the asymptotic variance, and

$$T_f = \frac{1}{n} LR_n(\hat{\beta}_n, \hat{\gamma}_n) - \hat{E}_f \left( \frac{1}{n} LR_n(\hat{\beta}_n, \hat{\gamma}_n) \right)$$

The expression  $\hat{E}_f(\cdot)$  stands for the conditional expectations operator under  $H_f$ ; it should be noted that  $E_f \left( \frac{1}{n} LR_n(\hat{\beta}_n, \hat{\gamma}_n) \right)$  is zero when we have nested models but does not vanish in the case of non-nested models. Due to the difficulties in computing this term (see Pesaran and Weeks (2000)), this test has not been widely applied outside the linear regression model. The difficulty lies in computing an estimate of  $E_f \left( \frac{1}{n} LR_n(\hat{\beta}_n, \hat{\gamma}_n) \right)$ ,



because it entails finding an estimate of the pseudo true value  $\gamma^*$ , i.e the value that maximizes  $E_f(\log g(y|x, \gamma))$ .

In the case of discrete choice models, P&P have derived a simulation method to compute the above statistic which we can apply to the case of the single bound CV model. P&P simulate  $R$  independent samples of  $n$  indicators (dependent variable) assuming that  $F$  is the true distribution; then for each one of the  $R$  simulated samples they compute the maximum likelihood estimate of  $\gamma$  using the c.d.f.  $G$ . Denoting by  $\hat{\gamma}_n^*(R)$  the average of the  $R$  estimates of  $\gamma$ , this is a consistent estimate of the pseudo true value. Finally we can estimate the expected value above as follows:

$$\hat{E}_f\left(\frac{1}{n}LR_n(\hat{\beta}_n, \hat{\gamma}_n)\right) = \frac{1}{n} \sum_{i=1}^n \left\{ (1 - F_i(\hat{\beta}_n)) \log\left(\frac{1 - F_i(\hat{\beta}_n)}{1 - G_i(\hat{\gamma}_n^*(R))}\right) + F_i(\hat{\beta}_n) \log\left(\frac{F_i(\hat{\beta}_n)}{G_i(\hat{\gamma}_n^*(R))}\right) \right\}.$$

Since there is no a priori reason why  $F$  should be the null hypothesis, P&P suggest to reverse the null and alternative hypothesis, i.e. testing  $G$  against  $F$ : therefore it will be necessary to find an estimate of the expression  $E_g\left(\frac{1}{n}LR_n(\hat{\beta}_n, \hat{\gamma}_n)\right)$  and this in turn will require finding an estimate of the value that maximizes the expected value of the log-likelihood using model  $F$  when  $G$  is the null model. Full details on the derivation of the Cox simulation based test statistic are given in Appendix 2.

### 3. Experimental Design

The dichotomous choice elicitation method for contingent valuation produces a dichotomous type of response to payment questions that are differentiated among individuals. This particular setting allows different modeling options: the latent dependent variable can be modeled either as a dichotomous variable, as in the random utility model (RUM) framework used in the utility differential model by Hanemann (1984); or as a censored variable, which is the approach proposed by Cameron and James (1987) and Cameron (1988). The latter produces separate estimates for the coefficients and the scale

parameter of the model, and their standard errors, and allows for a more straightforward calculation for the mean or median value of the public good, and was therefore chosen for this application.

Depending on the assumptions on the individuals' preferences, the latent variable, individual willingness to pay, or *wtp* for brevity, can be modeled as a linear or non linear function of the individual socioeconomic covariates. The econometric modeling involves further assumptions on the distribution of the error term, and its functional relationship with the deterministic part of the *wtp* model: the combination of the two components can possibly give rise to many modeling specifications, but in practice probit, logit, log-normal, log-logistic, weibull are the most commonly used. This choice may be due to the fact that they can easily be estimated with econometric modules available in popular statistical packages like Limdep, Stata, or Sas (cfr. Hanemann and Kanninen, 1999).

For our experiment, we considered different DGP for the *wtp*, obtained from a linear functional form for the deterministic part of the model, which is common to all experiments, and an additive error term that is varied across experiments. The general model is the following:

$$wtp_i = 27 + 1.5x_{i1} - 3x_{i2} + 0.5x_{i3} + \varepsilon_i,$$

where  $x_1$  and  $x_3$  are continuous variables respectively ranging from 4 to 75 and from 0.5 to 1.5; while  $x_2$  is a qualitative variable, taking values zero or one.

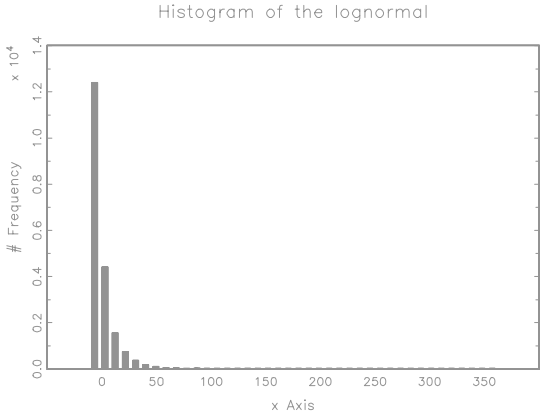
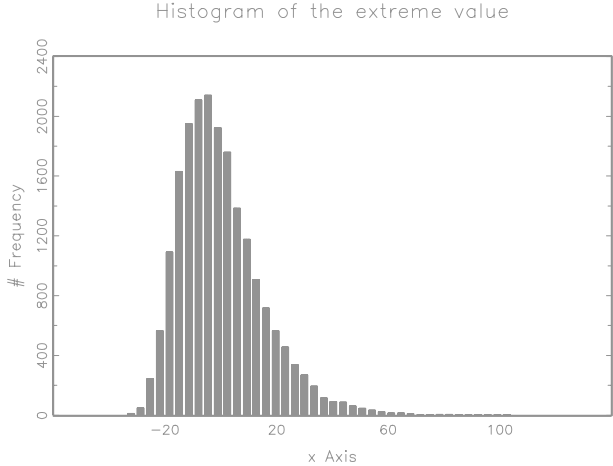
For the error term we investigate three different scenarios, allowing for differences in the degree of skewness of the distributions. In the first one the error term is distributed as a Normal with zero mean and standard deviation 15. In the second one we assume that the error term follows an extreme value distribution<sup>1</sup> with mean zero and standard deviation 15; while in the third scenario the error term is generated as a translated lognormal<sup>2</sup>, obtained from a lognormal by subtracting its mean, so that the resulting error has mean

---

<sup>1</sup> The parameters of the extreme value distribution are:  $a=-b\gamma$  and  $b=15\frac{\sqrt{6}}{\pi}$  where  $\gamma$  is Euler's constant which we approximate with the value 0.5772156649.

<sup>2</sup> The parameters of the lognormal are:  $\mu = \ln\left(\frac{15}{\sqrt{e \cdot (e-1)}}\right)$ ,  $\sigma=1$ .

zero and standard deviation 15. The shape of the last two distributions is depicted in the histograms below which have been produced by generating samples of size 20000 from the two distributions described above.



We hypothesize that the researcher assumes a model linear in the covariates, with an additive error term, obtaining the following econometric model: for each individual  $i$ ,

$$Y_i = x_i' \delta + \varepsilon_i,$$

where  $x$  is the vector of regressors. In this model the latent variable  $Y_i$  is unobserved: the observed variable is the answer YES or NO to the question regarding whether or not the individual would be willing to pay a given amount  $t_i$ .

For a given sample of  $n$  independent observation, the generic log-likelihood function is:

$$\begin{aligned}
LogL &= \sum_{i=1}^n \{I_i \log [1 - \Psi_i ((t_i - x_i' \delta) / \nu)] + (1 - I_i) \log [\Psi_i ((t_i - x_i' \delta) / \nu)]\} \\
&= \sum_{i=1}^n \{I_i \log [1 - \Psi_i (\theta)] + (1 - I_i) \log [\Psi_i (\theta)]\}
\end{aligned}$$

where  $\Psi$  represents generically one of the distributions hypothesized by the researcher;  $\theta = (\delta, \nu)$ , and  $I_i$  is a dummy variable assuming value one if the individual response to the bid question is positive, zero otherwise. Since bids are varied among individuals,  $\delta$  and  $\nu$  can be estimated separately.

A further assumption is that the researcher (righteously) thinks that the deterministic part of the model is correctly specified, but is unsure about the distribution of the random term, and tries different specifications: Normal, Logistic, Extreme Value, that combined with the linear function for the deterministic part of the model give rise to the Probit, Logit, and Weibit<sup>3</sup> models. As mentioned earlier, the first two models are frequently applied by contingent valuation practitioners: the underlying distributions for these two models are both symmetric, with fatter tails for the logistic. The weibit model is much less common; we choose it because it is an example of asymmetric distribution associated to a linear functional form for the deterministic part of the model, and this facilitates comparisons between models for our purposes.

Since the model checking methods under analysis are based on the KLIC, the models have to be compared in pairs. The Akaike criterion only requires maximization of the log-likelihood for each model, and then the calculation is straightforward; in Vuong's approach the calculation is slightly more involved, but still it only requires the computation of the maximum log-likelihood of each model. For the simulated Cox test the procedure is definitely more complex, since it involves the computation of the quasi maximum likelihood estimates of each model assuming that the other model is the true DGP.

The purpose of this study is to assess the three approaches in situations that are typically encountered by applied researchers, rather than to investigate the large sample

---

<sup>3</sup> Notice that this is a different specification from the non linear Weibull model frequently used in contingent valuation studies.

properties of the tests. The sample sizes considered in the experiments, 300, 600, and 1000 observations, are representative of a small, medium and large contingent valuation data set.

#### 4. Results

We now examine the results of our simulation for each of the three DGP (Normal, Extreme Value, translated Lognormal) and the three candidate models: probit, logit and weibit. As explained in the preceding section, the KLIC requires that the three models should be compared in pairs, so the tables report results for the pairs probit vs logit, probit vs weibit, and logit vs weibit.

In table 1 of each experiment we report the estimates of the parameters obtained from the two models to be compared. The number of replications in the Monte Carlo experiments was fixed at 300, but the actual number in each experiment depends on the rate of convergence –always very high, even though it never attains 100%; the differences in the number of successful replications explain the slight differences between estimates obtained under the same specification and the same data but in a different experiment.

It can be observed that parameter estimates do not differ very much across models. The parameter  $\nu$  is a scale parameter for the logit and the weibit, that should be multiplied respectively by  $\pi/\sqrt{3}$  and  $\pi/\sqrt{6}$  to obtain the estimated standard deviations for the two models. Moreover, since the mean of the extreme value distribution is not zero we have to add  $0.5772 \cdot \nu$  to the estimated constant by the weibit to compare it to the corresponding probit estimate. A general feature of all experiments is that the estimates of the parameters  $\delta_3$  and especially  $\delta_4$  are less efficient than the others; we are not able at this stage to provide an explanation of this discrepancy in the precision of the parameter estimates.

Similar parameter estimates produce similar estimates for the mean wtp value, as it can be observed in table 2 of each experiment; but it can be observed that the asymmetry of the extreme value distribution generally produces a lower value for the estimated median – whether it is more correct or biased depends on the true DGP. Choice of the median rather than the mean as a central tendency measure is often deemed to be preferable, both on statistics and economics grounds (cfr. Hanemann and Kanninen, cit., pp. 325-6): therefore, we will focus on the estimated value for this central tendency measure for the weibit

model. Closeness of the estimated mean and median wtp values to the true population value can be measured by the mean square error (MSE); this criterion can give us a measure of the goodness of fit of each specification. Since the sampling happened to be not well behaved, with the 600 sample being closer to the population values than the 1000 sample, in the same table we report also the MSE with respect to the sample WTP values.

Finally, table 3 of each experiment shows the conclusions drawn from each of the three methods under investigation: the Akaike information criterion, the Vuong test and the Cox test.

From table 3 of each experiment, we can derive the probability of rejecting the null for the Vuong and Cox tests when the nominal size is five percent. For the Vuong test, the null hypothesis is that the two models are *equivalent* and therefore the probability of rejecting the null should approach 1 if the null is false, while it should approach 0.05 if the null is true. The null hypothesis for each one of the two Cox tests is that one of the models is correctly specified, therefore the probability of rejecting each one of the two nulls should be computed considering the two cases when the corresponding null is rejected, either because the alternative is accepted, or because both are rejected.

To help the exposition, we use the following code for our experiments:

| <b>Pairs of models</b><br><b>DGP</b> | <b>Probit</b><br>vs<br><b>Logit</b> | <b>Probit</b><br>vs<br><b>Weibit</b> | <b>Logit</b><br>vs<br><b>Weibit</b> |
|--------------------------------------|-------------------------------------|--------------------------------------|-------------------------------------|
| <b>Normal</b>                        | <b>A-I</b>                          | <b>A-II</b>                          | <b>A-III</b>                        |
| <b>Extreme value</b>                 | <b>B-I</b>                          | <b>B-II</b>                          | <b>B-III</b>                        |
| <b>Trans. Lognormal</b>              | <b>C-I</b>                          | <b>C-II</b>                          | --                                  |

We analyze the results by first considering the experiments where one of the models was correctly specified: i.e. experiments A-I and A-II, and experiments B-II and B-III.

As recalled above, in these cases the probability for the Vuong test to reject the null should approach one. In experiment A (probit correctly specified) it falls quite short of one even for  $n=1000$ : the probability is 23% in I (probit vs logit), and 52% in II (probit vs weibit). In experiment B (weibit correctly specified) the probability is 42% in II and 57% in III (logit vs weibit). In terms of power of the test, these results are not very encouraging for the Vuong test; we are not aware of any simulation study where the behavior of the Vuong test has been investigated, so we cannot compare our results with other benchmarks.

Let's see now, for the same set of experiments, the performance of the Cox test. We can observe mixed results, but in general it performs better, in terms of power of the test, than the Vuong test, especially when the probit and the weibit are compared. In addition, the power improves, while the size of the tests decreases, with the sample size; the latter reaching levels of 10-13% for  $n=1000$ . In this sense our results are comparable to Weeks (2000)<sup>4</sup>. It remains to analyze the performance of the Akaike criterion: when one of the models is correctly specified, and the sample size is large enough (in our case, 600 or more) the Akaike criterion works very well when the weibit is compared with the other two; while when two similar models, probit and logit, are compared, the Akaike criterion is not able to pick up the correct model in more than 20% of the cases. For the small sample size this problem extends also to the comparisons involving the weibit model: with the exception of experiment B-III, the Akaike criterion may lead to choosing the wrong model in about the 20% of the cases.

Let's now consider the subset of experiments where both models are misspecified. These are experiments A-III, B-I, and both experiments C.

In experiment A-III, we compare the logit and the weibit, the DGP being Normal. The Akaike criterion will choose the logit over the weibit, with increasing frequency as the sample size increases. The Vuong test seems to point out that both models are equivalent in terms of KLIC; but as the sample size increases, the probability that the logit model is chosen increases as well. It can be noticed that the probability that the Vuong test selects

---

<sup>4</sup> It should be noted that Weeks compares alternative variants for the Cox statistic and does not use the simulated version of P&P for all sample sizes.

the “worse” model (in this case, the weibit) is close to zero. As for the Cox test, its behavior is not very satisfying for the small sample; while it improves with sample size, yet in this case the power of the first Cox test falls short of 53% even for  $n=1000$ .

Experiment B-I considers probit versus logit when the true DGP is Extreme Value. The Akaike criterion would choose logit or probit with almost the same probability when the sample size is small, but the frequency of the logit choice increases with sample size. Here the probability of rejecting the null for the Vuong test reaches 0.096 and both models are found to be equivalent 90% of the time, while the Cox test accepts both models over 65% of the time for  $n=1000$ . Here the Cox test lacks power since the probability of rejecting the null is below 23% for both tests for  $n=1000$ .

This situation is reversed when we consider the last experiments, under the translated Lognormal DGP. Here the Akaike criterion discriminates much more between models, leading to choosing the logit and the weibit over the probit. Also the Vuong test now shows a strong support for the same models and it chooses the logit around 85%, and the weibit 100% of the time for  $n=1000$ . While the selected models, especially the weibit, are indeed closer to the true DGP, none of them is correct: in this case the Vuong test, like the Akaike criterion, is not able to signal that both models are misspecified. Only the Cox test points out that in this situation none of the models is to be accepted, again improving with sample size, with the power of both Cox tests approaching one for the large samples.

Summing up, when both models are misspecified, the Akaike criterion is obviously dominated by the other two methods, since it does not give any information about the misspecification problem. For the other two methods, we get mixed results, depending on the DGP and the sample size. In general, the Cox test seems to perform better than the Vuong, especially for large samples; for the smaller sample size instead the Cox test is not so satisfying, since it often accepts the null when it is false.

## **5. Conclusions**

From an operative point of view, it is important that the selection method be able to signal a possible misspecification when its consequences are more serious. In our context this could be answered by considering the MSE of the estimated mean and median. The



most serious consequences of misspecification can be observed in experiment C, where the data generated by a strongly skewed distribution (the translated Lognormal) is fitted with symmetric (probit and logit), or slightly asymmetric (weibit) models. In this experiment we could see that the Cox test is the only one that, for an adequately large sample size, uncovers the misspecification problem.

Unfortunately, in other cases we could observe that the power of the Cox test is unacceptably low: this holds for example in experiment B-I, where misspecification leads to a relevant bias in the estimated median wtp. Also, the performance of the Cox test when the sample size is small is generally poor.

In conclusion, the results of our experiments can be interpreted as follows:

- the Akaike model selection criterion is unsatisfactory, since when both models are misspecified it cannot signal it; even when one model is correctly specified, for small size samples it presents a non negligible probability of choosing a wrong model. When the sample size is small, it may be preferable to use the Vuong test rather than the Akaike criterion, since it makes less mistakes in the choice of the model, and may suggest further checks of the specification in the (too frequent) case that the two models are deemed equivalent;
- when the sample size is large enough, it is preferable to use the Cox test, that, although more complex to implement, seems more effective and conclusive than the two model selection methods.

In this paper we have analyzed three methods for model selection based on the KLIC: as pointed out by Aznar Grasa (1989), this is just one of many alternative measures of closeness. Further research is required to find out if other methods are superior to those assessed in this study.

## REFERENCES:

- Akaike, H. (1973): "Information theory and an extension of the maximum likelihood principle," in *Proceedings of the 2<sup>nd</sup> International Symposium on Information Theory*, ed. by N. Petrov, and F. Csadki, pp. 267-281. Akademiai Kiado, Budapest.
- Aznar Grasa, A. (1989): *Econometric Model Selection: A New Approach*. Kluwer Academic Publishers, Spain.
- Cameron, T. A. (1988): "A new paradigm for valuing non-market goods using referendum data: maximum likelihood estimation by censored logistic regression," *Journal of Environmental Economics and Management*, **15**, 355-79.
- Cameron, T. A., and M.D. James (1987): "Efficient estimation methods for closed-ended contingent valuation surveys," *The Review of Economics and Statistics*, **69**, 269-76.
- Davidson, R., and J.G. MacKinnon (1981): "Several tests for model specification in the presence of alternative hypotheses," *Econometrica*, **49**, 781-793.
- Gourieroux, C., and A. Monfort (1995): *Statistics and Econometric Models*, Cambridge University Press.
- Granger, C., M. L. King, and H. White (1995): "Comments on Testing Economic Theories and the Use of Model Selection Criteria," *Journal of Econometrics*, **67**, 173-187.
- Hanemann W. M., and B. J. Kanninen (1999): "Statistical analysis of discrete response cv data", in Bateman, I. and K. Willis, *Valuing Environmental Preferences*, Oxford University Press.
- Pesaran, M. H., and B. Pesaran (1993): "A simulation approach to the problem of computing Cox's statistic for testing nonnested models," *Journal of Econometrics*, **57**, 377-392.
- Pesaran, M. H., and M. Weeks (2000): "Non-nested hypothesis testing: an overview", in *Companion to Theoretical Econometrics*, ed. Badi H. Baltagi, Basil-Blackwell, Oxford. (forthcoming).
- Vuong, Q.H. (1989): "Likelihood Ratio Tests for Model Selection and Non-Nested Hypothesis," *Econometrica*, **57(2)**, 307-333.
- Weeks, M. (2000): "Testing the Binomial and Multinomial Choice Models Using Cox's Non-Nested Test," in Mariano, R., Schuermann, T. and Weeks, M.J., *Simulation based inference in econometrics: methods and applications*, Cambridge University Press, Cambridge.

## APPENDIX 1

### A. NORMAL DGP

#### **Experiment A-I: Normal DGP, Probit vs Logit**

Table A-I.1 *Parameter estimates<sup>a</sup> for normal DGP using  $H_f$ (normal) and  $H_g$  (logistic) across 300 replications<sup>b</sup>.*

| Parameters | Sample Size       |                   |                   |                   |                   |                   |
|------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
|            | 300               |                   | 600               |                   | 1000              |                   |
|            | $H_f$             | $H_g$             | $H_f$             | $H_g$             | $H_f$             | $H_g$             |
| $\delta_1$ | 26.488<br>(9.056) | 26.745<br>(9.154) | 26.450<br>(5.747) | 26.688<br>(5.661) | 26.609<br>(4.647) | 26.853<br>(4.674) |
| $\delta_2$ | 1.508<br>(0.118)  | 1.502<br>(0.118)  | 1.501<br>(0.083)  | 1.497<br>(0.084)  | 1.503<br>(0.067)  | 1.497<br>(0.068)  |
| $\delta_3$ | -2.817<br>(4.264) | -2.803<br>(4.216) | -2.847<br>(2.907) | -2.872<br>(2.935) | -3.009<br>(2.120) | -2.965<br>(2.126) |
| $\delta_4$ | 0.531<br>(6.260)  | 0.497<br>(6.411)  | 0.796<br>(4.439)  | 0.783<br>(4.390)  | 0.762<br>(3.407)  | 0.722<br>(3.450)  |
| $v^c$      | 14.752<br>(2.018) | 8.290<br>(1.215)  | 14.829<br>(1.387) | 8.346<br>(0.894)  | 14.863<br>(1.021) | 8.325<br>(0.643)  |

- a) Mean values and standard deviations (in parenthesis) over 300 replications.  
b) The actual number of successful experiments was 284, 288 and 293 for the 300, 600 and 1000 sample size respectively.  
c) The estimated scale parameter of the logit should be multiplied by  $\pi/3^{1/2}$  for comparison with the corresponding probit estimate.

Table A-I.2 *Mean estimated wtp and their MSE*

|      | $H_f$                   |  | $H_g$                   |                    |
|------|-------------------------|--|-------------------------|--------------------|
|      | Mean-Median<br>(st.dev) | MSE <sup>a</sup><br><i>SMSE</i> <sup>b</sup> | Mean-Median<br>(st.dev) | MSE<br><i>SMSE</i> |
| 300  | 86.314<br>(1.869)       | 6.773<br>3.482                               | 86.301<br>(1.886)       | 6.788<br>3.545     |
| 600  | 84.351<br>(1.158)       | 1.359<br>1.343                               | 84.374<br>(1.167)       | 1.374<br>1.361     |
| 1000 | 85.455<br>(0.991)       | 1.892<br>0.979                               | 85.458<br>(0.990)       | 1.895<br>0.977     |

- a) MSE with respect to the population true mean wtp: 84.500  
b) MSE with respect to the sample true mean wtp: 86.292, 84.429, and 85.457 for the 300, 600 and 1000 sample size respectively.

Table A-I.3 *Conclusions drawn from each method (% frequency)*

| Method |      | Conclusion   |  |  |  |
|--------|------|--|--|--|--|
| Akaike |      | H <sub>f</sub> is better                           |  | H <sub>g</sub> is better                           |  |
|        | 300  | 0.711  |  | 0.288  |  |
|        | 600  | 0.781  |  | 0.218  |  |
|        | 1000 | 0.789  |  | 0.211  |  |
| Vuong  |      | H <sub>f</sub> is better                           | H <sub>g</sub> is better                           | H <sub>f</sub> and H <sub>g</sub> equivalent       |  |
|        | 300  | 0.204  | 0.003  | 0.792  |  |
|        | 600  | 0.260  | 0.000  | 0.740  |  |
|        | 1000 | 0.228  | 0.003  | 0.769  |  |
| Cox    |      | H <sub>f</sub> accepted<br>H <sub>g</sub> rejected | H <sub>g</sub> accepted<br>H <sub>f</sub> rejected | Both H <sub>f</sub> and<br>H <sub>g</sub> accepted | Both H <sub>f</sub> and<br>H <sub>g</sub> rejected |
|        | 300  | 0.373  | 0.024  | 0.429  | 0.172  |
|        | 600  | 0.423  | 0.013  | 0.423  | 0.138  |
|        | 1000 | 0.413  | 0.003  | 0.467  | 0.116  |

## Experiment A-II: Normal DGP, Probit vs Weibit

Table A-II.1 *Parameter estimates<sup>a</sup> for normal DGP using  $H_f$  (normal) and  $H_g$  (extreme value) across 300 replications<sup>b</sup>.*

| Parameters | Sample Size       |                   |                   |                   |                   |                   |
|------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
|            | 300               |                   | 600               |                   | 1000              |                   |
|            | $H_f$             | $H_g$             | $H_f$             | $H_g$             | $H_f$             | $H_g$             |
| $\delta_1$ | 26.688<br>(8.438) | 19.308<br>(9.631) | 26.875<br>(5.872) | 19.360<br>(6.644) | 27.231<br>(4.497) | 19.649<br>(5.019) |
| $\delta_2$ | 1.499<br>(0.106)  | 1.507<br>(0.118)  | 1.495<br>(0.084)  | 1.498<br>(0.089)  | 1.497<br>(0.063)  | 1.503<br>(0.065)  |
| $\delta_3$ | -3.357<br>(4.291) | -3.512<br>(4.756) | -2.903<br>(2.877) | -3.016<br>(3.212) | -2.888<br>(2.329) | -3.045<br>(2.541) |
| $\delta_4$ | 1.033<br>(6.268)  | 1.014<br>(6.736)  | 0.747<br>(4.629)  | 0.774<br>(4.920)  | 0.236<br>(3.269)  | 0.076<br>(3.633)  |
| $v^c$      | 14.573<br>(1.940) | 12.880<br>(2.053) | 14.676<br>(1.426) | 13.093<br>(1.584) | 14.761<br>(1.151) | 13.281<br>(1.149) |

- a) Mean values and standard deviations (in parenthesis) over 300 replications.  
b) The actual number of successful experiments was 292 for the 300, and 296 for the 600 and 1000 sample size.  
c) The estimated scale parameter of the weibit should be multiplied by  $\pi/6^{1/2}$  for comparison with the corresponding probit estimate and we should add the factor  $0.5772v$  to the constant of the weibit.

Table A-II.2 *Mean and Median<sup>a</sup> estimated wtp and their MSE*

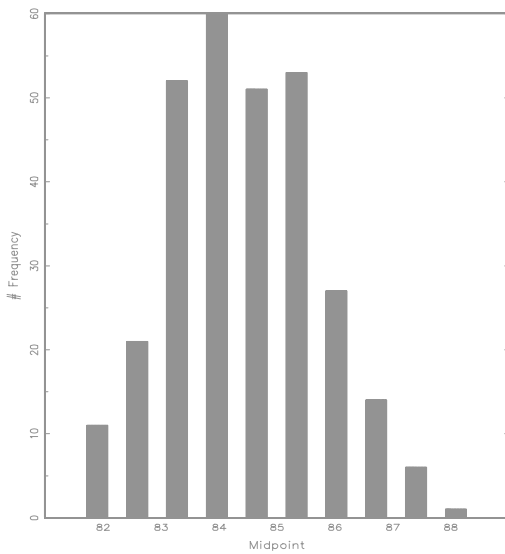
|      | $H_f$                   |  | $H_g$             |                    |                    |                    |
|------|-------------------------|--|-------------------|--------------------|--------------------|--------------------|
|      | Mean-Median<br>(st.dev) | MSE <sup>b</sup><br><i>SMSE</i> <sup>c</sup> | Mean<br>(st.dev)  | MSE<br><i>SMSE</i> | Median<br>(st.dev) | MSE<br><i>SMSE</i> |
| 300  | 86.227<br>(1.808)       | 6.242<br>3.262                               | 86.469<br>(2.165) | 8.549<br>4.703     | 83.755<br>(2.138)  | 5.110<br>10.988    |
| 600  | 84.425<br>(1.271)       | 1.616<br>1.611                               | 84.542<br>(1.416) | 2.000<br>2.011     | 81.783<br>2.138    | 9.324<br>8.942     |
| 1000 | 85.404<br>(0.924)       | 1.668<br>0.853                               | 85.477<br>(0.969) | 1.891<br>0.937     | 82.679<br>(0.964)  | 4.244<br>8.645     |

- a) For the probit model the two values coincide.  
b) MSE with respect to the population true mean wtp: 84.5  
c) MSE with respect to the sample true mean wtp: 86.292, 84.429, and 85.457 for the 300, 600 and 1000 sample size respectively.

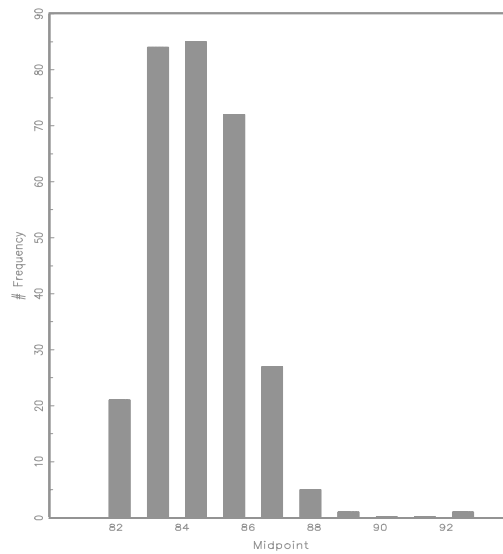
Table A-II.3 Conclusions drawn from each method (% frequency)

| Method |      | Conclusion   |  |  |  |
|--------|------|--|--|--|--|
| Akaike |      | H <sub>f</sub> is better                           |  | H <sub>g</sub> is better                           |  |
|        | 300  | 0.764  |  | 0.236  |  |
|        | 600  | 0.905  |  | 0.094  |  |
|        | 1000 | 0.983  |  | 0.017  |  |
| Vuong  |      | H <sub>f</sub> is better                           | H <sub>g</sub> is better                           |  | H <sub>f</sub> and H <sub>g</sub> equivalent       |
|        | 300  | 0.250  | 0.000  |  | 0.750  |
|        | 600  | 0.344  | 0.006  |  | 0.648  |
|        | 1000 | 0.523  | 0.000  |  | 0.476  |
| Cox    |      | H <sub>f</sub> accepted<br>H <sub>g</sub> rejected | H <sub>g</sub> accepted<br>H <sub>f</sub> rejected | Both H <sub>f</sub> and<br>H <sub>g</sub> accepted | Both H <sub>f</sub> and<br>H <sub>g</sub> rejected |
|        | 300  | 0.589  | 0.164  | 0.123  | 0.123  |
|        | 600  | 0.790  | 0.074  | 0.006  | 0.128  |
|        | 1000 | 0.885  | 0.013  | 0.003  | 0.097  |

Histogram of mean wtp  
for n=600-PROBIT



Histogram of mean wtp  
for n=600-WEIBIT



### Experiment A-III: Normal DGP, Logit vs Weibit

Table A.III.1: Parameter estimates<sup>a</sup> for normal DGP using  $H_f$  (logistic) and  $H_g$  (extreme value) across 300 replications<sup>b</sup>.

| Parameters | Sample Size       |                   |                   |                   |                   |                   |
|------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
|            | 300               |                   | 600               |                   | 1000              |                   |
|            | $H_f$             | $H_g$             | $H_f$             | $H_g$             | $H_f$             | $H_g$             |
| $\delta_1$ | 26.903<br>(8.476) | 19.276<br>(9.615) | 26.800<br>(6.417) | 19.160<br>(7.037) | 27.383<br>(4.583) | 19.649<br>(5.265) |
| $\delta_2$ | 1.512<br>(0.122)  | 1.527<br>(0.136)  | 1.494<br>(0.085)  | 1.497<br>(0.094)  | 1.490<br>(0.061)  | 1.506<br>(0.068)  |
| $\delta_3$ | -3.351<br>(4.137) | -3.638<br>(4.255) | -2.924<br>(2.953) | -2.755<br>(3.009) | -2.982<br>(2.279) | -3.251<br>(2.538) |
| $\delta_4$ | 0.347<br>(6.075)  | 0.185<br>(6.364)  | 0.999<br>(4.531)  | 0.908<br>(4.816)  | 0.538<br>(3.443)  | 0.238<br>(3.892)  |
| $v^c$      | 8.264<br>(1.020)  | 13.085<br>(2.052) | 8.273<br>(0.942)  | 13.137<br>(1.508) | 8.276<br>(0.660)  | 13.293<br>(1.124) |

- a) Mean values and standard deviations (in parenthesis) over 300 replications.  
b) The actual number of successful experiments was 288, 290 and 294 for the 300, 600 and 1000 sample size respectively.  
c) The estimated scale parameter of the logit and the weibit should be multiplied by  $\pi/3^{1/2}$  and by  $\pi/6^{1/2}$  respectively, for comparison with the corresponding probit estimate. We should add as well the factor  $0.5772v$  to the constant of the weibit.

Table A-III.2: Mean and Median<sup>a</sup> estimated wtp and their MSE

|      | $H_f$                   |  | $H_g$             |                       |                    |                        |
|------|-------------------------|--|-------------------|-----------------------|--------------------|------------------------|
|      | Mean-Median<br>(st.dev) | MSE <sup>b</sup><br><i>SMSE</i> <sup>c</sup> | Mean<br>(st.dev)  | MSE<br><i>SMSE</i>    | Median<br>(st.dev) | MSE<br><i>SMSE</i>     |
| 300  | 86.284<br>(1.772)       | 6.315<br><i>3.130</i>                        | 86.459<br>(2.044) | 8.001<br><i>4.190</i> | 83.702<br>(2.014)  | 4.678<br><i>10.748</i> |
| 600  | 84.548<br>(1.393)       | 1.936<br><i>1.948</i>                        | 84.548<br>(1.476) | 2.199<br><i>2.228</i> | 81.900<br>(1.441)  | 8.828<br><i>8.463</i>  |
| 1000 | 85.526<br>(0.940)       | 1.932<br><i>0.885</i>                        | 85.526<br>(0.978) | 2.154<br><i>0.973</i> | 82.795<br>(0.988)  | 3.881<br><i>8.060</i>  |

- a) For the logit model the two values coincide.  
b) MSE with respect to the population true mean wtp: 84.5  
c) MSE with respect to the sample true mean wtp: 86.292, 84.429, and 85.457 for the 300, 600 and 1000 sample size respectively.

Table A-III.3 *Conclusions drawn from each method (% frequency)*

| Method |      | Conclusion   |  |  |  |
|--------|------|--|--|--|--|
| Akaike |      | H <sub>f</sub> is better                           |  | H <sub>g</sub> is better                           |  |
|        | 300  | 0.753  |  | 0.247  |  |
|        | 600  | 0.813  |  | 0.186  |  |
|        | 1000 | 0.921  |  | 0.079  |  |
| Vuong  |      | H <sub>f</sub> is better                           | H <sub>g</sub> is better                           |  | H <sub>f</sub> and H <sub>g</sub> equivalent       |
|        | 300  | 0.118  | 0.034  |  | 0.847  |
|        | 600  | 0.217  | 0.020  |  | 0.762  |
|        | 1000 | 0.333  | 0.003  |  | 0.663  |
| Cox    |      | H <sub>f</sub> accepted<br>H <sub>g</sub> rejected | H <sub>g</sub> accepted<br>H <sub>f</sub> rejected | Both H <sub>f</sub> and<br>H <sub>g</sub> accepted | Both H <sub>f</sub> and<br>H <sub>g</sub> rejected |
|        | 300  | 0.392  | 0.395  | 0.038  | 0.173  |
|        | 600  | 0.493  | 0.189  | 0.006  | 0.310  |
|        | 1000 | 0.486  | 0.040  | 0.000  | 0.472  |



## **B) EXTREME VALUE DGP**

### **Experiment B-I: Extreme Value DGP, Probit vs Logit**

Table B-I.1: *Parameter estimates<sup>a</sup> for extreme value DGP using  $H_f$ (normal) and  $H_g$ (logistic) across 300 replications<sup>b</sup>.*

| Parameters | Sample Size       |                   |                   |                   |                   |                   |
|------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
|            | 300               |                   | 600               |                   | 1000              |                   |
|            | $H_f$             | $H_g$             | $H_f$             | $H_g$             | $H_f$             | $H_g$             |
| $\delta_1$ | 26.270<br>(8.009) | 26.195<br>(7.870) | 26.596<br>(5.539) | 26.432<br>(5.463) | 26.183<br>(4.576) | 26.152<br>(4.489) |
| $\delta_2$ | 1.522<br>(0.117)  | 1.515<br>(0.114)  | 1.514<br>(0.078)  | 1.507<br>(0.077)  | 1.515<br>(0.065)  | 1.506<br>(0.064)  |
| $\delta_3$ | -2.783<br>(4.037) | -2.725<br>(3.987) | -3.140<br>(2.837) | -3.145<br>(2.769) | -3.226<br>(2.197) | -3.152<br>(2.163) |
| $\delta_4$ | 0.414<br>(6.088)  | 0.325<br>(5.965)  | 0.593<br>(4.318)  | 0.634<br>(4.184)  | 1.0<br>(3.422)    | 0.918<br>(3.388)  |
| $v^c$      | 14.624<br>(2.316) | 8.013<br>(1.226)  | 14.858<br>(1.475) | 8.127<br>(0.798)  | 14.839<br>(1.202) | 8.103<br>(0.682)  |

- a) Mean values and standard deviations (in parenthesis) over 300 replications.  
b) The actual number of successful experiments was 274 for the 300, and 291 for the 600 and 1000 sample size.  
c) The estimated scale parameter of the logit should be multiplied by  $\pi/3^{1/2}$  for comparison with the corresponding probit estimate.

Table B-I.2: *Mean, median estimated wtp and their MSE*

|      | $H_f$                   |  |  | $H_g$                   |                            |                              |
|------|-------------------------|--|--|-------------------------|----------------------------|------------------------------|
|      | Mean-Median<br>(st.dev) | MSE <sup>a</sup><br><i>SMSE</i> <sup>b</sup><br>Mean | MSE <sup>c</sup><br><i>SMSE</i> <sup>d</sup><br>Median | Mean-Median<br>(st.dev) | MSE<br><i>SMSE</i><br>Mean | MSE<br><i>SMSE</i><br>Median |
|      | 300                     | 86.531<br>(1.775)                                    | 7.268<br>3.197   | 23.354<br>10.452        | 86.150<br>(1.729)          | 5.704<br>3.001               |
| 600  | 84.594<br>(1.250)       | 1.568<br>1.586                                       | 8.107<br>8.476   | 84.198<br>(1.247)       | 1.641<br>1.603             | 6.228<br>6.541               |
| 1000 | 85.583<br>(1.032)       | 2.234<br>1.077                                       | 13.645<br>7.772  | 85.166<br>(1.015)       | 1.471<br>1.112             | 10.829<br>5.753              |

- a) MSE with respect to the population true mean wtp: 84.5  
b) MSE with respect to the sample true mean wtp: 86.292, 84.429, and 85.457 for the 300, 600 and 1000 sample size respectively.  
c) MSE with respect to the population true median wtp: 82.03.  
d) MSE with respect to the sample true median wtp: 83.827, 81.964, and 82.992 for the 300, 600 and 1000 sample size respectively.

Table B-I.3: *Conclusions drawn from each method (% frequency)*

| Method |      | Conclusion   |  |  |  |
|--------|------|--|--|--|--|
| Akaike |      | H <sub>f</sub> is better                           |  | H <sub>g</sub> is better                           |  |
|        | 300  | 0.485  |  | 0.514  |  |
|        | 600  | 0.344  |  | 0.656  |  |
|        | 1000 | 0.251  |  | 0.749  |  |
| Vuong  |      | H <sub>f</sub> is better                           | H <sub>g</sub> is better                           |  | H <sub>f</sub> and H <sub>g</sub> equivalent       |
|        | 300  | 0.124  | 0.036  |  | 0.839  |
|        | 600  | 0.048  | 0.054  |  | 0.896  |
|        | 1000 | 0.031  | 0.065  |  | 0.904  |
| Cox    |      | H <sub>f</sub> accepted<br>H <sub>g</sub> rejected | H <sub>g</sub> accepted<br>H <sub>f</sub> rejected | Both H <sub>f</sub> and<br>H <sub>g</sub> accepted | Both H <sub>f</sub> and<br>H <sub>g</sub> rejected |
|        | 300  | 0.248  | 0.106  | 0.500  | 0.146  |
|        | 600  | 0.158  | 0.140  | 0.642  | 0.058  |
|        | 1000 | 0.120  | 0.168  | 0.656  | 0.054  |

## Experiment B-II: Extreme value DGP, Probit vs Weibit

Table B-II.1: *Parameter estimates<sup>a</sup> for the extreme value DGP using  $H_f$  (normal) and  $H_g$  (extreme value) across 300 replications<sup>b</sup>.*

| Parameters | Sample Size       |                   |                   |                   |                   |                   |
|------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
|            | 300               |                   | 600               |                   | 1000              |                   |
|            | $H_f$             | $H_g$             | $H_f$             | $H_g$             | $H_f$             | $H_g$             |
| $\delta_1$ | 26.905<br>(8.729) | 20.531<br>(8.384) | 26.662<br>(5.726) | 20.475<br>(5.777) | 26.904<br>(4.490) | 20.323<br>(4.314) |
| $\delta_2$ | 1.503<br>(0.128)  | 1.504<br>(0.119)  | 1.513<br>(0.080)  | 1.506<br>(0.076)  | 1.499<br>(0.064)  | 1.497<br>(0.060)  |
| $\delta_3$ | -2.935<br>(4.056) | -3.206<br>(3.811) | -3.005<br>(2.782) | -3.085<br>(2.584) | -2.954<br>(2.165) | -3.039<br>(1.908) |
| $\delta_4$ | 0.486<br>(6.071)  | 0.488<br>(5.550)  | 0.507<br>(4.572)  | 0.490<br>(4.270)  | 0.709<br>(3.480)  | 0.713<br>(3.143)  |
| $v^c$      | 14.326<br>(2.351) | 11.229<br>(1.737) | 14.707<br>(1.422) | 11.518<br>(1.211) | 14.847<br>(1.268) | 11.677<br>(1.030) |

- a) Mean values and standard deviations (in parenthesis) over 300 replications.  
b) The actual number of successful experiments was 294 for the 300 and 600 sample sizes and 290 for the 1000 sample size.  
c) the 1000 sample size.  
d) The estimated scale parameter of the weibit should be multiplied by  $\pi/6^{1/2}$  for comparison with the corresponding probit estimate and we should add the factor  $0.5772v$  to the constant of the weibit.

Table B-II.2: *Mean and Median<sup>a</sup> estimated wtp and their MSE*

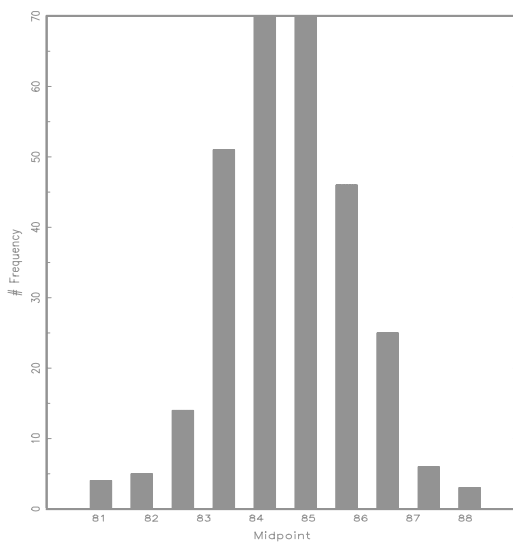
|      | $H_f$                   |                                  |                                    | $H_g$             |                            |                    |                              |
|------|-------------------------|----------------------------------|------------------------------------|-------------------|----------------------------|--------------------|------------------------------|
|      | Mean-Median<br>(st.dev) | MSE <sup>b</sup>                 | MSE <sup>d</sup>                   | Mean<br>(st.dev)  | MSE<br><i>SMSE</i><br>Mean | Median<br>(st.dev) | MSE<br><i>SMSE</i><br>Median |
|      |                         | <i>SMSE</i> <sup>c</sup><br>Mean | <i>SMSE</i> <sup>e</sup><br>Median |                   |                            |                    |                              |
| 300  | 86.361<br>(1.573)       | 5.932<br>2.472                   | 21.18<br>8.888                     | 86.312<br>(1.704) | 6.181<br>2.895             | 83.947<br>(1.654)  | 6.38<br>2.742                |
| 600  | 84.623<br>(1.262)       | 1.605<br>1.627                   | 8.288<br>8.661                     | 84.722<br>(1.881) | 3.575<br>3.612             | 82.295<br>(1.775)  | 3.208<br>3.250               |
| 1000 | 85.583<br>(0.960)       | 2.091<br>0.934                   | 13.502<br>7.629                    | 85.619<br>(1.477) | 3.426<br>2.200             | 83.159<br>(1.391)  | 3.190<br>1.956               |

- a) For the probit model the two values coincide.  
b) MSE with respect to the population true mean wtp: 84.5  
c) MSE with respect to the sample true mean wtp: 86.292, 84.429, and 85.457 for the 300, 600 and 1000 sample size respectively.  
d) MSE with respect to the population true median wtp: 82.03  
e) MSE with respect to the sample true mean wtp: 83.827, 81.964, and 82.992 for the 300, 600 and 1000 sample size respectively.

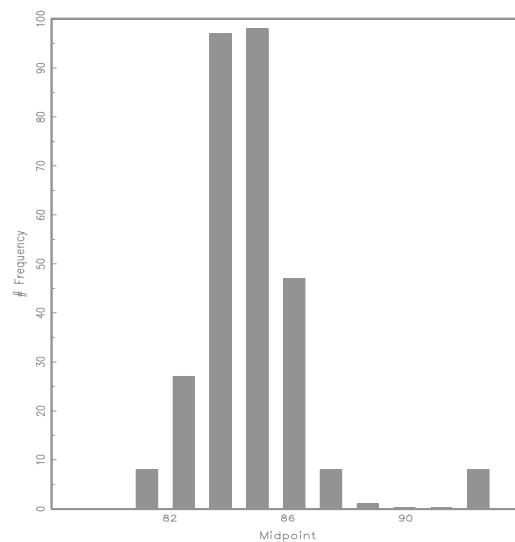
Table B-II.3: *Conclusions drawn from each method (% frequency)*

| Method |      | Conclusion   |  |  |  |
|--------|------|--|--|--|--|
| Akaike |      | H <sub>f</sub> is better                           |  | H <sub>g</sub> is better                           |  |
|        | 300  | 0.190  |  | 0.810  |  |
|        | 600  | 0.092  |  | 0.908  |  |
|        | 1000 | 0.059  |  | 0.941  |  |
| Vuong  |      | H <sub>f</sub> is better                           | H <sub>g</sub> is better                           | H <sub>f</sub> and H <sub>g</sub> equivalent       |  |
|        | 300  | 0.020  | 0.204  | 0.775  |  |
|        | 600  | 0.017  | 0.353  | 0.629  |  |
|        | 1000 | 0.017  | 0.400  | 0.583  |  |
| Cox    |      | H <sub>f</sub> accepted<br>H <sub>g</sub> rejected | H <sub>g</sub> accepted<br>H <sub>f</sub> rejected | Both H <sub>f</sub> and<br>H <sub>g</sub> accepted | Both H <sub>f</sub> and<br>H <sub>g</sub> rejected |
|        | 300  | 0.119  | 0.717  | 0.068  | 0.095  |
|        | 600  | 0.030  | 0.826  | 0.020  | 0.122  |
|        | 1000 | 0.010  | 0.862  | 0.003  | 0.124  |

Histogram of mean wtp  
for n=600-PROBIT



Histogram of mean wtp  
for n=600-WEIBIT



### Experiment B-III: Extreme value DGP, Logit vs Weibit

Table B-III.1: *Parameter estimates<sup>a</sup> for the extreme value DGP using  $H_f$  (logistic) and  $H_g$  (extreme value) across 300 replications<sup>b</sup>.*

| Parameters | Sample Size       |                   |                   |                   |                   |                   |
|------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
|            | 300               |                   | 600               |                   | 1000              |                   |
|            | $H_f$             | $H_g$             | $H_f$             | $H_g$             | $H_f$             | $H_g$             |
| $\delta_1$ | 26.364<br>(8.567) | 20.052<br>(8.479) | 26.291<br>(5.553) | 20.280<br>(5.648) | 26.798<br>(4.523) | 20.269<br>(4.389) |
| $\delta_2$ | 1.515<br>(0.112)  | 1.519<br>(0.110)  | 1.506<br>(0.081)  | 1.506<br>(0.078)  | 1.496<br>(0.062)  | 1.502<br>(0.059)  |
| $\delta_3$ | -3.055<br>(4.152) | -3.248<br>(3.974) | -3.010<br>(2.783) | -3.087<br>(2.670) | -2.995<br>(2.168) | -3.049<br>(1.966) |
| $\delta_4$ | 0.351<br>(6.283)  | 0.386<br>(5.884)  | 0.712<br>(4.439)  | 0.596<br>(4.168)  | 0.574<br>(3.443)  | 0.556<br>(3.193)  |
| $v^c$      | 8.033<br>(1.205)  | 11.485<br>(1.656) | 8.065<br>(0.787)  | 11.501<br>(1.183) | 8.145<br>(0.737)  | 11.656<br>(0.995) |

- a) Mean values and standard deviations (in parenthesis) over 300 replications.  
b) The actual number of successful experiments was 282, 294 and 286 for the 300, 600 and 1000 sample size respectively.  
c) The estimated scale parameter of the logit and the weibit should be multiplied by  $\pi/3^{1/2}$  and by  $\pi/6^{1/2}$  respectively, for comparison with the corresponding probit estimate. We should add as well the factor  $0.5772v$  to the constant of the weibit.

Table B-III.2: *Mean and Median<sup>a</sup> estimated wtp and their MSE*

|      | $H_f$                   |                                 |                                   | $H_g$             |                       |                    |                       |
|------|-------------------------|---------------------------------|-----------------------------------|-------------------|-----------------------|--------------------|-----------------------|
|      | Mean-Median<br>(st.dev) | MSE <sup>a</sup>                | MSE <sup>c</sup>                  | Mean<br>(st.dev)  | MSE                   | Median<br>(st.dev) | MSE                   |
|      |                         | <i>SMSE<sup>b</sup></i><br>Mean | <i>SMSE<sup>d</sup></i><br>Median |                   | <i>SMSE</i><br>Mean   |                    | <i>SMSE</i><br>Median |
| 300  | 86.077<br>(1.836)       | 5.851<br><i>3.407</i>           | 19.7<br><i>8.425</i>              | 86.468<br>(2.037) | 8.010<br><i>4.166</i> | 84.048<br>(1.950)  | 7.841<br><i>3.838</i> |
| 600  | 84.195<br>(1.243)       | 1.634<br><i>1.595</i>           | 6.206<br><i>6.518</i>             | 84.651<br>(1.740) | 3.043<br><i>3.069</i> | 82.227<br>(1.646)  | 2.740<br><i>2.772</i> |
| 1000 | 85.189<br>(0.949)       | 1.373<br><i>0.969</i>           | 10.844<br><i>5.724</i>            | 85.587<br>(1.406) | 3.155<br><i>1.989</i> | 83.132<br>(1.309)  | 2.911<br><i>1.729</i> |

- a) For the logit model the two values coincide.  
b) MSE with respect to the population true mean wtp: 84.5  
c) MSE with respect to the sample true mean wtp: 86.292, 84.429, and 85.457 for the 300, 600 and 1000 sample size respectively.  
d) MSE with respect to the population true median wtp: 82.03  
e) MSE with respect to the sample true median wtp: 83.827, 81.964, and 82.992 for the 300, 600 and 1000 sample size respectively.

Table B-III.3: *Conclusions drawn from each method (% frequency)*

| Method |      | Conclusion   |  |  |  |
|--------|------|--|--|--|--|
| Akaike |      | H <sub>f</sub> is better                           |  | H <sub>g</sub> is better                           |  |
|        | 300  | 0.160  |  | 0.840  |  |
|        | 600  | 0.061  |  | 0.938  |  |
|        | 1000 | 0.031  |  | 0.969  |  |
| Vuong  |      | H <sub>f</sub> is better                           | H <sub>g</sub> is better                           |  | H <sub>f</sub> and H <sub>g</sub> equivalent       |
|        | 300  | 0.003  | 0.330  |  | 0.667  |
|        | 600  | 0.017  | 0.476  |  | 0.507  |
|        | 1000 | 0.013  | 0.555  |  | 0.430  |
| Cox    |      | H <sub>f</sub> accepted<br>H <sub>g</sub> rejected | H <sub>g</sub> accepted<br>H <sub>f</sub> rejected | Both H <sub>f</sub> and<br>H <sub>g</sub> accepted | Both H <sub>f</sub> and<br>H <sub>g</sub> rejected |
|        | 300  | 0.025  | 0.830  | 0.014  | 0.131  |
|        | 600  | 0.006  | 0.877  | 0.010  | 0.105  |
|        | 1000 | 0.003  | 0.888  | 0.000  | 0.108  |

### C. LOGNORMAL DGP

#### Experiment C-I: Lognormal DGP, Probit vs Logit

Table C-I.1: Parameter estimates<sup>a</sup> for the lognormal DGP using  $H_f$  (normal) and  $H_g$  (logistic) across 300 replications<sup>b</sup>.

| Parameters | Sample Size       |                   |                   |                   |                   |                   |
|------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
|            | 300               |                   | 600               |                   | 1000              |                   |
|            | $H_f$             | $H_g$             | $H_f$             | $H_g$             | $H_f$             | $H_g$             |
| $\delta_1$ | 26.572<br>(7.758) | 26.306<br>(6.804) | 25.246<br>(5.379) | 25.326<br>(4.768) | 25.952<br>(4.258) | 25.930<br>(3.763) |
| $\delta_2$ | 1.511<br>(0.103)  | 1.500<br>(0.092)  | 1.535<br>(0.07)   | 1.515<br>(0.062)  | 1.529<br>(0.057)  | 1.510<br>(0.052)  |
| $\delta_3$ | -2.813<br>(3.766) | -2.808<br>(3.328) | -2.952<br>(2.672) | -2.965<br>(2.348) | -3.095<br>(2.057) | -3.019<br>(1.769) |
| $\delta_4$ | 0.529<br>(5.985)  | 0.510<br>(5.196)  | 1.018<br>(4.420)  | 0.946<br>(3.857)  | 0.695<br>(3.459)  | 0.623<br>(2.963)  |
| $v^c$      | 12.202<br>(3.073) | 6.091<br>(1.390)  | 13.147<br>(2.143) | 6.468<br>(0.962)  | 13.216<br>(1.835) | 6.455<br>(0.814)  |

- a) Mean values and standard deviations (in parenthesis) over 300 replications.  
b) The actual number of successful experiments was 258, 261 and 247 for the 300, 600 and 1000 sample size respectively.  
c) The estimated scale parameter of the logit should be multiplied by  $\pi/3^{1/2}$  for comparison with the corresponding probit estimate.

Table C-I.2: Mean and Median estimated wtp and their MSE

|      | $H_f$                   |  |  | $H_g$                   |                            |                              |
|------|-------------------------|--|--|-------------------------|----------------------------|------------------------------|
|      | Mean-Median<br>(st.dev) | MSE <sup>a</sup><br><i>SMSE</i> <sup>b</sup><br>Mean | MSE <sup>c</sup><br><i>SMSE</i> <sup>d</sup><br>Median | Mean-Median<br>(st.dev) | MSE<br><i>SMSE</i><br>Mean | MSE<br><i>SMSE</i><br>Median |
|      | 300                     | 86.483<br>(1.591)                                    | 6.456<br>2.559   | 44.588<br>24.555        | 85.773<br>(1.463)          | 3.755<br>2.401               |
| 600  | 84.633<br>(1.096)       | 1.214<br>1.238                                       | 22.688<br>23.353                                       | 83.851<br>(1.044)       | 1.507<br>1.419             | 15.934<br>16.488             |
| 1000 | 85.715<br>(0.846)       | 2.192<br>0.781                                       | 33.412<br>23.385                                       | 84.92<br>(0.764)        | 0.758<br>0.869             | 24.818<br>16.312             |

- a) MSE with respect to the population true mean wtp: 84.5  
b) MSE with respect to the sample true mean wtp: 86.292, 84.429, and 85.457 for the 300, 600 and 1000 sample size respectively.  
c) MSE with respect to the population true median wtp: 80.  
d) MSE with respect to the sample true median wtp: 81.789, 79.926, and 80.954 for the 300, 600 and 1000 sample size respectively.

Table C-I.3: *Conclusions drawn from each method (% frequency)*

| Method |      | Conclusion   |  |  |  |
|--------|------|--|--|--|--|
| Akaike |      | H <sub>f</sub> is better                           |  | H <sub>g</sub> is better                           |  |
|        | 300  | 0.097  |  | 0.903  |  |
|        | 600  | 0.015  |  | 0.984  |  |
|        | 1000 | 0.008  |  | 0.992  |  |
| Vuong  |      | H <sub>f</sub> is better                           | H <sub>g</sub> is better                           |  | H <sub>f</sub> and H <sub>g</sub> equivalent       |
|        | 300  | 0.011  | 0.422  |  | 0.565  |
|        | 600  | 0.000  | 0.704  |  | 0.295  |
|        | 1000 | 0.000  | 0.854  |  | 0.146  |
| Cox    |      | H <sub>f</sub> accepted<br>H <sub>g</sub> rejected | H <sub>g</sub> accepted<br>H <sub>f</sub> rejected | Both H <sub>f</sub> and<br>H <sub>g</sub> accepted | Both H <sub>f</sub> and<br>H <sub>g</sub> rejected |
|        | 300  | 0.050  | 0.259  | 0.232  | 0.457  |
|        | 600  | 0.004  | 0.210  | 0.100  | 0.678  |
|        | 1000 | 0.000  | 0.133  | 0.024  | 0.842  |



## Experiment C-II: Lognormal DGP, Probit vs Weibit

Table C-II.1: *Parameter estimates<sup>a</sup> for lognormal DGP using  $H_f$ (normal) and  $H_g$  (extreme value) across 300 replications<sup>b</sup>.*

| Parameters | Sample Size       |                   |                   |                   |                   |                   |
|------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
|            | 300               |                   | 600               |                   | 1000              |                   |
|            | $H_f$             | $H_g$             | $H_f$             | $H_g$             | $H_f$             | $H_g$             |
| $\delta_1$ | 26.343<br>(8.456) | 21.789<br>(6.617) | 25.832<br>(5.324) | 22.214<br>(3.985) | 26.391<br>(4.191) | 22.240<br>(3.169) |
| $\delta_2$ | 1.511<br>(0.106)  | 1.502<br>(0.085)  | 1.530<br>(0.079)  | 1.501<br>(0.058)  | 1.527<br>(0.055)  | 1.505<br>(0.042)  |
| $\delta_3$ | -2.825<br>(3.892) | -2.846<br>(2.918) | -2.876<br>(2.722) | -3.047<br>(1.947) | -2.848<br>(1.934) | -2.771<br>(1.386) |
| $\delta_4$ | 0.646<br>(5.784)  | 0.786<br>(4.058)  | 0.592<br>(4.186)  | 0.475<br>(2.866)  | 0.225<br>(3.372)  | 0.162<br>(2.462)  |
| $v^c$      | 12.513<br>(3.005) | 8.046<br>(1.636)  | 13.305<br>(1.978) | 8.428<br>(1.150)  | 13.210<br>(1.581) | 8.259<br>(0.768)  |

- a) Mean values and standard deviations (in parenthesis) over 300 replications.  
b) The actual number of successful experiments was 285, 281 and 287 for the 300, 600 and 1000 sample sizes respectively.  
c) The estimated scale parameter of the weibit should be multiplied by  $\pi/6^{1/2}$  for comparison with the corresponding probit estimate and we should add the factor  $0.5772v$  to the constant of the weibit.

Table C-II.2: *Mean and Median<sup>a</sup> estimated wtp and their MSE*

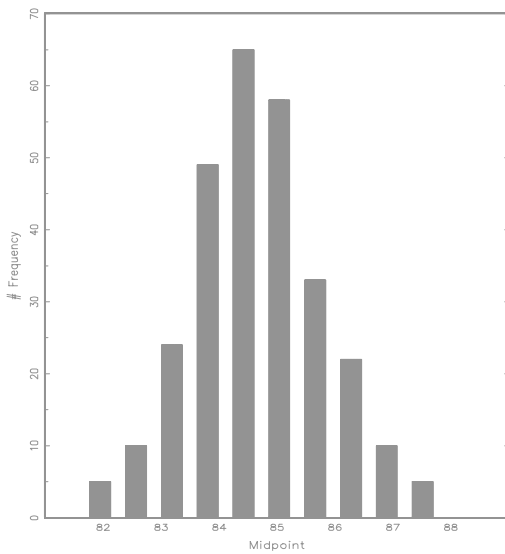
|      | $H_f$                   |   |   | $H_g$             |                       |                    |                        |
|------|-------------------------|---|---|-------------------|-----------------------|--------------------|------------------------|
|      | Mean-Median<br>(st.dev) | MSE <sup>a</sup><br><i>SMSE<sup>b</sup></i><br>Mean | MSE <sup>c</sup><br><i>SMSE<sup>c</sup></i><br>Median | Mean<br>(st.dev)  | MSE<br><i>SMSE</i>    | Median<br>(st.dev) | MSE<br><i>SMSE</i>     |
| 300  | 86.376<br>(1.704)       | 6.416<br><i>2.901</i>                               | 43.588<br><i>23.937</i>                               | 86.215<br>(1.631) | 5.596<br><i>2.658</i> | 84.52<br>(1.455)   | 22.571<br><i>9.571</i> |
| 600  | 84.642<br>(1.138)       | 1.312<br><i>1.337</i>                               | 22.866<br><i>23.533</i>                               | 84.518<br>(1.319) | 1.734<br><i>1.742</i> | 82.742<br>(1.167)  | 8.892<br><i>9.288</i>  |
| 1000 | 85.780<br>(0.813)       | 2.298<br><i>0.763</i>                               | 34.102<br><i>23.951</i>                               | 85.511<br>(0.753) | 1.588<br><i>0.568</i> | 83.771<br>(0.671)  | 14.689<br><i>8.383</i> |

- a) For the probit model the two values coincide.  
b) MSE with respect to the population true mean wtp: 84.5  
c) MSE with respect to the sample true mean wtp: 86.292, 84.429, and 85.457 for the 300, 600 and 1000 sample size respectively.  
d) MSE with respect to the population true median wtp: 80.  
e) MSE with respect to the sample true median wtp: 81.789, 79.926, and 80.954 for the 300, 600 and 1000 sample size respectively.

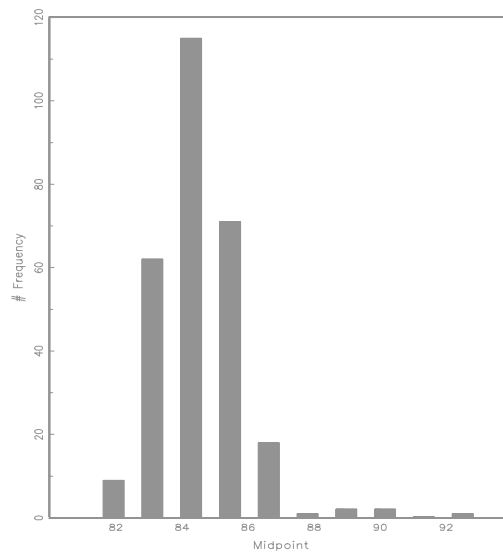
Table C-II.3: *Conclusions drawn from each method (% frequency)*

| Method |      | Conclusion   |  |  |  |
|--------|------|--|--|--|--|
| Akaike |      | H <sub>f</sub> is better                           |  | H <sub>g</sub> is better                           |  |
|        | 300  | 0.014  |  | 0.985  |  |
|        | 600  | 0.011  |  | 0.989  |  |
|        | 1000 | 0.000  |  | 1.000  |  |
| Vuong  |      | H <sub>f</sub> is better                           | H <sub>g</sub> is better                           |  | H <sub>f</sub> and H <sub>g</sub> equivalent       |
|        | 300  | 0.003  | 0.905  |  | 0.091  |
|        | 600  | 0.003  | 0.982  |  | 0.014  |
|        | 1000 | 0.000  | 1.000  |  | 0.000  |
| Cox    |      | H <sub>f</sub> accepted<br>H <sub>g</sub> rejected | H <sub>g</sub> accepted<br>H <sub>f</sub> rejected | Both H <sub>f</sub> and<br>H <sub>g</sub> accepted | Both H <sub>f</sub> and<br>H <sub>g</sub> rejected |
|        | 300  | 0.003  | 0.245  | 0.003  | 0.747  |
|        | 600  | 0.000  | 0.078  | 0.003  | 0.918  |
|        | 1000 | 0.000  | 0.020  | 0.000  | 0.980  |

Histogram of mean wtp  
for n=600-PROBIT



Histogram of mean wtp  
for n=600-WEIBIT



## APPENDIX 2

### **Derivation of the simulation based Cox test statistic.**

Let  $F$  and  $G$  be two candidate cumulative distribution functions for the single bound model, with corresponding density functions  $f$  and  $g$ , which give rise to the two models for the willingness to pay below:

$$H_f : Y_i = \mathbf{x}'_i \boldsymbol{\delta}_1 + \varepsilon_{1i}, \varepsilon_{1i} \sim F \text{ and } P(Y_i > t_1 | \mathbf{x}_i) = 1 - F((t_1 - \mathbf{x}'_i \boldsymbol{\delta}_1)/v_1) = 1 - F_i(\boldsymbol{\beta}), \text{ and}$$

$$H_g : Y_i = \mathbf{x}'_i \boldsymbol{\delta}_2 + \varepsilon_{2i}, \varepsilon_{2i} \sim G \text{ and } P(Y_i > t_1 | \mathbf{x}_i) = 1 - G((t_1 - \mathbf{x}'_i \boldsymbol{\delta}_2)/v_2) = 1 - G_i(\boldsymbol{\gamma}).$$

These two models differ only in the specification of the distribution of the error term. The test statistic for the null hypothesis that the true DGP belongs to  $H_f$  against the alternative that it belongs to  $H_g$  is given by

$$S_f(\hat{\boldsymbol{\beta}}_n, \hat{\boldsymbol{\gamma}}_n) = \sqrt{n} T_f / \hat{v}_f, \text{ where}$$

$$T_f = \frac{1}{n} LR_n(\hat{\boldsymbol{\beta}}_n, \hat{\boldsymbol{\gamma}}_n) - \hat{E}_f \left( \frac{1}{n} LR_n(\hat{\boldsymbol{\beta}}_n, \hat{\boldsymbol{\gamma}}_n) \right)$$

Under  $H_f$  the log-likelihood is given by,

$$\log \mathfrak{L}(f(\boldsymbol{\beta})) = \sum_{i=1}^n I_i \log [1 - F_i(\boldsymbol{\beta})] + (1 - I_i) \log F_i(\boldsymbol{\beta}), \text{ while under } H_g \text{ we have,}$$

$$\log \mathfrak{L}(g(\boldsymbol{\gamma})) = \sum_{i=1}^n I_i \log [1 - G_i(\boldsymbol{\gamma})] + (1 - I_i) \log G_i(\boldsymbol{\gamma}). \text{ Let } \hat{\boldsymbol{\beta}}_n \text{ and } \hat{\boldsymbol{\gamma}}_n \text{ be the corresponding}$$

maximum likelihood estimators, then the log-likelihood ratio is given by,

$$LR_n(\hat{\boldsymbol{\beta}}_n, \hat{\boldsymbol{\gamma}}_n) = \sum_{i=1}^n I_i \log \left( \frac{1 - F_i(\hat{\boldsymbol{\beta}}_n)}{1 - G_i(\hat{\boldsymbol{\gamma}}_n)} \right) + (1 - I_i) \log \left( \frac{F_i(\hat{\boldsymbol{\beta}}_n)}{G_i(\hat{\boldsymbol{\gamma}}_n)} \right).$$

In order to compute the second term in the numerator of the Cox statistic we need to find an estimator of that value of  $\gamma$  - call it  $\gamma^*$  - that maximizes the expected log-likelihood of  $H_g$  when  $H_f$  is true. The simulation method proposed by P&P(1993) consists of generating  $R$  samples of  $n \times 1$  indicators according to the distribution  $F(\hat{\beta}_n)$ , i.e. for each of the  $R$  replications we have  $I_j = (I_{1j}, \dots, I_{nj})$ ,  $j=1, \dots, R$  where

$$I_{ij} = \begin{cases} 1 & \text{if } 1 - F_i(\hat{\beta}_n) > U \\ 0 & \text{otherwise} \end{cases}, \text{ and } U \text{ is drawn from a uniform distribution from } (0,1).$$

For each one of the  $R$  samples  $I_j$  is used as the new data to compute the estimator  $\hat{\gamma}_j^*$  which maximizes the expression below,

$$\sum_{i=1}^n I_{ij} \log[1 - G_i(\gamma)] + (1 - I_{ij}) \log G_i(\gamma).$$

The average of the estimates  $\hat{\gamma}_j^*$  over the  $R$  replications, i.e.  $\hat{\gamma}^*(R) = \frac{1}{R} \sum_{j=1}^R \hat{\gamma}_j^*$ , can be used

as an estimate of  $\gamma^*$ , therefore we have,

$$\hat{E}_f \left( \frac{1}{n} \sum_{i=1}^n I_{ij} \log \left[ \frac{1 - F_i(\hat{\beta}_n)}{1 - G_i(\hat{\gamma}^*(R))} \right] + F_i(\hat{\beta}_n) \log \left[ \frac{F_i(\hat{\beta}_n)}{G_i(\hat{\gamma}^*(R))} \right] \right).$$

The denominator of the Cox statistic -  $\hat{v}_f$  - is obtained as the standard error (not corrected for the loss of degrees of freedom) of a regression of

$$I_i \log \left[ \frac{1 - F_i(\hat{\beta}_n)}{1 - G_i(\hat{\gamma}_n)} \right] + (1 - I_i) \log \left[ \frac{F_i(\hat{\beta}_n)}{G_i(\hat{\gamma}_n)} \right] \text{ on a constant and the vector of derivatives}$$

$$\left\{ \frac{\partial (I_i \log(1 - F_i(\hat{\beta}_n)) + (1 - I_i) \log(F_i(\hat{\beta}_n)))}{\partial \beta} \right\}.$$

The Cox statistic when we reverse the roles of the null and alternative hypotheses can be derived in a similar fashion by reversing the roles of  $F$  and  $G$ , and it requires simulating samples of indicators under the assumption that  $G$  is correct.