

# 1                   **Preference and WTP stability for public forest management**

2

3   **Abstract:** The assumption of the stability of preferences is fundamental to consumer theory and the use  
4 of cost-benefit analysis. Many papers within the stated preferences literature have tested this assumption,  
5 and have found mixed results. Individuals may become more sure of their preferences as they repeat a  
6 valuation task or purchase decision; they may also learn more about prices and quantities of substitutes  
7 or complements over time, or about other relevant characteristics of both the good being valued and  
8 alternatives in their choice sets. In this paper, we test for the stability of preferences and willingness to  
9 pay for attributes of forest management both within one survey and between two different moments of  
10 time. The “within survey” test compares a set of responses from individuals over the sequence of the  
11 first 12 and the second 12 choices in a stated preference survey; the “between two different moments of  
12 time” test compares responses from the same people over a period of 6 months. Non-parametric analysis  
13 reveals little clear trending in choices across these sets, and higher consistency for status quo choices  
14 than for enhanced environmental management choices. Overall, we reject the strictest test of the  
15 equivalence of WTP distributions between choice sets. However, we also find that respondents’ mean  
16 willingness to pay is fairly stable both within survey and between moments of time. Such differences as  
17 emerge are mainly driven by the changes in variances of WTP and by imperfect correlations of  
18 individual-specific WTP between choice sets.

19   **JEL classification:** D01, H4, Q23, Q51

20   **Keywords:** preference stability, test-retest, discrete choice experiment, contingent valuation, stated  
21 preferences, forestry

## 22   **Highlights:**

- 23       • Preference stability is an important assumption in welfare economics and the economic theory  
24       of value. We test this assumption using evidence from a discrete choice experiment study of  
25       forest management.
- 26       • Using the same group of people, we compare willingness to pay estimates derived from two sets  
27       of choice tasks presented in one survey, and with estimates obtained from the same people six  
28       months later.
- 29       • We analyse the dynamics of responses to choice tasks, showing how and when people changed  
30       their minds.
- 31       • We find that average willingness to pay does not significantly change both within survey and  
32       between two different moments of time. However, the shapes of the distributions may differ.

# 1 **1. Introduction**

2 [Brouwer \(2012\)](#) notes that whilst “... in micro-economic theory, it is assumed that individuals know  
3 their preferences and *that these preferences are stable ...*”, the consensus from behavioural psychology  
4 is that individuals are continually (re-) constructing their preferences in a context-dependent manner.  
5 This implies that preferences for the same good, and willingness to pay (WTP) for a particular change  
6 in that good, might well vary over time for an individual, even if the time span over which preferences  
7 are observed is very short. Standard economic theory allows for WTP estimates to change as variables  
8 which co-determine one’s demand for a good change, or as one learns more about the characteristics of  
9 a good ([Munro and Hanley, 2002](#)) or one’s preferences for experience goods ([Czajkowski, Hanley and  
10 LaRiviere, 2014](#)). However, in the standard model preference parameters are supposed to be stable  
11 ([McFadden, 2001](#)). This is a crucial assumption when valuation of a public good is conducted in order  
12 to inform policy makers. If preferences are unstable such that willingness to pay for a specific change  
13 in the quantity of a public or private good varies even though there is no change in any of the standard  
14 economic drivers of welfare measures, then benefit-cost analysis is no longer informative as to the  
15 efficiency implications of policy change or changes in environmental management. For example,  
16 changes in stated willingness to pay due to variations in the emotional condition of a respondent would  
17 mean that the Kaldor-Hicks potential compensation test could no longer be applied (since whether gains  
18 exceeded losses would depend on un-observable variations in context). Our study sheds some light on  
19 validity of valuation methods with regard to preference stability assumption, since we test both the  
20 stability of an individual’s willingness to pay for a good across a sequence of choice tasks in an initial  
21 survey, and across a 6-month period between this initial survey and a follow-up survey.

22 Specific tests for preference stability over environmental goods can be found in both contingent  
23 valuation (CV) and discrete choice experiment (DCE) settings. CV test-retest procedures were  
24 conducted among others by [Loomis \(1989\)](#), [Carson et al. \(1997\)](#), [Brouwer \(2006\)](#) and [Brouwer \(2012\)](#).  
25 In all cases two surveys were carried out over an interval ranging from two weeks to two years. The  
26 results in all cases indicate that the average WTP is stable.

27 Test-retest procedures have also been applied within DCE. [Bliem, Getzner and Rodiga-Laßnig \(2012\)](#)  
28 estimate multinomial and mixed logit models on samples from two surveys of river restoration options  
29 in Australia, where the two surveys were undertaken one year apart. The model coefficients were  
30 compared using a Chow test. This indicated that there was no difference between preferences in these  
31 samples. [Liebe, Meyerhoff and Hartje \(2012\)](#) used an Error Component model to compare preference  
32 and WTP estimates in two samples collected 11 months apart. Choices over on-shore wind power  
33 options were reasonably consistent over the interval, but WTP estimates differed significantly for around  
34 half of the attribute values. [Schaafsma et al. \(forthcoming\)](#) used a one-year interval to conduct a test-

1 retest CE survey, and found that there were no significant changes in either preference parameters or  
2 WTP over this interval. However, the estimated error variance of choices fell over time. Most recently,  
3 [Mørkbak and Olsen \(2014\)](#) used DCE to compare responses over a 2 week period for a market good  
4 (apples) with “real economic incentives”. They thus sought to undertake a test-retest experiment in an  
5 incentive-compatible setting. They found “very good agreement” between the DCE estimates of  
6 preferences over this rather short time interval. However, their sample consisted of 25 persons only.

7 Further relevant contributions include [Dupont, Price and Adamowicz \(2014\)](#), who compare estimates of  
8 WTP for health end points related to water quality in Canada between surveys undertaken in 2004 and  
9 2012, using both CV and DCE. The health end points relate to illness and death cases from microbial  
10 infections and bladder cancer. They found that whilst there was a significant change in estimated WTP  
11 values across time when values were elicited using CV, there was no such significant change for the  
12 same values elicited using CE. A similar methodological comparison was undertaken by [Brouwer and](#)  
13 [Logar \(2014\)](#), who survey the same sample of people in Switzerland at a 6-month interval using both  
14 CV and CE. Their study relates to WTP for upgrading of waste water treatment plants in Switzerland to  
15 remove micro-pollutants. Some 20% of CE responses and some 30% of CV responses showed no change  
16 in preferences over the 6 month interval. There was no significant difference, however, between WTP  
17 estimates over time for the sample as a whole, and no significant difference between CV and CE in this  
18 respect. Beyond environmental applications, [Ryan et al. \(2006\)](#) and [Skjoldborg, Lauridsen and Junker](#)  
19 [\(2009\)](#) provide test-retest analysis for preferences regarding health care.

20 Unfortunately, within-survey tests of preference stability within a DCE setting may also reflect fatigue  
21 or learning effects. As people progress through a series of choice tasks, they may learn more what they  
22 like or do not like, so that they become more precise in their preferences in the sense that the distribution  
23 of their preference type becomes narrower as experience in choosing increases ([Czajkowski, Hanley and](#)  
24 [LaRiviere, 2014](#)). As people repeat choices, they may also find that a choice task becomes simpler; or  
25 else they may become bored and start using heuristics more frequently ([Swait and Adamowicz, 2001](#)).  
26 Any of these effects could show up as a change in the estimated values implied by choices, whereas in  
27 fact there has been no shift in underlying “true” preferences. Such fatigue or learning effects could also  
28 show up in the random component of utility ([Czajkowski, Giergiczny and Greene, 2014](#)). A review of  
29 multiple such “ordering effects” as well as their empirical testing can be found in [Day et al. \(2012\)](#).  
30 There have been a number of papers which also demonstrate a related “time to think” effect on WTP for  
31 changes in an environmental good ([Whittington et al., 1992](#); [MacMillan, Hanley and Lienhoop, 2006](#)).

32 In this paper, we conduct both within survey and between two different moments of time tests of the  
33 stability of choices and estimated distributions of WTP. These tests are based on observations of the  
34 same individuals. The within-survey test considers responses to the first 12 and then second 12 choice  
35 questions in a survey on options for forest management. The between moments of time test compares

1 these choices with responses from a similar (and for one subsample identical), 12-question DCE carried  
2 out six months later. This design provides a contribution to most of the test-retest literature, which as  
3 noted above has focussed on between moments of time tests only. In a within-survey experiment,  
4 individuals may become more precise in stating their preferences, or may discover them as they gain  
5 experience in choosing between different bundles of a good. This can confuse any signal about  
6 preference stability. This perspective stands in contrast to between moments of time tests, but here the  
7 researcher must confront a different set of problems, such as whether an individuals' socio-economic  
8 conditions changed, or where they may learn more about the good (rather than learning their preferences)  
9 over the interval. By investigating both issues jointly, our study provides an insight into the extent of  
10 the changes which may result from each of them. Although the two phenomena may be caused by  
11 different behavioural and economic effects, researchers' interest is basically the same in both cases –  
12 whether the hypothesis of stable welfare measures can be rejected (within a sequence of choices in a  
13 survey or between two different moments of time).

14

## 15 **2. Study Design and Data**

### 16 **2.1. The setting of the study – the Białowieża Forest**

17 The Białowieża Forest in Poland is an ancient woodland straddling at the border between Belarus and  
18 Poland, located in north east-central Poland on the border with Belarus. It is one of the last and largest  
19 remaining parts of the immense primeval forest which once spread across the European Plain. The  
20 Białowieża Forest is one of the most recognized and ecologically valuable forests in Poland  
21 ([Czajkowski, Buszko-Briggs and Hanley, 2009](#)). Despite some visible signs of human activity, it is still  
22 commonly considered the last natural lowland forest in temperate Europe. It is especially regarded for  
23 its natural dynamics as well as its species richness, and its ecological structures and functions  
24 ([Wesołowski, 2005](#)).

25 From the early 1990s biologists, environmentalists and various NGOs have been trying to convince  
26 decision makers to enlarge BNP over the entire territory of the Białowieża Forest; so far, unsuccessfully.  
27 One of the aims of conducting our study was to provide arguments in public discussions regarding the  
28 enlargement of the Białowieża National Park and possible changes in the forest management. In  
29 addition, our survey was constructed in a way which enabled testing preference stability, which is the  
30 main purpose of this paper.

31 A few one-to-one in-depth interviews were conducted by the research team members to fine-tune the  
32 survey instruments (structure, wording, visual materials – maps and photos). After consultations with

1 biologists<sup>1</sup> working in the Białowieża Forests two possible management levels for the forests outside  
2 BNP and the reserve have been considered that is:

- 3 1) maintain the current management typical for managed forest or
- 4 2) enlarge the passive protection zone, that is to allow for rewilding<sup>2</sup> of the managed part of the  
5 Białowieża Forest.

6 It was explained that these options would result in low or high level of forest naturalness respectively.

7 The differences between managed forests (low level of naturalness) and close to natural forests (high  
8 level of naturalness) were explained to the respondents with the use of photographs, drawings and  
9 written descriptions presented in Figure 1.

10 The Białowieża Forest can be divided into three relatively homogenous parts which differ in naturalness  
11 levels. The short characteristic and possible changes in the management in each part were explained to  
12 the respondents.

- 13 1) **The Białowieża National Park and the nature reserve Natural Forests of the Białowieża**  
14 **Primeval Forest** – the core of the Białowieża Forests is 10,500 ha protected in the Białowieża  
15 National Park (16% of the Polish part of the forest). In addition, 12,000 ha (19% of the Polish  
16 part of the forest) of natural forests are protected in the reserve outside the Białowieża National  
17 Park. Almost the entire area of the Białowieża National Park and forests in the reserve are  
18 protected forests under non-intervention regime. The Białowieża National Park and the reserve  
19 comprise the best preserved part of the Białowieża Forest. Some species of lichens, fungi,  
20 insects which depend on high volumes of dead wood are present only here ([Wesołowski, 2005](#)).  
21 In the questionnaire we referred to this part as having ‘High naturalness level’.
- 22 2) **Second growth forests** – approximately 6,000 ha (15% of the Polish part of the Białowieża  
23 Forest) are so called remains after large scale clear-cuts made in 1920’s by the British European  
24 Century Timber Corporation known as Centura which made no attempt to renew these clear-  
25 cuts<sup>3</sup> ([Directorate General of the State Forests, 2011](#)). These plots have been naturally  
26 regenerated and in most of these places no significant human intervention have taken place over  
27 the last 90 years. However, this may change at any time as these stands are considered by the  
28 State Forests as commercial forests. It was explained in the questionnaire that the second growth

---

<sup>1</sup> We are very grateful to prof. Bogdan Jaroszewicz a director of Białowieża Geobotanical Station for his comments on an early draft of our questionnaire.

<sup>2</sup> By rewilding we mean the whole process of returning ecosystems to a state of ecological health and dynamic balance, making them self-sustaining, without the need for ongoing human management ([Navarro and Pereira, 2012](#)).

<sup>3</sup> Centura was obliged to use 10-hectare clear-cuts and retain seed trees. In fact, the clear-cuts often exceeded a hundred hectares each and no attempt was made at their renewal ([Directorate General of the State Forests, 2011](#)).

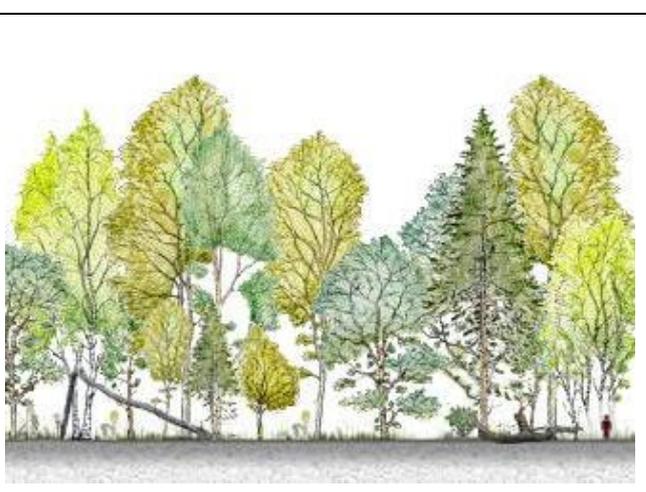
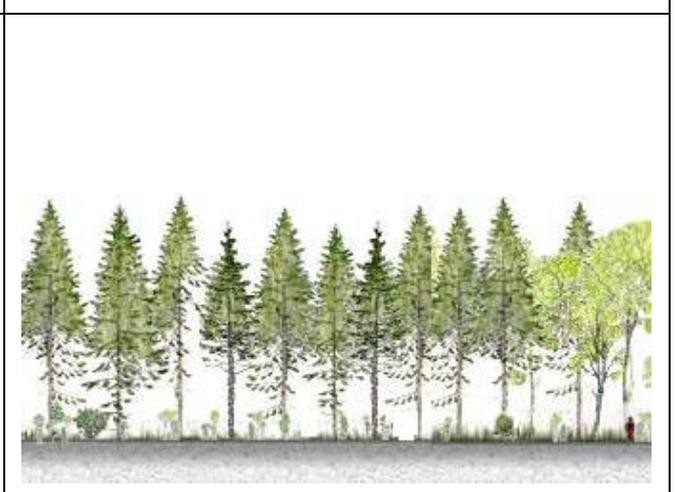
1 forests are under dynamic changes now as the climax status has not been reached yet and that  
2 their future status will depend on the chosen management. If passive protection is extended this  
3 forest will have 'High naturalness level' as the forests in the Białowieża National Park and the  
4 reserve outside the Park currently. After consultations with forest botanists we informed the  
5 respondents that this status will be reached in about 150 years. However, if these forests are  
6 turned into standard managed forests they will have 'Low naturalness level' as the managed  
7 forests in the Białowieża Forest.

8 **Managed forests** – approximately 33,000 ha (50% of the Polish part of the Białowieża Forest)  
9 are standard managed forests. In the questionnaire we referred to this part as having 'Low  
10 naturalness level'. After consultations with forest botanists we informed the respondents that it  
11 is possible to rewild these forests that is if passive protection is extended over the managed  
12 forests they will have 'High naturalness level' as the forests in the Białowieża National Park  
13 and the reserve outside the Park currently. The respondents that this status will be reached in  
14 about 250 years

15 Finally, the Białowieża Forest is a popular destination for both Polish and foreign visitors. For several  
16 years, an increasing number of visitors to the Białowieża Forest has been noted, especially in summer  
17 and autumn during public holidays. Currently around 100,000 people visit the Białowieża Forest each  
18 year. During the public holidays in May and in the summer the forest is visited much more frequently,  
19 with up to 10,000 people on exceptionally busy days. In the future, these numbers may increase. On the  
20 one hand, a large number of visitors may make the achievement of a higher level of forest naturalness  
21 difficult, and may even reduce the current level of ecological quality. For many years, environmentalists  
22 and various non-government organizations have been trying to convince decision makers to extend  
23 National Park designation to the entire Białowieża Forest (so far unsuccessfully, mainly due to local  
24 opposition), and to restrict visitor numbers.

25

1 Figure 1. The differences between natural forests and managed forests

Natural forest (high level of naturalness)	Managed forest (low level of naturalness)
	
	
<ol style="list-style-type: none"> <li>1) All trees are left in the forest until they die and decay.</li> <li>2) Trees regenerate naturally by self-sown seeds</li> <li>3) Forest has a multi-age structure</li> <li>4) There are usually many tree species</li> <li>5) High volume of dead wood (100 m<sup>3</sup>/ha or more).</li> <li>6) There is a greater diversity of species of plants, animals and fungi. Many rare species can live only in the forests with a large quantity of old rotting trees.</li> </ol>	<ol style="list-style-type: none"> <li>1) After attaining a certain age, the forest is logged. Old trees are met rarely.</li> <li>2) Trees regenerate artificially by direct seeding or <i>planting</i>.</li> <li>3) Forest has an even-age structure</li> <li>4) Usually scots pine or spruce monoculture</li> <li>5) There is a small volume of dead wood (less than 10 m<sup>3</sup>/ha).</li> <li>6) There is a much smaller diversity of species of plants, animals and fungi. Rare species do not have good conditions to live here.</li> </ol>
<p>Close to natural forests cover about 0.6% of all Polish forests</p>	<p>About 99% of forests in Poland are managed forests</p>

1 **2.2. Experimental design and data**

2 The DCE comprised five attributes, of which one – forest naturalness in the National Park and the nature  
 3 reserve – could take only one level and was presented for completeness. We decided to include this  
 4 attribute in the choice sets to show respondents the entire picture of the current and future forest  
 5 management of the Białowieża forest. It is very unlikely that in future the area currently under passive  
 6 protection could be subject to less stringent protection, and thus the level of this attribute was held  
 7 constant across all choice tasks and alternatives. The next two attributes referred to future management  
 8 programs (and the resulting naturalness level of that part of the Białowieża Forest) for the areas which  
 9 are currently under commercial management, and the second-growth sections of the Białowieża Forest.  
 10 These management programs (attribute levels) were “commercial use” or “passive protection”, resulting  
 11 in a low and high level of naturalness in the future, respectively. The fourth attribute was the maximum  
 12 visitor numbers allowed for the entire Białowieża Forest, i.e. limiting the tourism pressure on the forest.  
 13 Lastly, the design included the cost attribute. The payment vehicle was described as the increase in  
 14 income taxes paid annually by each household in Poland. The attribute and attribute levels are  
 15 summarized in Table 1.

16

17 Table 1. Attributes and levels used in the DCE study

Attributes	Levels	Variables
National Park and Nature Reserve (35% of the Białowieża forest)	High level of naturalness	<i>SQ</i>
Commercial forests (50% of the Białowieża forest)	Low level of naturalness	<i>SQ</i>
	High level of naturalness in 250 years	<i>COM</i>
Second-growth forests (15% of the Białowieża forest)	Low level of naturalness	<i>SQ</i>
	High level of naturalness in 150 years	<i>SGR</i>
Number of visitors	No restrictions	<i>SQ</i>
	5,000 visitors per day	<i>VIS<sub>1</sub></i>
	7,500 visitors per day	<i>VIS<sub>2</sub></i>
Annual cost per household	0, 25, 50, 75, 100 PLN	<i>COST</i>

18

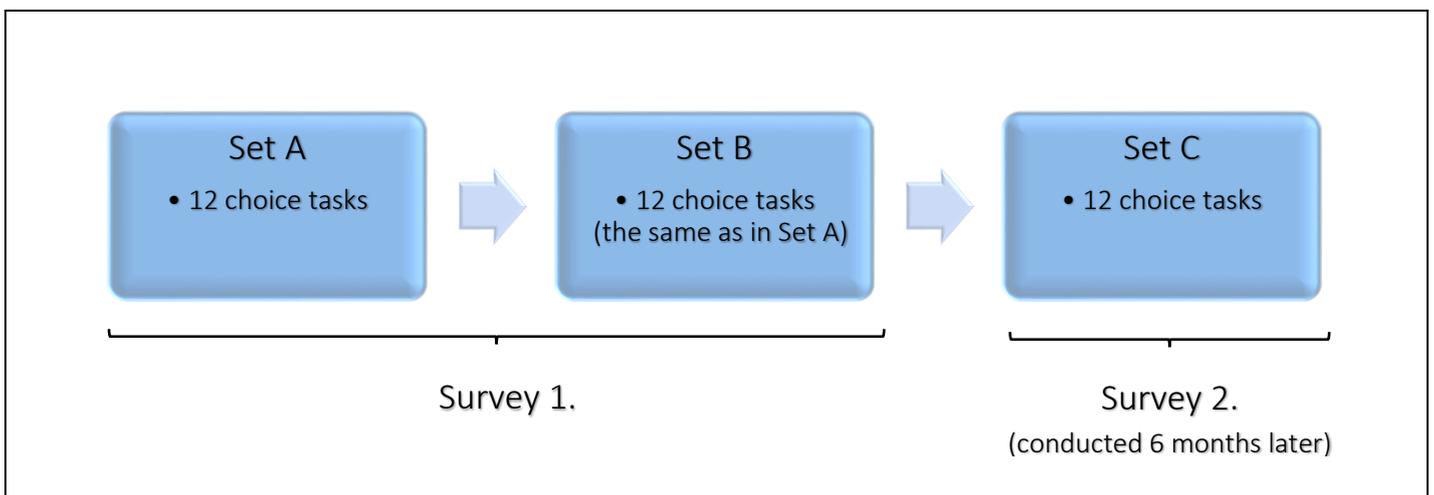
19 The contingent scenario description, attributes and their levels were developed in the process of  
 20 extensive qualitative testing to make sure the survey is understandable and the program credible. The  
 21 main survey was preceded with a pilot study, aimed at making sure the instrument worked well and for  
 22 collecting priors for the main study design.

23 The data used in this study was collected through a national online survey of the Polish population using  
 24 two surveys carried out by a professional polling agency. The first survey was conducted in December  
 25 2011. Of the 3,016 respondents who were invited to take part in the web panel survey 1,459 responded,

1 and 1,302 of them provided usable responses. The final sample was quota-controlled for sex, age, region  
2 and city size. In June 2012, 998 of the participants who successfully completed the first survey were  
3 invited to participate in a follow-up survey, and 789 of them successfully completed the survey. Since  
4 we were able to identify specific respondents within the data set, this provides an opportunity for us to  
5 test how their preferences changed over the course of the 6-month period (December to June) while  
6 controlling for the panel structure of the data.<sup>4</sup>

7 The first survey (December 2011) included 24 choice tasks, which were in fact the same 12 choice tasks  
8 repeated. The first 12 choices are hereafter referred to as set A and the second 12 choices as set B. This  
9 design provides an opportunity for us to undertake a within-sample test of preference stability, and  
10 indeed to compare on a choice-task by choice-task basis. The second survey (June 2012) included 12  
11 choice tasks (set C). The design of the study is illustrated in Figure 2. Comparing set A with set C for  
12 each person is thus a “between two different moments of time” comparison of preference stability.<sup>5</sup>

13



14 Figure 2. Illustration of the design of the preference stability study

15

16 The experimental design of our study utilized both an optimal-in-difference ([Street, Burgess and](#)  
17 [Louviere, 2005](#); [Street and Burgess, 2007](#)) and an efficient design ([Ferrini and Scarpa, 2007](#); [Scarpa and](#)  
18 [Rose, 2008](#)) approach.<sup>6</sup> The designs were generated and applied separately for half of the participants

<sup>4</sup> The questionnaire and dataset used in this study are made available online at [czaj.org](http://czaj.org). The models were estimated using custom code developed in Matlab which is made available from [github.com/czaj/DCE](https://github.com/czaj/DCE) under Creative Commons BY 4.0 license.

<sup>5</sup> Even though the same choice tasks were repeated in the same survey, respondents of the pretesting phase and the main survey appeared unaware of this, or they were not disturbed and did not comment on this.

<sup>6</sup> The efficient design was generated using priors obtained from an MNL model estimated on the results of a pilot study conducted on a sample of 100 respondents. In order to account for uncertainty associated with our priors we

1 of the first survey. In the second survey, half of the participants received the same optimal-in-difference  
 2 design choice tasks, while the other half received updated efficient design choice tasks. For ¼ of our  
 3 respondents ( $n = 193$ ) we are thus additionally able to test if their *choices* changed between the first and  
 4 the second survey, in addition to testing if their *implied preferences* changed.

5 Each choice task consisted of 3 alternatives, one of which was the Status Quo (*SQ*, meaning no changes  
 6 in current management program). In order to control for possible choice task- and position-specific  
 7 ordering effects, each respondent was presented with a counterbalanced design in which (1) the order of  
 8 choice tasks and (2) the order of alternatives (including the status quo alternative) was randomized. We  
 9 have taken steps to ensure that each choice task and each alternative was presented in every position in  
 10 the sequence a comparable number of times. An example of a choice card is presented in Figure 3.

	<b>Program A</b>	<b>Program B</b>	<b>Program C</b>
	Continuation of current management program	Changes in current management program	Changes in current management program
<b>National Park and Natural Reserves (35% of the Białowieża forest)</b>	High level of naturalness	High level of naturalness	High level of naturalness
<b>Commercial forests (50% of the Białowieża forest)</b>	Low level of naturalness	Low level of naturalness	High level of naturalness
<b>Second-growth forests (15% of the Białowieża forest)</b>	Low level of naturalness	High level of naturalness	High level of naturalness
<b>Number of visitors (per day)</b>	No limit	No limit	5,000
<b>Cost for your household (per year)</b>	0 PLN	50 PLN	100 PLN
<b>Your preferred program:</b>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>

11 Figure 3. An example of a choice card presented to the respondents.

12

### 13 3. Non-parametric analysis

14 In what follows, non-parametric tests of respondents' preference stability between sets A, B and C are  
 15 presented. Because the choice tasks in sets A and B for each respondent were exactly the same, we are  
 16 able to test if any respondents changed their answers and if so, then how and when. In addition, for the  
 17 193 respondents who were presented with the same choice tasks in sets A, B and C we are able to test

---

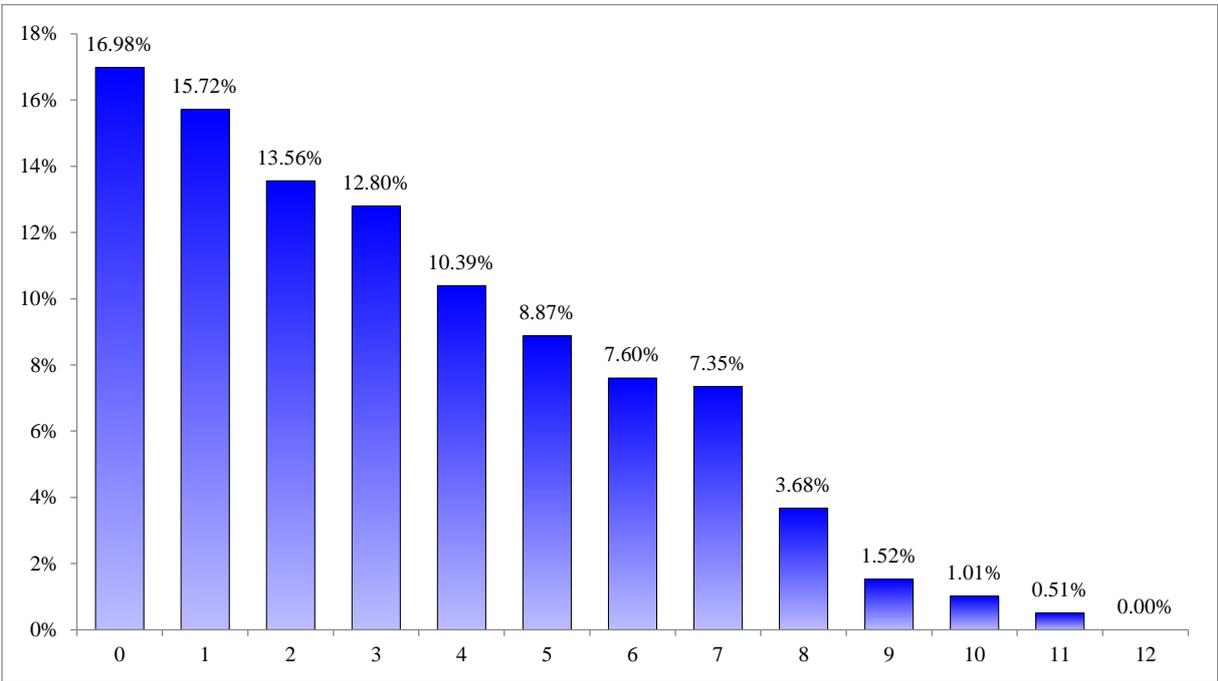
employed Bayesian approach ([Bliemer, Rose and Hess, 2008](#)) in which we assumed all priors to be normally distributed with means estimated from the MNL model and standard deviations equal to 0.25 of each parameter mean.

1 if choices differed in any pair of sets. Additionally we checked whether a pattern of choosing the *SQ*  
2 alternative differ across sets.

3

4 **3.1. Within survey comparison**

5 Figure 4 shows the distribution of the number of respondents who changed their choices between the  
6 sets A and B a particular number of times. Overall, only a minority of respondents (16.98%) made the  
7 same choice in every one of the repeated 12 choice tasks (and nearly 70 % of these always chose the *SQ*  
8 alternative). 59% of respondents changed decisions no more than 3 times, while 6.72% of respondents  
9 changed their decisions 8 times or more.<sup>7</sup>



10

11 Figure 4. The distribution of respondents who changed their decisions a particular number of times in  
12 the series of 12 choice tasks repeated twice (sets A and B)

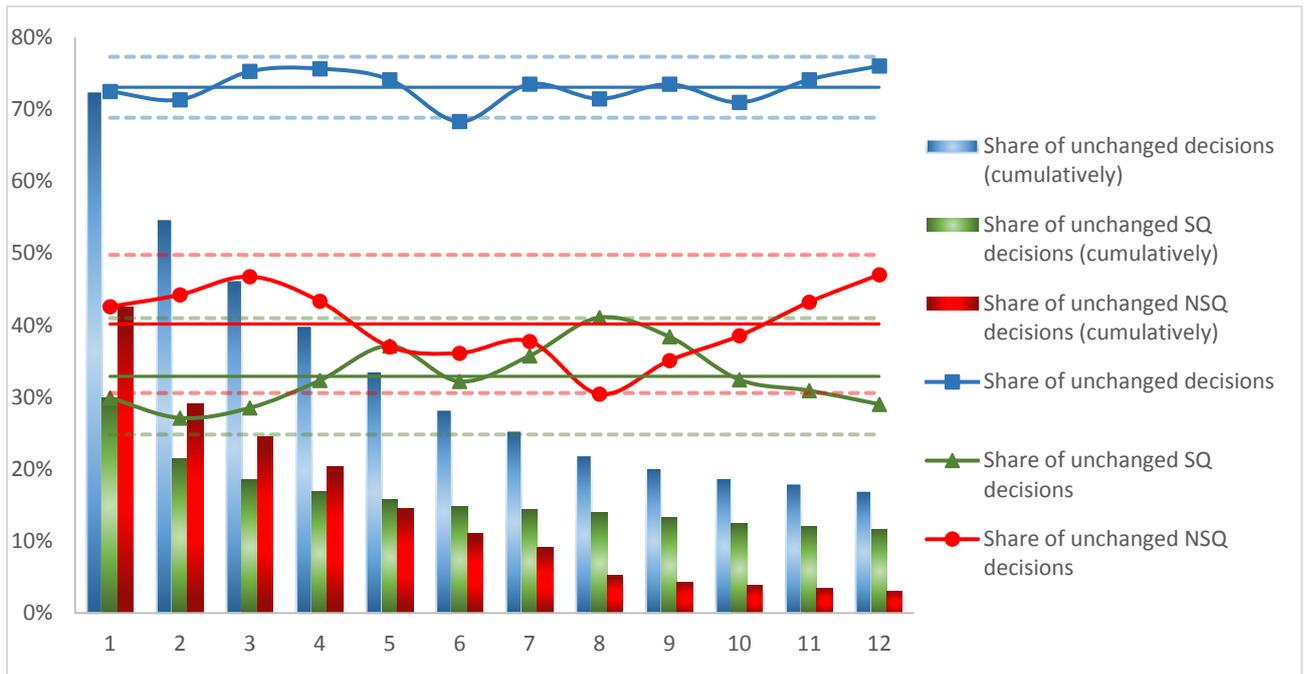
13

14 Additional insights are provided by an analysis of the dynamics of decision changes illustrated in Figure  
15 4. In addition to the share of unchanged decisions in each choice task, we present their sample means  
16 and their confidence intervals. No dynamics appear significant – the changes of decisions appear  
17 relatively uniformly distributed across choice tasks. This can be an indication of a preference stability  
18 when combined with the relatively high degree of variability in respondents' choices (a low scale

---

<sup>7</sup> Note that by making decisions completely randomly a respondent would, on average, make 8 different choices over the course of 12 choice tasks.

1 parameter, i.e. a high utility function error term variance in relation to its deterministic component). If  
 2 respondents' preferences change between choice tasks (e.g., due to preference learning), it does not  
 3 appear to change their decisions in a significant way. In addition, we find no difference between  
 4 respondents who chose the *SQ* and non-*SQ* alternatives in the sense that both groups seemed equally  
 5 likely to change their decisions at every choice task. However, respondents who chose non-status quo  
 6 alternatives were less likely to be fully consistent in all 12 choice tasks – it seems easier for the  
 7 respondents' who chose the *SQ* alternative in every choice to be consistent, as indicated by the shares  
 8 of cumulatively unchanged decisions presented as bar plots in Figure 5.



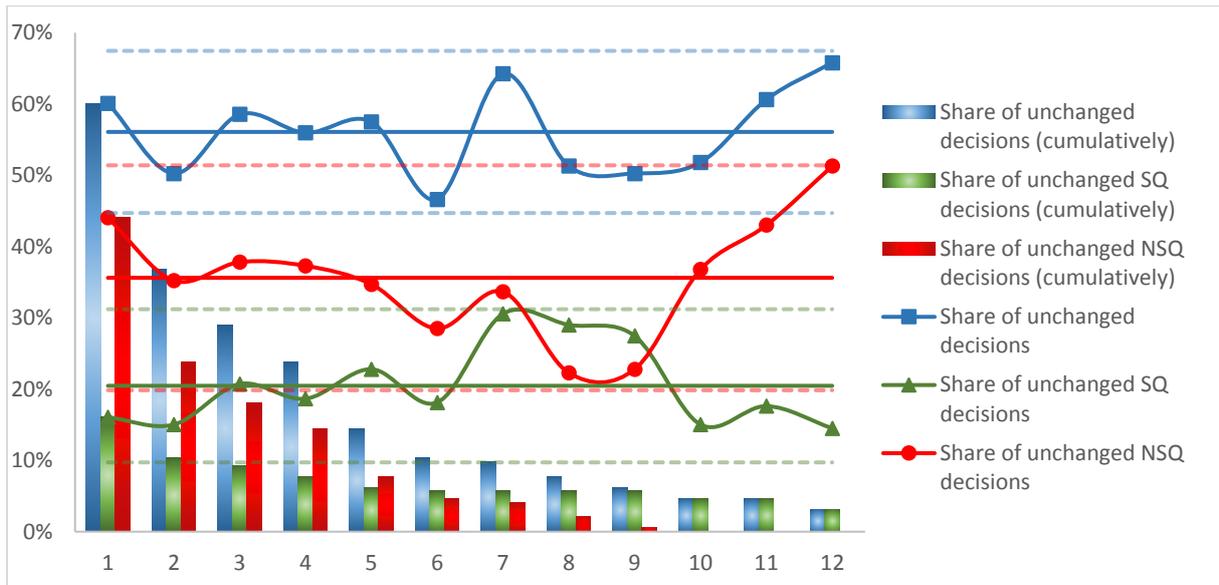
9  
 10 Figure 5. Dynamics of the number of decision changes in different choice tasks (set A vs. B;  $n = 789$ )

11  
 12 In Annex A we additionally included analogous plot for decision changes dynamics between sets A and  
 13 B for respondents who did not participate in second survey. The plots are very similar, there seems to  
 14 be no significant differences between them. To formally confirm this, we applied the Kruskal-Wallis  
 15 test of the differences in shares of unchanged decisions (overall, *SQ* and *NSQ*) for each choice task. The  
 16 results show no apparent patterns, the only significant differences were observed for the fifth choice task  
 17 for *SQ* shares and the eighth choice task for the *NSQ* shares. We therefore conclude that there are no  
 18 major differences between these samples with regard to preference dynamics and there is no apparent  
 19 self-selection bias, which could influence our results.

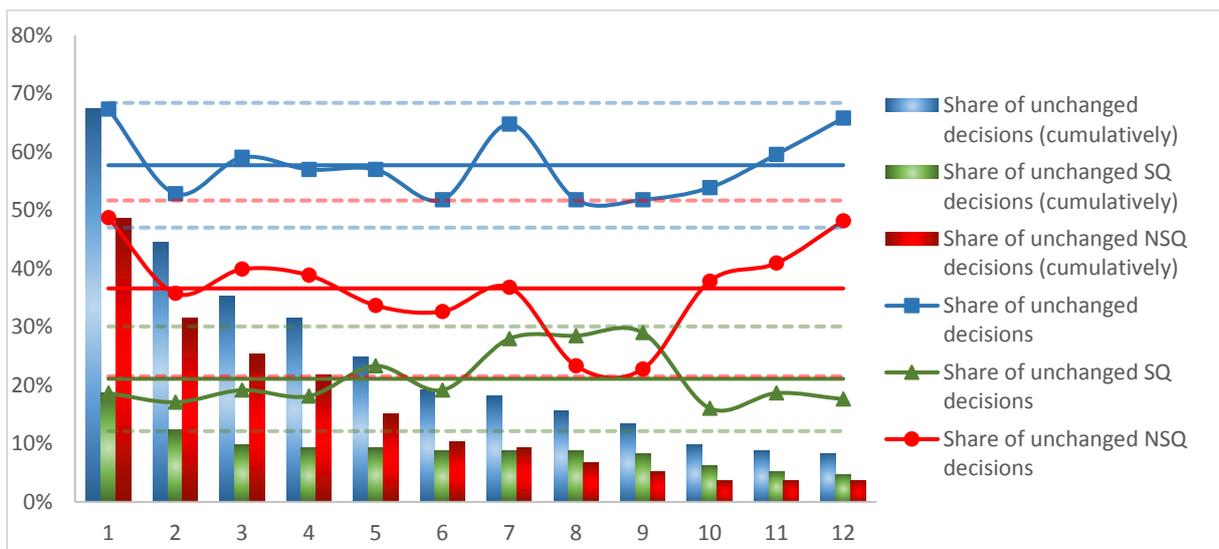
20

1 **3.2. Between two different moments of time differences**

2 Figures 6 and 7 present the dynamics of decision changes between surveys (i.e., between sets A and C,  
 3 and between sets B and C, respectively). Overall, there is more variation in these cases, but this is partly  
 4 a result of a smaller sample, since in this case there were only 193 respondents who were presented with  
 5 the same choices in both surveys. Like in the within survey case, the dynamics of respondents' choices  
 6 do not reveal any significant patterns. Interestingly, the only respondents who answered consistently  
 7 between the sets A and C were those who always chose the *SQ* alternative.



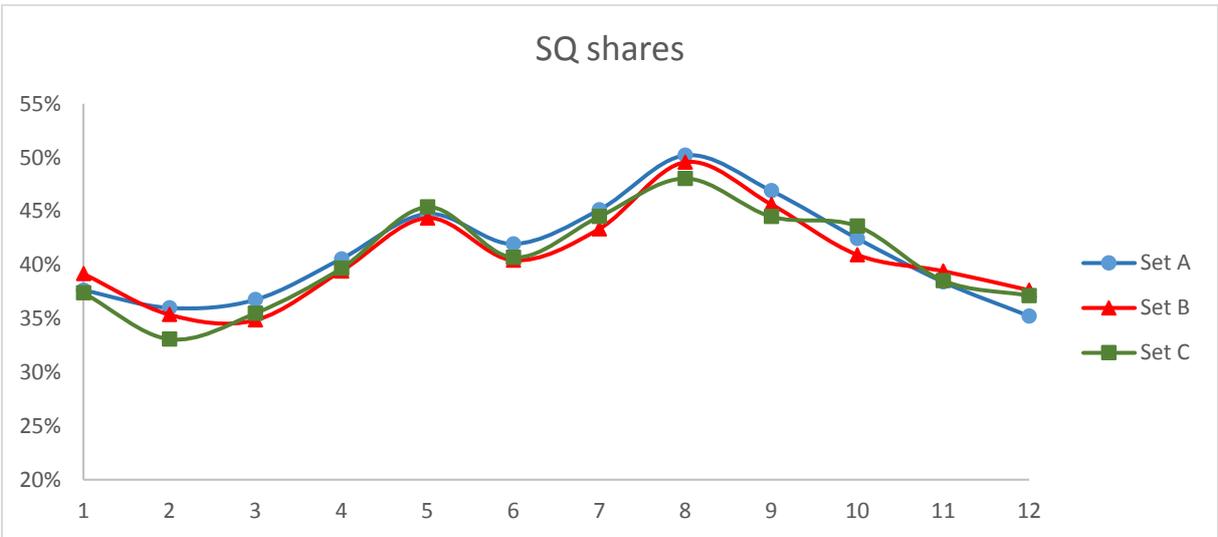
8  
 9 Figure 6. Dynamics of the number of decision changes in different choice tasks (set A vs. C;  $n = 193$ )



11  
 12 Figure 7. Dynamics of the number of decision changes in different choice tasks (set B vs. C;  $n = 193$ )

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14

Finally, we analyze the share of the SQ alternative choices. The respondents who had selected the SQ alternative in every choice task constitute 14.70%, 19.52% and 15.46% of the sample for sets A, B and C, respectively. The SQ alternative could have been selected by some respondents in some choice tasks only. Figure 8 presents the shares of the SQ alternative choices in each choice task. As argued by [Day et al. \(2012\)](#), the occurrence of fatigue can be characterized by, among others, an increase in the probability of selecting the SQ alternative. In our case, the share of respondents choosing the SQ alternative in each choice task is very similar. Despite some increase between the third and the eighth choice task, and subsequent decrease in the probability, there is no significant difference in the share of SQ choice between set A and B. We interpret it as an indication of no apparent fatigue effect of this type, although it should be noted that in set B share of respondents who have chosen only SQ increased when compared with set A.



15  
16 Figure 8. Shares of respondents who've chosen status quo alternative. (sets A, B and C; n = 789)

17  
18

#### 4. Econometric analysis

19 This section presents the results of parametric analysis of respondents' choices based on the random  
20 utility framework ([McFadden, 1974](#)). We test if respondents' marginal WTP distributions were stable  
21 between groups of choice sets A and B presented in survey 1, and between these choice sets and choice

1 set C presented in survey 2 six months later. We focus on stability of welfare measures instead of  
 2 stability of preferences, as it is usually done, for two reasons. First, WTP is the main result of valuation  
 3 studies, the reason most of these studies are conducted, and in the focus of policymakers who wish to  
 4 undertake a cost-benefit analysis. Second, WTP measures are scale-free and therefore can be directly  
 5 compared. For the comparison of preference parameters [Swait and Louviere \(1993\)](#) or [Czajkowski,  
 6 Hanley and LaRiviere \(2016\)](#) approach should be applied to take the possible differences in scale  
 7 between different sets into account. Having 3 sets and using this procedure for mixed logit is  
 8 computationally cumbersome.<sup>8</sup>

9 To conduct our analysis we used the mixed logit model ([MXL, Revelt and Train, 1998](#)) which was  
 10 estimated in the so called WTP-space ([Train and Weeks, 2005](#)), which means that utility of  $n$ -th  
 11 individual from choosing  $i$ -th alternative in  $j$ -th choice task is reparametrized in following way:

$$12 \quad U_{ijn} = \mathbf{X}_{ijn} \mathbf{b}_n - \alpha_n c_{ijn} + \varepsilon_{ijn} = \alpha_n \left( \mathbf{X}_{ijn} \mathbf{b}_n / \alpha_n - c_{ijn} \right) + \varepsilon_{ijn} = \alpha_n \left( \mathbf{X}_{ijn} \boldsymbol{\beta}_n - c_{ijn} \right) + \varepsilon_{ijn}$$

13 where,  $c_{ijn}$  denotes cost of given alternative,  $\mathbf{X}_{ijn}$  is a vector of other attributes and  $\varepsilon_{ijn}$  is a random  
 14 term. In this framework, the parameters for all attributes ( $\boldsymbol{\beta}_n$ ) but the cost can directly be interpreted  
 15 as WTP (in PLN per year)<sup>9</sup>. The advantage of this approach is that no additional calculations are needed  
 16 to estimate some characteristics of WTP distribution, e.g., means of normally distributed random  
 17 parameters can simply be interpreted as mean WTP. In this model we assumed that all parameters are  
 18 random, assuming that WTP for all attributes follow normal distribution, while the negative cost  
 19 parameter ( $\alpha_n$ ) is log-normally distributed.<sup>10</sup> The cost parameters are necessarily normalized to one in  
 20 WTP-space models; instead we present their preference-space counterparts (coefficients of the  
 21 underlying normal distribution).

22 In order to analyse stability of WTP we interacted every attribute with binary variables for sets A, B and  
 23 C, which lead to 18 random parameters in the model. Evaluating the statistical significance of the  
 24 differences of the estimates for different sets provides a convenient way of testing if respondents'  
 25 preferences as revealed in sets A, B and C are stable. In addition, we also allowed for full correlation  
 26 between all random parameters. This is important for two reasons: (i) it is very likely that WTP for the  
 27 same attributes in different sets are highly correlated (in the case of no changes in respondents' choices  
 28 these correlations should be equal to one), and (ii) the (possibly) heteroskedastic variance of error terms

---

<sup>8</sup> Particularly, because a meaningful likelihood ratio test comparisons of models estimated with the simulated maximum likelihood method requires using very many draws ([Andersen, 2014](#)).

<sup>9</sup> 1 PLN  $\approx$  0.23 EUR  $\approx$  0.25 USD

<sup>10</sup> We also tested the specification in which WTP followed log-normal distributions. It was not better fitted to the data, however (Voung test statistic = -0.4028).

1 could be different between sets A, B and C, but it is also very likely that its distribution is correlated  
2 between them.<sup>11</sup>

3 Table 2 presents the results of the MXL model. Estimation was performed using Matlab 8.1. The  
4 maximum likelihood function was simulated using 10,000 scrambled Sobol draws ([Czajkowski and](#)  
5 [Budziński, 2015](#)). Standard errors of coefficients associated with standard deviations of random  
6 parameters were simulated using Krinsky and Robb method with  $10^6$  draws ([Krinsky and Robb, 1986](#)).

7 The results presented in Table 2 show that, overall, respondents were generally in favor of extending  
8 passive protection to commercial (*COM*) and second-growth (*SGR*) areas of the Białowieża Forest, with  
9 average WTP in the range of 12.65 to 17.12 PLN for the former and 21.77 to 29.99 PLN for the latter.  
10 The distributions of these WTP were highly heterogeneous, as indicated by significant standard  
11 deviations of the random parameters. We can also observe a dislike regarding the *SQ* alternative  
12 (negative but heterogeneous WTP associated with the alternative specific constant for the no-change  
13 alternative). WTP for restricting the number of visitors ( $VIS_1$ ,  $VIS_2$ ) were highly heterogeneous as well,  
14 but mostly insignificant with different signs for different sets. As for the dynamics, we can observe that  
15 mean WTP for the *SQ*, *COM* and *SGR* is higher in set B than in set A, but then decreases to an even  
16 lower level in set C. WTP associated with  $VIS_1$  and  $VIS_2$  presents the opposite pattern, but due to the  
17 insignificance of these parameters it is not very informative. We revisit these conclusions later, in more  
18 formal way, using Wald test.

19 Inspecting the correlation matrix of random parameters (which are equivalent to WTP for all attributes  
20 except cost), estimated with the MXL model, presented in Table 3, provides additional insights. The  
21 greyed-out cells correspond to correlations between parameters of the same attributes in different sets.  
22 The correlations between WTP in set A and their counterparts in set B are very high, ranging from 92.3%  
23 to 99.6%. This means that WTP in set B are approximately linear transformations of their set A  
24 equivalents.<sup>12</sup> The correlations between sets A and C as well as sets B and C are much lower, but still  
25 positive and significantly different from zero.<sup>13</sup> WTPs are highly correlated within every set as well,  
26 especially WTP for  $VIS_1$  and  $VIS_2$ , which was expected.

27

---

<sup>11</sup> For example, respondents who were observed to have had a high variance of the error term in set A are also likely to have had a high variance in set C.

<sup>12</sup> Note that it is not a “scale effect”, as WTP’s are scale free.

<sup>13</sup> The correlations between the random *COST* parameters is lower than what we observe for WTP. However, because these parameters are a product of scale (set specific) and the parameter representing cost sensitivity, and it is not possible to disentangle the two effects, this observation offers limited insight and is difficult to interpret – we cannot identify the driver of the observed changes in heterogeneity structure.

1 Table 2. The results of the mixed logit model accounting for preference and scale differences between  
 2 sets A, B and C (in WTP-space).

	Mean WTP (s.e.)			Standard deviation of WTP (s.e.)		
	Set A (Survey 1)	Set B (Survey 1)	Set C (Survey 2)	Set A (Survey 1)	Set B (Survey 1)	Set C (Survey 2)
	<i>SQ</i>	-33.5013*** (1.7048)	-31.0284*** (1.9827)	-38.2758*** (2.4345)	93.4651*** (4.9807)	95.8237*** (4.8928)
<i>COM</i>	14.4530*** (1.3746)	17.1193*** (1.7130)	12.6514*** (1.6731)	30.5130*** (1.4613)	45.2908*** (2.2754)	38.1746*** (2.0322)
<i>SGR</i>	28.5644*** (1.4058)	29.9556*** (1.8671)	21.7744*** (1.7015)	35.5747*** (1.5599)	49.9512*** (2.1923)	36.6213*** (1.6475)
<i>VIS1</i>	0.7744 (1.5689)	-2.9075 (1.9450)	0.4340 (1.9776)	19.9192*** (2.0502)	25.5335*** (2.0286)	19.3573*** (1.6864)
<i>VIS2</i>	-3.0014* (1.7697)	-2.5435 (2.1626)	-2.6190 (2.4219)	41.1142*** (2.1927)	47.2967*** (2.0489)	40.4539*** (2.1286)
<i>COST</i> (pref. space equivalent)	-2.0903*** (0.0470)	-1.7978*** (0.0595)	-1.9872*** (0.0520)	1.1387*** (0.0474)	1.4160*** (0.0693)	1.1986*** (0.0647)
Log-likelihood (constants only)	-30,833.5642					
Log-likelihood	-18,496.7692					
McFadden R <sup>2</sup>	0.4001					
Ben-Akiva R <sup>2</sup>	0.5449					
AIC/n	1.3158					
n (observations)	28,404					
k (parameters)	189					

3 Note: \*\*\*, \*\*, \* represent statistical significance at the 1%, 5% and 10% levels, respectively.

4

5 Table 4 presents the results of the Wald tests conducted for different hypotheses of parameter equality.  
 6 It consists of 4 panels. In the first three, we tested whether WTP means, variances and means and  
 7 variances jointly are equal between every pair of sets, and all sets jointly for each of attributes, and  
 8 additionally – if they are equal for all the attributes jointly. The fourth panel provides the results for joint  
 9 equality of WTP means and variances and their respective correlations between sets being equal to 1.

10 Formally, we could say that the distributions of WTP are stable if we were not able to reject the null  
 11 hypothesis of the equality of random parameters means and variances for each of attribute between all  
 12 sets ( $A = B = C$ ) and if the correlations of random parameters were perfectly correlated between the 3  
 13 sets. The Wald test statistic for such a hypothesis is 582.46 which leads to its rejection (d.f. = 30). A less  
 14 restrictive definition of stability involves the equality of means and variances of WTP.<sup>14</sup> In this case, the  
 15 Wald statistic is 80.30 which also indicates the null hypothesis of equality can be rejected (d.f. = 20). In  
 16 conclusion, we cannot say that the distributions of welfare measures are stable between all sets.

<sup>14</sup> This hypothesis involves the equality of population-level characteristics of the WTP distributions, but does not require individual-specific WTP to be equal between sets.

1 Table 3. Correlation matrix of random parameters implied by the MXL model

		Set A (Survey 1)					Set B (Survey 1)					Set C (Survey 2)							
		<i>SQ</i>	<i>COM</i>	<i>SGR</i>	<i>VIS<sub>1</sub></i>	<i>VIS<sub>2</sub></i>	<i>COST</i>	<i>SQ</i>	<i>COM</i>	<i>SGR</i>	<i>VIS<sub>1</sub></i>	<i>VIS<sub>2</sub></i>	<i>COST</i>	<i>SQ</i>	<i>COM</i>	<i>SGR</i>	<i>VIS<sub>1</sub></i>	<i>VIS<sub>2</sub></i>	<i>COST</i>
Set A (Survey 1)	<i>SQ</i>	1.0000	-0.5601	-0.5288	-0.8579	-0.8507	0.1352	<b>0.9963</b>	-0.5742	-0.6219	-0.8346	-0.8732	0.1091	<b>0.6463</b>	-0.3478	-0.3181	-0.4337	-0.4662	0.0380
	<i>COM</i>		1.0000	0.4898	0.1080	0.1874	-0.0150	-0.5613	<b>0.9804</b>	0.5620	0.1454	0.2027	-0.0201	-0.3048	<b>0.3569</b>	0.2738	0.0492	0.0950	0.2212
	<i>SGR</i>			1.0000	0.4734	0.4756	-0.1366	-0.5108	0.4883	<b>0.9798</b>	0.4223	0.5308	0.0044	-0.3195	0.3028	<b>0.3894</b>	0.1329	0.1092	0.0535
	<i>VIS<sub>1</sub></i>				1.0000	0.9075	-0.2050	-0.8441	0.1360	0.5252	<b>0.9231</b>	0.9434	-0.1722	-0.5584	0.2366	0.2693	<b>0.4328</b>	0.4504	-0.1456
	<i>VIS<sub>2</sub></i>					1.0000	-0.0979	-0.8255	0.2197	0.5549	0.9237	<b>0.9647</b>	0.0087	-0.5755	0.2636	0.2365	0.5360	<b>0.5920</b>	-0.0372
	<i>COST</i>						1.0000	0.1191	0.0636	-0.0534	-0.0723	-0.1833	<b>0.7854</b>	0.0966	-0.0285	0.2900	0.1109	-0.0448	<b>0.4147</b>
Set B (Survey 1)	<i>SQ</i>							1.0000	-0.5748	-0.6029	-0.8119	-0.8475	0.1009	<b>0.6513</b>	-0.3328	-0.3162	-0.4336	-0.4539	0.0593
	<i>COM</i>								1.0000	0.5792	0.1844	0.2164	0.0180	-0.2892	<b>0.3380</b>	0.2945	0.0680	0.0944	0.1928
	<i>SGR</i>									1.0000	0.4970	0.5926	0.0697	-0.3688	0.3316	<b>0.4221</b>	0.1982	0.1705	0.0782
	<i>VIS<sub>1</sub></i>										1.0000	0.9568	-0.0220	-0.5005	0.1442	0.2554	<b>0.4422</b>	0.4506	0.0299
	<i>VIS<sub>2</sub></i>											1.0000	-0.0680	-0.5654	0.2556	0.2470	0.4604	<b>0.5024</b>	-0.0043
	<i>COST</i>												1.0000	0.0074	0.0132	0.2125	0.0857	-0.0298	<b>0.3722</b>
Set C (Survey 2)	<i>SQ</i>													1.0000	-0.5714	-0.4637	-0.5963	-0.7142	0.1380
	<i>COM</i>														1.0000	0.3651	-0.0601	0.1903	-0.0569
	<i>SGR</i>															1.0000	0.3075	0.4493	0.1567
	<i>VIS<sub>1</sub></i>																1.0000	0.8506	0.0075
	<i>VIS<sub>2</sub></i>																	1.0000	-0.0583
	<i>COST</i>																		1.0000

2

Note: Coefficients in bold are significant at at least 5% level.

1 Table 4. Wald test statistics (degrees of freedom in brackets) for equality of different sets of parameters  
 2 in the MXL model

	Means (Panel I)				Variances (Panel II)			
	A = B	A = C	B = C	A = B = C	A = B	A = C	B = C	A = B = C
<i>SQ</i>	1.0138 (1)	2.7056* (1)	5.6417** (1)	5.6518* (2)	0.1118 (1)	0.5788 (1)	0.2236 (1)	0.5788 (2)
<i>COM</i>	2.1128 (1)	0.7258 (1)	3.7172* (1)	4.0337 (2)	26.2436*** (1)	8.3517*** (1)	5.2064** (1)	32.1968*** (2)
<i>SGR</i>	0.4459 (1)	9.7694*** (1)	10.6682*** (1)	12.7026*** (2)	29.3349*** (1)	0.1169 (1)	22.3449*** (1)	30.0529*** (2)
<i>VIS<sub>1</sub></i>	2.3144 (1)	0.0186 (1)	1.5134 (1)	2.5617 (2)	3.5469* (1)	0.2496 (1)	5.9240** (1)	6.0582** (2)
<i>VIS<sub>2</sub></i>	0.0313 (1)	0.0176 (1)	0.0006 (1)	0.0365 (2)	4.5015** (1)	0.1045 (1)	5.6861** (1)	6.7840** (2)
All WTP	6.9604 (5)	10.4590* (5)	18.0242*** (5)	22.1515** (10)	37.9044*** (5)	11.3648** (5)	26.9242*** (5)	52.5862*** (10)
	Means & Variances (Panel III)				Means & Variances & Correlations (Panel IV)			
	A = B	A = C	B = C	A = B = C	A = B	A = C	B = C	A = B = C
<i>SQ</i>	1.0173 (2)	3.5412 (2)	5.9852* (2)	6.2237 (4)	5.0550 (3)	136.6064*** (3)	135.3316*** (3)	171.0386*** (6)
<i>COM</i>	29.8657*** (2)	8.4317** (2)	9.6575*** (2)	37.2137*** (4)	32.3881*** (3)	141.8678*** (3)	150.9603*** (3)	161.9812*** (6)
<i>SGR</i>	29.8734*** (2)	9.8580*** (2)	35.0110*** (2)	42.7886*** (4)	31.5231*** (3)	117.3079*** (3)	118.7153*** (3)	129.6190*** (6)
<i>VIS<sub>1</sub></i>	4.2198 (2)	0.3403 (2)	6.0724** (2)	6.7376 (4)	10.0569** (3)	20.7889*** (3)	29.7390*** (3)	34.6797*** (6)
<i>VIS<sub>2</sub></i>	6.8311** (2)	0.1048 (2)	6.9301** (2)	9.1927* (4)	11.9285*** (3)	75.1860*** (3)	94.1143*** (3)	102.3977*** (6)
All WTP	43.9139*** (10)	30.1523*** (10)	43.4875*** (10)	80.2975*** (20)	57.1342*** (15)	463.9029*** (15)	482.8585*** (15)	582.4573*** (30)

Note: \*\*\*, \*\*, \* represent statistical significance at the 1%, 5% and 10% levels, respectively.

Further investigation of the results presented in Table 4 reveals some interesting facts. First of all, WTP for the *SQ* (alternative specific constant) seems to be the most stable between all sets. As can be seen in the third panel, we cannot reject the null hypothesis of equality of mean and variance between all 3 sets (Wald test statistic = 6.22; d.f. = 4). In addition, the fourth panel reveals that the correlation of WTP for the *SQ* between sets A and B is also not statistically different from 1. This does not hold between sets A and C or B and C – the correlations are significantly lower. This result is in line with what we observed in section 3, i.e., that *SQ* choices appeared more stable (as indicated by cumulative plots and the shares of the *SQ* alternative choices). It is contrast with the result of [Dekker, Koster and Brouwer \(2014\)](#) who found that WTP for the *SQ* is likely to be decreasing across choice tasks.

1 The Wald tests results indicate that the changes in mean WTP between sets A and B are not statistically  
2 significant. Nevertheless, there are significant differences between the variances of random parameters  
3 (see Panel II of Table 4). As illustrated by Table 2, standard deviations of WTP are higher in set B than  
4 in set A. As a result, although the means of WTP do not seem to change, their standard deviations  
5 increase. In addition, we note very high correlations between the observed WTP (above 92%). This  
6 indicates that even though the distributions of WTP are getting more disperse in set B than in set A,  
7 individual respondents' WTP generally stay close.

8 Mean WTP do not significantly change between sets A and C either, except for the second-growth areas  
9 (*SGR*). The hypothesis of equal variances cannot be rejected on 1% confidence level in this case, due to  
10 significant difference of variances of WTP for commercial forests (*COM*). The differences between sets  
11 B and C are even more evident – means and variances are significantly different irrespectively of  
12 whether they are tested separately or jointly. Overall, we cannot conclude that WTP remain stable  
13 between sets A, B and C, although the differences seem to be mainly driven by the changes in variances  
14 and imperfect correlations of individual-specific WTP.

## 15 **5. Conclusions**

16 The temporal stability of preferences is a key working hypothesis in cost-benefit analysis and the theory  
17 of value. It is thus something which is important to test. For environmental goods, several such tests are  
18 reported in the literature using stated preference methods. Tests are typically made using samples from  
19 the sample individuals or repeated, non-identical samples from the same population over a time interval  
20 such as one month or one year. If we observe that preferences or WTP has changed, then this can be  
21 attributed to several factors which are hard to separately identify. Over time, people can change the  
22 relative importance they attach to goods and services, possibly through learning more about substitutes  
23 and complements for this good and/or the characteristics of these goods. Incomes and other socio-  
24 economic attributes of individuals can also change over time, impacting on WTP. Between two different  
25 moments of time tests examine such changes.

26 But we can also look at how stated preferences and WTP evolve within an experiment if a DCE approach  
27 is taken, where people complete multiple choice tasks. Here, testing for preference stability is  
28 complicated by a different set of considerations. Apparent differences in estimated preferences could be  
29 due to people learning their true preference type with more precision as they repeat the choice tasks,  
30 learning how to choose better with repetition, or becoming bored and thus falling back onto heuristics  
31 whilst making choices. Unfortunately, the data collected in the present survey do not allow us to isolate  
32 any learning or fatigue effects in the within-sample tests; whilst there may be changes in respondents'

1 circumstances and knowledge over the 6 months which elapsed between surveys 1 and 2 which we  
2 cannot control for.<sup>15</sup>

3 Our contribution is to combine within- and between two different moments of time tests of preference  
4 stability for the same set of individuals. The within-survey test compares WTP between the first 12  
5 choices and the second set of 12 choices within the same survey. The between two different moments  
6 of time test involved re-surveying the same individuals six months later, and repeating the same first 12  
7 choice tasks. This allowed us to examine the dynamics of choices. Results show that respondents often  
8 changed their choices, but there is no obvious pattern in terms of *when* these changes occur. Deviations  
9 from previous choices were relatively uniformly distributed. Having said that, respondents who  
10 consistently choose the *SQ* alternative were less likely to change their choices than those who did not  
11 prefer the *SQ*.

12 The comparison of mean WTP revealed that this key variable is relatively stable, both within- and  
13 between two different moments of time. However, we found that the variances of WTP distributions did  
14 differ significantly, especially within survey. In addition, the analysis of correlations of individual-  
15 specific WTP showed that they were not perfectly correlated, although the correlation coefficients within  
16 survey were very high (above 92%) and between two different moments of time (observed from the  
17 same respondents 6 months later) remained positive (above 35%) and significant, indicating a high level  
18 of self-consistency. Overall, our findings led us to rejecting the hypothesis of perfect preference stability,  
19 which would here imply failing to reject the equivalence of the WTP distributions in choice sets A, B  
20 and C. This may be seen as being in contrast with the results of some earlier studies. However, this null  
21 hypothesis is more stringent than that usually applied (e.g., a simple comparison of mean WTP). We  
22 also take preference heterogeneity and correlations into account.

23 As a more general conclusion, despite observing statistically significant differences in WTP  
24 distributions and imperfect correlation of individual-specific WTP, our results indicate that the extent  
25 of changes is relatively low. That is, when compared with the level of uncertainty associated with  
26 welfare measures which results from sampling, survey administration, model specification or the  
27 valuation method used, the instability in willingness to pay values which are measured in our survey  
28 seem rather small. This is probably the most relevant aspect of our work for policy and management  
29 applications, and can be considered reassuring, since it implies that mean WTP remained relatively  
30 stable, both within and between surveys. This suggests some support for the continued use of cost-  
31 benefit analysis to help inform environmental management and policy decisions.

---

<sup>15</sup> Another potential source of differences was the time of the year when the two surveys were conducted. While the first one was administered in winter, the second elicited respondents' preferences 6 months later – in the summer. We acknowledge that this could also potentially influence the relative importance respondents attach to forests, and hence their WTP. This is a caveat future studies should avoid.

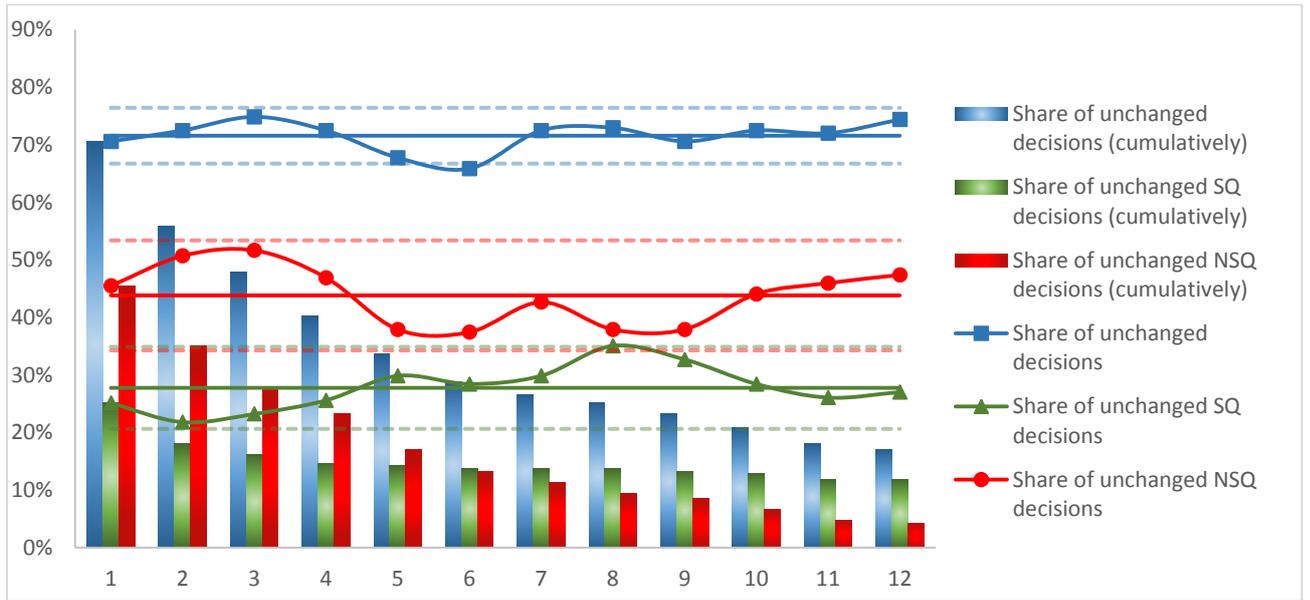
## References

- 1
- 2 Andersen, L. M., 2014. Obtaining Reliable Likelihood Ratio Tests from Simulated Likelihood  
3 Functions. *PLoS ONE*, 9(10):e106136.
- 4 Bliem, M., Getzner, M., and Rodiga-Laßnig, P., 2012. Temporal stability of individual preferences for  
5 river restoration in Austria using a choice experiment. *Journal of Environmental Management*,  
6 103(0):65-73.
- 7 Bliemer, M. C. J., Rose, J. M., and Hess, S., 2008. Approximation of Bayesian Efficiency in  
8 Experimental Choice Designs. *Journal of Choice Modelling*, 1(1):98-127.
- 9 Brouwer, R., 2006. Do stated preference methods stand the test of time? A test of the stability of  
10 contingent values and models for health risks when facing an extreme event. *Ecological  
11 Economics*, 60(2):399-406.
- 12 Brouwer, R., 2012. Constructed preference stability: a test–retest. *Journal of Environmental Economics  
13 and Policy*, 1(1):70-84.
- 14 Brouwer, R., and Logar, I., 2014. Do choice experiments produce more stable welfare measures than  
15 contingent valuation? A test-retest. paper presented at the Fifth World Congress of  
16 Environmental and Resource Economists, 28 June - 2 July 2014, Istanbul, Turkey.
- 17 Carson, R. T., Hanemann, W. M., Kopp, R. J., Jon, A. K., Mitchell, R. C., Presser, S., Rudd, P. A.,  
18 Smith, V. K., Conaway, M., and Martin, K., 1997. Temporal Reliability of Estimates from  
19 Contingent Valuation. *Land Economics*, 73(2):151-163.
- 20 Czajkowski, M., and Budziński, W. (2015). "An insight into the numerical simulation bias – a  
21 comparison of efficiency and performance of different types of quasi Monte Carlo simulation  
22 methods under a wide range of experimental conditions." In: *Environmental Choice Modelling  
23 Conference*, Copenhagen.
- 24 Czajkowski, M., Buszko-Briggs, M., and Hanley, N., 2009. Valuing Changes in Forest Biodiversity.  
25 *Ecological Economics*, 68(12):2910–2917.
- 26 Czajkowski, M., Giergiczny, M., and Greene, W. H., 2014. Learning and fatigue effects revisited.  
27 Investigating the effects of accounting for unobservable preference and scale heterogeneity.  
28 *Land Economics*, 90(2):323-350.
- 29 Czajkowski, M., Hanley, N., and LaRiviere, J., 2014. The Effects of Experience on Preferences: Theory  
30 and Empirics for Environmental Public Goods. *American Journal of Agricultural Economics*,  
31 97(1):333-351.
- 32 Czajkowski, M., Hanley, N., and LaRiviere, J., 2016. Controlling for the effects of information in a  
33 public goods discrete choice model. *Environmental and Resource Economics*, 63(3):523-544.
- 34 Day, B., Bateman, I. J., Carson, R. T., Dupont, D., Louviere, J. J., Morimoto, S., Scarpa, R., and Wang,  
35 P., 2012. Ordering effects and choice set awareness in repeat-response stated preference studies.  
36 *Journal of Environmental Economics and Management*, 63(1):73-91.
- 37 Dekker, T., Koster, P., and Brouwer, R., 2014. Changing with the Tide: Semiparametric Estimation of  
38 Preference Dynamics. *Land Economics*, 90(4):717-745.
- 39 Directorate General of the State Forests, 2011. Bialowieza primeval forest - restoration. Warsaw.
- 40 Dupont, D., Price, J., and Adamowicz, W., 2014. Temporal stability of water quality values across stated  
41 preference question formats. paper presented at the Fifth World Congress of Environmental and  
42 Resource Economists, 28 June - 2 July 2014, Istanbul, Turkey.
- 43 Ferrini, S., and Scarpa, R., 2007. Designs with a priori information for nonmarket valuation with choice  
44 experiments: A Monte Carlo study. *Journal of Environmental Economics and Management*,  
45 53(3):342-363.
- 46 Krinsky, I., and Robb, A. L., 1986. On approximating the statistical properties of elasticities. *The Review  
47 of Economics and Statistics*, 68(4):715-719.
- 48 Liebe, U., Meyerhoff, J., and Hartje, V., 2012. Test–Retest Reliability of Choice Experiments in  
49 Environmental Valuation. *Environmental and Resource Economics*, 53(3):389-407.
- 50 Loomis, J. B., 1989. Test-Retest Reliability of the Contingent Valuation Method: A Comparison of  
51 General Population and Visitor Responses. *American Journal of Agricultural Economics*,  
52 71(1):76-84.

- 1 MacMillan, D., Hanley, N., and Lienhoop, N., 2006. Contingent valuation: Environmental polling or  
2 preference engine? *Ecological Economics*, 60(1):299-307.
- 3 McFadden, D., 1974. Conditional Logit Analysis of Qualitative Choice Behaviour. In: *Frontiers in*  
4 *Econometrics*, P. Zarembka, ed., Academic Press, New York, NY, 105-142.
- 5 McFadden, D., 2001. Economic Choices. *The American Economic Review*, 91(3):351-378.
- 6 Mørkbak, M. R., and Olsen, S. B., 2014. A within-sample investigation of test–retest reliability in choice  
7 experiment surveys with real economic incentives. *Australian Journal of Agricultural and*  
8 *Resource Economics*:n/a-n/a.
- 9 Munro, A., and Hanley, N. D., 2002. Information, Uncertainty, and Contingent Valuation. In: *Valuing*  
10 *Environmental Preferences*, I. J. Bateman and K. G. Willis, eds., Oxford University Press.
- 11 Navarro, L. M., and Pereira, H. M., 2012. Rewilding Abandoned Landscapes in Europe. *Ecosystems*,  
12 15(6):900-912.
- 13 Revelt, D., and Train, K., 1998. Mixed Logit with Repeated Choices: Households' Choices of Appliance  
14 Efficiency Level. *Review of Economics and Statistics*, 80(4):647-657.
- 15 Ryan, M., Netten, A., Skåtun, D., and Smith, P., 2006. Using discrete choice experiments to estimate a  
16 preference-based measure of outcome—An application to social care for older people. *Journal*  
17 *of Health Economics*, 25(5):927-944.
- 18 Scarpa, R., and Rose, J. M., 2008. Design Efficiency for Non-Market Valuation with Choice Modelling:  
19 How to Measure it, What to Report and Why. *Australian Journal of Agricultural and Resource*  
20 *Economics*, 52(3):253-282.
- 21 Schaafsma, M., Brouwer, R., Liekens, I., and De Nocker, L., forthcoming. Temporal stability of  
22 preferences and willingness to pay for natural areas in choice experiments: A test-retest.  
23 *Resource and Energy Economics*.
- 24 Skjoldborg, U. S., Lauridsen, J., and Junker, P., 2009. Reliability of the Discrete Choice Experiment at  
25 the Input and Output Level in Patients with Rheumatoid Arthritis. *Value in Health*, 12(1):153-  
26 158.
- 27 Street, D. J., and Burgess, L., 2007. The Construction of Optimal Stated Choice Experiments: Theory  
28 and Methods. Wiley-Interscience, Hoboken, NJ.
- 29 Street, D. J., Burgess, L., and Louviere, J. J., 2005. Quick and easy choice sets: Constructing optimal  
30 and nearly optimal stated choice experiments. *International Journal of Research in Marketing*,  
31 22(4):459–470.
- 32 Swait, J., and Adamowicz, W., 2001. Choice Environment, Market Complexity, and Consumer  
33 Behavior: A Theoretical and Empirical Approach for Incorporating Decision Complexity into  
34 Models of Consumer Choice. *Organizational Behavior and Human Decision Processes*,  
35 86(2):141-167.
- 36 Swait, J., and Louviere, J., 1993. The Role of the Scale Parameter in the Estimation and Comparison of  
37 Multinomial Logit Models. *Journal of Marketing Research*, 30(3):305-314.
- 38 Train, K. E., and Weeks, M., 2005. Discrete Choice Models in Preference Space and Willingness-to-  
39 pay Space. In: *Applications of Simulation Methods in Environmental and Resource Economics*,  
40 R. Scarpa and A. Alberini, eds., Springer, Dordrecht, 1-16.
- 41 Wesolowski, T., 2005. Virtual Conservation: How the European Union is Turning a Blind Eye to Its  
42 Vanishing Primeval Forests. *Conservation Biology*, 19(5):1349-1358.
- 43 Whittington, D., Smith, V. K., Okorafor, A., Okore, A., Liu, J. L., and McPhail, A., 1992. Giving  
44 respondents time to think in contingent valuation studies: A developing country application.  
45 *Journal of Environmental Economics and Management*, 22(3):205-225.

46

1 **Annex A**



2  
3 Figure A1. Dynamics of the number of decision changes in different choice tasks for respondents who  
4 did not participate in the second survey (set A vs. B;  $n = 211$ )

5  
6 Table A1. Kruskal-Wallis test statistics.

	Share of unchanged decisions	Share of unchanged SQ decisions	Share of unchanged NSQ decisions
1	0.2924	1.8592	0.5748
2	0.1092	2.4483	2.8132*
3	0.0145	2.3435	1.5943
4	0.8821	3.5247*	0.8615
5	3.4101*	3.8423**	0.0584
6	0.4526	1.0896	0.1250
7	0.0848	2.5477	1.6712
8	0.1855	2.4935	4.2951**
9	0.7054	2.3143	0.5708
10	0.1917	1.2374	2.1369
11	0.3808	1.8733	0.5118
12	0.2426	0.3290	0.0092

8