

A Formal and Empirical Comparison of Two Score Measures for Best-Worst
Scaling

A. A. J. Marley^{*,1,2}, T. Islam³ and G. E. Hawkins⁴

*(corresponding author)

¹Department of Psychology
University of Victoria
PO Box 1700 STN CSC
Victoria
BC V8W 2Y2
Canada
email: ajmarley@uvic.ca

²Research Professor
Institute for Choice
UniSA Business School
Level 13, 140 Arthur Street
North Sydney
NSW 2060
Australia

³Department of Marketing
and Consumer Studies
College of Business and Economics
University of Guelph
50 Stone Road East
Guelph
ON N1G 2W1
Canada
email: islam@uoguelph.ca

⁴Amsterdam Brain and Cognition Center
University of Amsterdam
Nieuwe Actergracht 129
1018 WS Amsterdam
The Netherlands
email: guy.e.hawkins@gmail.com

Abstract

Best-worst scaling (BWS) is a method that asks individuals to choose the most and the least preferred option from a set of available options. There has been extensive discussion and evaluation of the use of *scores* (data summaries) in the analysis of such data. Here we motivate, summarize, and compare the usefulness of two such score measures: the *analytical closed form solution* (Lipovetsky & Conklin, 2014, *Journal of Choice Modelling*) and *normalized best-worst scores* (Louviere, Flynn, & Marley, 2015, *Cambridge University Press*). We conclude that both have underlying motivations in the *maxdiff model* of best-worst choice and that the analytical closed form solution provides better fits to the aggregate choices in several best-worst choice data sets.

Keywords: best-worst scaling; best-worst choice; discrete choice experiment; scores; choice model

1 Introduction

Best-worst scaling (BWS) is a method that asks individuals to choose the most and least preferred option from a set of available options. Best-worst choice data is obtained using a design that leads to the presentation of each of several subsets of options, with each possible option being presented a number of times (within or between participants). Each participant chooses the best, and the worst, option in each presented set; thus, each option is chosen a number of times as best, as worst, or neither. Such count data have been summarized in various forms, conventionally known as *scores*. Recently, Lipovetsky and Conklin (2014; abbreviated to L&C in various places in the remainder of the paper) developed a closed-form score for best-worst choice, henceforth referred to as the *analytical best-worst (ABW) score (estimator)*. For each option x in the study (design), let N_x be the number of times x is presented, N_x^b the number of times x is selected as best, and N_x^w the number of times x is selected as worst; therefore $N_x - N_x^b - N_x^w$ is the number of times x is not selected as either best or worst. L&C develop the ABW score

$$\ln \frac{1 + \frac{N_x^b - N_x^w}{N_x}}{1 - \frac{N_x^b - N_x^w}{N_x}}, \quad (1)$$

where the quantity

$$\frac{N_x^b - N_x^w}{N_x} \quad (2)$$

is the *normalized best minus worst (NBW) score* for x . The latter closed-form measure has been used frequently in the literature on best-worst choice, where it is often linearly related to the estimated parameters of various models (Louviere, Flynn, & Marley, 2015) plus it has numerous interesting theoretical properties (summarized by Flynn & Marley, 2014). For instance, the set of (normalized) best minus worst scores is a sufficient statistic for the parameters of the *maxdiff model of best-worst choice* [below, Sect. 3, (9)].¹ Since the function in (1) is a strictly increasing function of the NBW score in (2), the set of ABW scores is also a sufficient statistic for that model.

Figure 1 shows the mathematical relation between ABW, (1), and NBW, (2), scores.² Our aim in this paper is to further understand the properties of, and relations between, these two measures, and to evaluate each against data. Our conclusions are that the maxdiff model of best-worst choice underlies the usefulness of each of the measures; and that the ABW measure gives better prediction of (aggregate) best (respectively, worst, best-worst) choices than does

¹The term *maxdiff* was coined for this model as its representation can be described as a result of a decision maker choosing as the best-worst pair in a set the pair of options with the largest (random) utility difference, with the random component having an extreme value distribution.

²The ABW measure is clearly highly linearly related to the NBW measure for much of the latter's range. For instance, with NBW restricted to $(-.5, .5)$, the slope is 2.11 with an R^2 of .999.

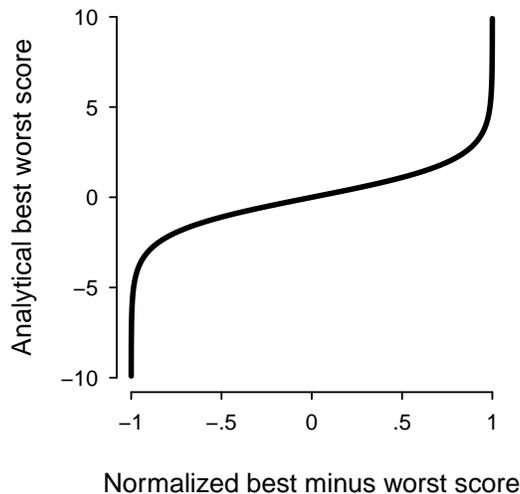


Figure 1: The relationship between the analytical best worst score (y -axis) and the normalized best minus worst score (x -axis).

33 the NBW score, when the (aggregate) best (resp., worst) choices are the relevant
 34 marginals of (aggregate) best-worst choices.

35 The remainder of the paper is as follows. Section 2 completely redevelops
 36 L&C’s ABW score measure in a way that allows us, in Section 3, to show the
 37 measure’s relationship with the likelihood of data generated by the maxdiff
 38 model of best-worst choice [below, (9)]; this link was not made in L&C. The
 39 corresponding previously known link for the NBW measure essentially puts
 40 the two measures on an equal footing. Then, in Section 4, we evaluate each
 41 measure against data. Section 5 summarizes and discusses further results and
 42 open questions.

43 2 Lipovetsky and Conklin’s (2014) Matrix of Ex- 44 tended Best-Worst Data

45 As did L&C, we analyze best-worst choice data aggregated across participants –
 46 that is, for each choice set in the design, we calculate the proportion of partici-
 47 pants who chose a particular pair of options as best-worst; we also calculate the
 48 appropriate marginals of those best-worst proportions to obtain the (marginal)
 49 best (resp., worst) choice proportions. Nonetheless, we describe the methods as
 50 if the data are for an individual decision maker. We present the results in terms
 51 of choice sets with four options; they are easily extended to other set sizes.

52 Let P denote a set of items (17 in L&C) and let $D(P)$ denote the design, that

53 is, the set of (sub)sets of items that occur in the study; without loss of generality,
 54 we assume each set in $D(P)$ has four items. We now present an argument for
 55 the *extended best-worst data* and *choice probabilities* of Table 1, followed by a
 56 demonstration that L&C's ABW score measures are the maximum likelihood
 57 estimates for the assumed model of the extended data. Along the way, we point
 58 out the close link between our and L&C's framework. Section 3 extends our
 59 approach to a process model of the actual best-worst data, in contrast to the
 60 modeling of the extended best-worst data considered in this section.

61 We assume that each person has a (common) referent item ϕ across all
 62 choice sets in the design; the referent is an inferred status quo option, not a
 63 presented option. Assume that the choice set $\{e, f, g, h\}$ is presented to the
 64 decision maker and e is chosen as best, h as worst. These observed choices are
 65 used to produce the *extended choice sets* and *implied choices For Best (resp.,*
 66 *For Worst)* in the first column of Table 1. The underlying concept is that
 67 each option is independently compared with the referent ϕ , under a “best”
 68 condition, and, separately, under a “worst” condition. When selecting the best
 69 option, the observed choice of e best is interpreted as corresponding to e being
 70 better than ϕ , and ϕ being better than each of f, g , and h . When selecting the
 71 worst option, the observed choice of h worst is interpreted as corresponding to
 72 h being worse than ϕ , and ϕ being worse than each of e, f , and g . We then
 73 assume that, for each item x in the set $\{e, f, g, h\}$, $b(x, \phi)$ is the probability
 74 that x is better than ϕ , and $w(\phi, x)$ is the probability that ϕ is worse than
 75 x ; this gives the final column of Table 1. As we discuss later, in detail, there
 76 are many possible combinations of choices that do *not* occur under the assumed
 77 expansion process; the latter combinations will be considered in Section 3, where
 78 we present a *process* interpretation of the extended best-worst data.

79 Table 1. Component binary choice probabilities on extended best-worst choice
 80 data on a set $\{e, f, g, h\}$, when e is chosen as best, h is chosen as worst; ϕ is a
 81 hypothetical referent item.

	Choice Set					
	e	f	g	h		
Observed Best choice	e					
Observed Worst choice				h		
Extended Choice Sets:					Implied Choices:	Choice Probability:
For Best						
e	1	0	0	0	e	$b(e, \phi)$
f	0	1	0	0	ϕ	$b(\phi, f)$
g	0	0	1	0	ϕ	$b(\phi, g)$
h	0	0	0	1	ϕ	$b(\phi, h)$
For Worst						
e	1	0	0	0	ϕ	$w(\phi, e)$
f	0	1	0	0	ϕ	$w(\phi, f)$
g	0	0	1	0	ϕ	$w(\phi, g)$
h	0	0	0	1	h	$w(h, \phi)$

83 Now assume that, for each item x in the design P , there is a finite real
 84 (utility) value $u(x)$, and, with no loss of generality, set $u(\phi) = 0$. We also
 85 assume that

$$b(x, \phi) = \frac{e^{u(x)}}{1 + e^{u(x)}} \text{ and } w(\phi, x) = \frac{1}{1 + e^{-u(x)}}. \quad (3)$$

86 These equations imply that³: $b(x, \phi) = w(\phi, x)$; the probabilities are in the open
 87 interval $(0, 1)$ for real values of $u(x)$; and

$$u(x) = \ln \frac{b(x, \phi)}{1 - b(x, \phi)} = \ln \frac{w(\phi, x)}{1 - w(\phi, x)}.$$

88 We show, below, that the assumptions in (3), with those in the *data expansion*
 89 of Table 1, lead to L&C's ABW score measure⁴.

90 The probability of the pattern of extended choices in Table 1 is

$$b(e, \phi)b(\phi, f)b(\phi, g)b(\phi, h)w(\phi, e)w(\phi, f)w(\phi, g)w(h, \phi).$$

91 By (3), we have $w(\phi, x) = b(x, \phi)$ for each item x , so the expression becomes

$$b(e, \phi)b(\phi, f)b(\phi, g)b(\phi, h)b(e, \phi)b(f, \phi)b(g, \phi)b(\phi, h), \quad (4)$$

92 which equals

$$b(e, \phi)^2b(\phi, f)b(f, \phi)b(\phi, g)b(g, \phi)b(\phi, h)^2. \quad (5)$$

93 *Comment:* The sum of (4) over all possible patterns allowed by the reasoning
 94 leading to Table 1 for the choice set $X = \{e, f, g, h\}$ can be less than 1 (see
 95 **Appendix A**). The fact that the sum can be less than one means that those
 96 expressions do not form a process model of the actual choices in a standard
 97 best-worst task, where the decision maker has to select a best-worst pair. This
 98 is not an issue at this stage as we are presenting an argument that leads to
 99 L&C's ABW score measure. Section 3 presents a process model of best-worst
 100 choice suggested by the logic of Table 1 that does not have this defect.

101 We now continue the reasoning based on Table 1 and show that it leads to
 102 the ABW score measure (parameter estimates) proposed by L&C.

103 For each (four element) set $X \in D(P)$, and $x, y \in X$, $x \neq y$, let $\widetilde{bw}_X(x, y)$
 104 be the aggregate number of times the best-worst choice pair in X is x best and

³There is ongoing empirical study of when the “mirror image” property of best and worst holds, including in the binary case of (3). The currently available data are summarized and discussed in Louviere et al. (2015, Sect. 6.2) and Section 5 of the current paper considers the relevance of our work to those discussions.

⁴This is because our referent ϕ corresponds to the 0's in L&C's Table 4, and our formulae (3) correspond to L&C's pairwise logistic formulae (their (5)).

105 y worst⁵; clearly, $\widetilde{bw}_X(x, y) = 0$ if either x or y is not in X . Then the likelihood
 106 of the extended best-worst choice data of Table 1 is

$$\prod_{X \in D(P)} \prod_{\substack{\{e, h\} \in X \\ e \neq h}} [b(e, \phi)^2 b(\phi, f) b(f, \phi) b(\phi, g) b(g, \phi) b(\phi, h)^2]^{\widetilde{bw}_X(e, h)} \quad (6)$$

107 Now consider a generic option $x \in P$ and the sets $X \in D(P)$ in which x is
 108 present. In agreement with L&C's notation, let N_x be the number of times x is
 109 presented in the design $D(P)$, and let

$$N_x^b = \sum_{X \in D(P)} \sum_{r \in X - \{x\}} \widetilde{bw}_X(x, r), \quad N_x^w = \sum_{X \in D(P)} \sum_{s \in X - \{x\}} \widetilde{bw}_X(s, x);$$

110 that is, x is best N_x^b times, worst N_x^w times, and neither best or worst $N_x -$
 111 $N_x^b - N_x^w$ times.

112 Collecting the terms for sets in which x appears in the likelihood (6), we
 113 obtain the expression

$$\begin{aligned} & b(x, \phi)^{2N_x^b} b(\phi, x)^{2N_x^w} [b(\phi, x) b(x, \phi)]^{N_x - N_x^b - N_x^w} \\ &= b(x, \phi)^{N_x + (N_x^b - N_x^w)} b(\phi, x)^{N_x - (N_x^b - N_x^w)} \\ &= b(x, \phi)^{N_x + (N_x^b - N_x^w)} [1 - b(x, \phi)]^{N_x - (N_x^b - N_x^w)} \end{aligned} \quad (7)$$

114 Differentiating the log of this likelihood w.r.t $b(x, \phi)$, and setting the result to
 115 zero, gives

$$\frac{N_x + [N_x^b - N_x^w]}{b(x, \phi)} = \frac{N_x - [N_x^b - N_x^w]}{1 - b(x, \phi)},$$

116 which, with the form of $b(x, \phi)$ from (3), gives

$$u(x) = \ln \frac{b(x, \phi)}{1 - b(x, \phi)} = \ln \frac{N_x + [N_x^b - N_x^w]}{N_x - [N_x^b - N_x^w]} = \ln \frac{1 + \frac{N_x^b - N_x^w}{N_x}}{1 - \frac{N_x^b - N_x^w}{N_x}},$$

117 that is, the ABW estimator.

⁵With N_X the number of times the set X is presented in the design $D(P)$, we have

$$N_X = \sum_{\substack{x, y \in X \\ x \neq y}} \widetilde{bw}_X(x, y).$$

118 **3 A Process Model Reinterpretation of Lipovet-**
119 **sky and Conklin’s (2014) Matrix of Extended**
120 **Best-Worst Data**

121 Section 2 presented our development of L&C’s ABW measure using their ex-
122 tended matrix form, which is based on the original best-worst data. We now
123 discuss why their solution gives an approximation, not only to the maximum
124 likelihood parameter estimates for a model of the best-worst data, as stated by
125 L&C, but also to the likelihood of the original best-worst data.

126 **Appendix A** shows that the formulae (4) do not form a complete model
127 of the original best-worst data because the sum of those expressions over each
128 choice set $\{e, f, g, h\}$ is not equal to one. We now extend the process that leads
129 to (4) in the most natural way, and show that the resulting model is the *maxdiff*
130 *model* of best-worst choice.⁶

131 The process described in Table 1 can be extended to a model of the actual
132 best-worst data as follows: the respondent independently compares each option
133 in $\{e, f, g, h\}$ with the referent ϕ in terms of which is better (the option or the
134 referent); if the results of all four comparisons is exactly one option being better
135 than ϕ , call it x , then x is noted as “best”. The respondent then carries out the
136 parallel comparison of each option in $\{e, f, g, h\}$ with the referent ϕ in terms of
137 which is worse (the option or the referent); if the result of all four comparisons
138 is exactly one option being worse than ϕ , call it y , then y is noted as “worst”.
139 In such a case, and $x \neq y$, these (single items) become the best, x , and worst,
140 y , choices for that set. However, if the process does *not* result in a single best
141 x and worst y , $x \neq y$, the respondent restarts the whole process.

142 Continuing with $X = \{e, f, g, h\}$, let $BW_X(e, h)$ denote the probability that
143 e is chosen as best, and h is chosen as worst, $e \neq h$, by the above process. Then
144 a standard argument (paralleling Marley & Louviere, 2005, Sect. 4.1.2, Case 2),
145 with the form (5) for the (unconditional) selection of the best-worst pair (e, h) ,
146 gives

$$BW_X(e, h) = \frac{b(e, \phi)^2 b(\phi, f) b(f, \phi) b(\phi, g) b(g, \phi) b(\phi, h)^2}{\sum_{\{r, s, t, u\} = X} b(r, \phi)^2 b(\phi, s) b(s, \phi) b(\phi, t) b(t, \phi) b(\phi, u)^2}. \quad (8)$$

147 Note that, in the denominator, every set $\{r, s, t, u\} = X = \{e, f, g, h\}$. Dividing
148 the numerator and denominator by

$$b(\phi, e) b(e, \phi) b(\phi, f) b(f, \phi) b(\phi, g) b(g, \phi) b(\phi, h) b(h, \phi),$$

149 which is symmetric (i.e., independent of the positions of e, f, g, h), and nonzero
150 by (3) with the assumption of finite valued u , gives the form

⁶The following extension of the L&C approach is not the simplest process interpretation of the maxdiff model; see Marley & Louviere (2005, Sect. 4.1.2, Case 2).

$$BW_X(e, h) = \frac{\frac{b(e, \phi)}{\bar{b}(\phi, e)} \cdot \frac{b(\phi, h)}{\bar{b}(h, \phi)}}{\sum_{\substack{\{r, u\} \in X \\ r \neq u}} \frac{b(r, \phi)}{\bar{b}(\phi, r)} \cdot \frac{b(\phi, u)}{\bar{b}(u, \phi)}}.$$

151 Substituting in the forms (3), this reduces to

$$BW_X(e, h) = \frac{e^{[u(e) - u(h)]}}{\sum_{\substack{\{r, u\} \in X \\ r \neq u}} e^{[u(r) - u(u)]}}, \quad (9)$$

152 which is the *maxdiff model of best worst choice* (Marley & Louviere, 2005).⁷
 153 The set of NBW scores for the options is a sufficient statistic for the maxdiff
 154 model (Marley & Pihlens, 2012); thus, as already noted in Section 1, the set of
 155 ABW scores is also a sufficient statistic for that model.

156 For each (four element) set $X \in D(P)$, and $x, y \in X$, $x \neq y$, let $\widetilde{bw}_X(x, y)$
 157 be the aggregate number of times the best-worst choice pair in X is given by x
 158 best and y worst; also let N_X be the number of times the set X occurs in the
 159 design.⁵ Then, according to the model in (8), the likelihood of the data is

$$\begin{aligned} & \prod_{X \in D(P)} BW_X(e, h) \\ &= \prod_{X \in D(P)} \prod_{\substack{\{e, h\} \in X \\ e \neq h}} \left(\frac{b(e, \phi)^2 b(\phi, f) b(f, \phi) b(\phi, g) b(g, \phi) b(\phi, h)^2}{\sum_{\{r, s, t, u\} = X} b(r, \phi)^2 b(\phi, s) b(s, \phi) b(\phi, t) b(t, \phi) b(\phi, u)^2} \right)^{\widetilde{bw}_X(e, h)} \\ &= \prod_{X \in D(P)} \prod_{\substack{\{e, h\} \in X \\ e \neq h}} [b(e, \phi)^2 b(\phi, f) b(f, \phi) b(\phi, g) b(g, \phi) b(\phi, h)^2]^{\widetilde{bw}_X(e, h)} \\ & \quad \times \prod_{X \in D(P)} \left(\frac{1}{\left[\sum_{\{r, s, t, u\} = X} b(r, \phi)^2 b(\phi, s) b(s, \phi) b(\phi, t) b(t, \phi) b(\phi, u)^2 \right]^{N_X}} \right) \end{aligned}$$

160 Comparing this likelihood function with the corresponding form (6), we see
 161 that the numerator is exactly the term in (6), and the denominator does not
 162 appear in (6). Thus, L&C's derivation of their score measure does not include
 163 the terms given by the denominators of the terms in (8) in the likelihood func-
 164 tion. Such an approach is sometimes advocated, based on the statement that
 165 the product of such denominator terms in the likelihood function can be treated
 166 as a "constant". This is incorrect, as that product depends on the values of
 167 the options in the design; and ignoring them may lead to bias in the estimation
 168 of parameters and incorrect fit measures based on log-likelihoods. Nonetheless,
 169 as we show in the next section, the ABW estimator performs extremely well

⁷The maxdiff model can also be developed in a similar manner using L&C's extended matrix method, but with best and worst extended separately, rather than together (see Marley & Louviere, 2005, Sect. 4.1.2. Case 2).

170 in fitting aggregate (marginal) best choice data derived from best-worst choice
 171 data.

172 4 Empirical Comparison of Analytical Best-Worst 173 Scores and Normalized Best-Minus-Worst Scores

174 We tested the ability of the ABW and NBW scores to fit the aggregate (marginal)
 175 best choice probabilities derived from numerous sets of best-worst data obtained
 176 in discrete choice experiments; the details of the designs are provided in **Ap-**
 177 **pendix B.** Four studies used Case 3 best-worst choice (choice between pro-
 178 files), and two used Case 1 best-worst choice (choice between objects).⁸ Table
 179 2 shows the root mean square error (RMSE) between the observed (aggregate)
 180 best choices and the (aggregate) best choices predicted using the NBW (resp.,
 181 ABW) scores as the parameters in a multinomial logit (MNL) model. In every
 182 case, the ABW scores had smaller RMSE; in almost all cases, the RMSE using
 183 the NBW scores was about twice the size of that using the ABW scores.

184 Table 2. Root Mean Square Error (RMSE) of in-sample prediction of
 185 aggregate (marginal) best choices using the NBW score (resp., ABW score).
 186 The RMSE of a set of L predictions (number of choice sets times alternatives
 187 in each choice set) is calculated as $RMSE = \sqrt{\frac{\sum_1^L (Best_{obs} - Best_{pred})^2}{L}}$. W_i
 188 denotes wave i of the design, where $i = 1, 2, 3, 4$.

	Product	Measure	RMSE			
			W ₁	W ₂	W ₃	W ₄
Case 3	Detergent	NBW	.095	.098	.100	.102
		ABW	.046	.039	.042	.044
	Toothpaste	NBW	.111	.111	.111	.103
		ABW	.054	.054	.054	.053
	Pizza	NBW	.101	.098	.103	.101
		ABW	.037	.036	.040	.041
	Solar	NBW	.096			
		ABW	.030			
Case 1	Budget Saving	NBW	.024			
		ABW	.014			
	Budget Spending	NBW	.021			
		ABW	.012			

190 Figure 2 plots a typical pair of results from Table 2. The range of predicted
 191 (marginal) best choice probabilities using NBW scores is smaller than the cor-
 192 responding range using ABW scores, and of the data. This is largely due to

⁸In Case 3, a person is presented with a set of profiles (multiattribute options), such as possible vacation packages, and has to choose the best and worst of those. In Case 1, a person is presented with several objects (e.g., brand names), and has to choose the best and worst object. See Marley and Louviere (2005, Sect. 1.3) for details on the three cases of best-worst scaling.

193 the fact that the ABW scores are approximately two to three times larger than
 194 the corresponding NBW scores when the latter are not at the extremes of their
 195 range; that is, not near -1 or 1. (cf. Figure 1).

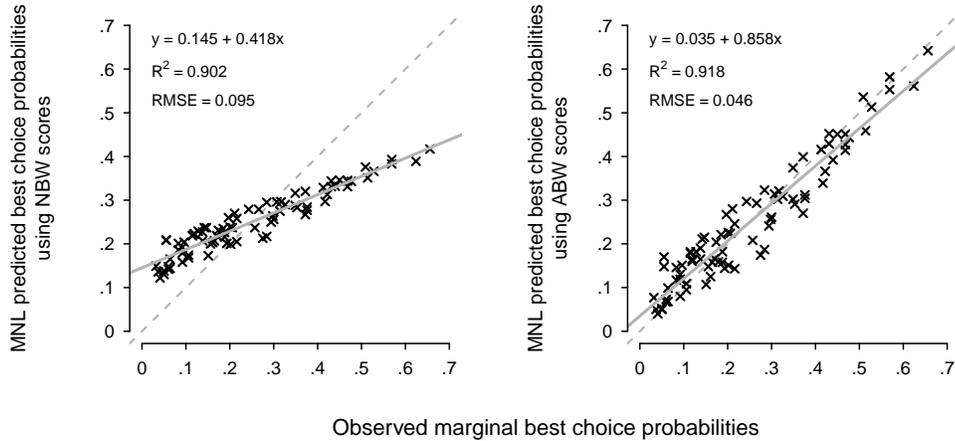


Figure 2: Observed (marginal) best choice probabilities from wave 1 of the detergent data set (x -axes) versus best choice probabilities predicted by an MNL model using the NBW scores (left panel) or ABW scores (right panel) as the (utility) parameters (y -axes). Solid lines show linear regression lines of best fit. Dashed lines show where the data would fall if the two variables were perfectly related.

196 Given the plot in Figure 1 and the results in Table 2, one would expect
 197 the following regarding the relation of each score measure to the parameters
 198 of a maximum likelihood fit of an MNL model to the best choice probabilities:
 199 For each data set, plot the ABW (resp., NBW) scores on the x -axis and the
 200 estimated parameters of the MNL model on the y -axis. Then the slope of the
 201 linear regression will be lower and closer to 1 for each ABW plot. We have
 202 carried out the relevant analyses, and this is the case for every data set in Table
 203 2.

204 5 Summary, Further Results, and Open Questions

205

206 L&C focussed on predicting in-sample (marginal) best choices, as we have done
 207 thus far. However, given the close relation we have demonstrated between the
 208 ABW scores, the NBW scores, and the maxdiff model, it is worth extending
 209 the analyses to (marginal) worst choices, and best-worst choices. The results of
 210 these analyses are summarized in **Appendix B**. The in-sample fits to (marginal)
 211 worst choices are as good as those already presented for (marginal) best choices;

212 and, if anything, the corresponding fits to the best-worst choices are better than
213 those to the (marginal) best or (marginal) worst choices⁹; in all cases, the ABW
214 scores perform better than the NBW scores. Three of the best-worst data
215 sets had four waves of data, which allowed us to test the ability of ABW and
216 NBW scores to predict out-of-sample (marginal) best, (marginal) worst, and
217 best-worst choices. These predictive tests supported a similar conclusion to the
218 in-sample tests (see **Appendix B**). Thus, these results suggest that, for these
219 aggregate data, best and worst choices are “mirror images” of each other.

220 Our results are very encouraging in terms of using the ABW scores as de-
221 scriptive statistics for aggregate best-worst choice probabilities, and they build
222 on earlier uses of the NBW scores (see Louviere et al., 2015). Nonetheless, we
223 do not have a clear understanding of why these measures work so well, other
224 than that each gives a sufficient statistic for the maxdiff model. Also, there is
225 a growing empirical literature demonstrating that best choices do not always
226 provide the same information (demonstrate the “same” preferences) as worst
227 choices; examples are the model-based analyses of Giergiczny et al. (2013) and
228 the purely data-based analyses of Greenacre et al (2015). Rose (2014) proposes
229 that these differences are not an issue for the best-worst method if one is inter-
230 ested in choices processes, per se, but can be if best and worst data are naively
231 pooled to, say, explain or predict best choices. Louviere et al. (2015, Chap. 6)
232 summarize and discuss various of these data and the issues they raise. Also,
233 there is an increasing tendency to model the (best-worst) choices of individual
234 decision makers by collecting more data from each individual (e.g., Hawkins et
235 al., 2014) and/or using Bayesian methods when there are a limited amount of
236 data for each individual (e.g., Dumont et al., 2015). Score measures (such as
237 ABW and NBW) play an important role in both exploratory and final analyses
238 in these situations (see, e.g., Louviere, 2015). In particular, our results show
239 that the ABW scores may provide a closer approximation to the relevant data
240 than NBW scores.

241 As a final consideration, assume that we have a maxdiff model for choice
242 among profiles (Case 3, choice between profiles) where the utility of an option
243 is the sum of the (component) utilities of its attribute-levels. Count the number
244 of times an attribute-level is “chosen”¹⁰ as best minus the number of times it
245 is “chosen” as worst, and normalize that number by the number of times the
246 attribute-level appears in the design. Then Marley and Pihlens (2012) show
247 that, given the maxdiff model, these score differences are a sufficient statis-
248 tic. Future work should study how these attribute-level NBW scores (and the
249 corresponding ABW scores) are related to maximum likelihood estimates of
250 attribute-level utilities (results in Louviere et al., 2015, suggest each relation
251 will be linear) and how well each score measure predicts (aggregate or individ-
252 ual) choice probabilities.

⁹These fits are consistent with the fact that, if the maxdiff model is the “true” model of the best-worst choices, then an MNL model is an approximation to the (marginal) best (resp., worst) choices (Marley & Louviere, 2005).

¹⁰Of course, in Case 3, an attribute-level is “chosen” (only) as a consequence of its being a component of a profile chosen on a given choice opportunity.

Acknowledgments

254 This research has been supported by Natural Science and Engineering Re-
 255 search Council Discovery Grant 8124-98 to the University of Victoria for Marley
 256 and by Social Sciences and Humanities Research Council (SSHRC) Grant No.
 257 430060 for Islam. The work was carried out whilst Marley was a Research Pro-
 258 fessor (part-time) in the Institute for Choice, University of South Australia
 259 Business School.

Appendix A

261 To show that the sum of (4) over all possible patterns allowed by the reason-
 262 ing leading to Table 1 can be less than 1, assume every binary choice probability
 263 is in the open interval $(0, 1)$ and consider a participant making the assumed *in-*
 264 *dependent* choices in each of the eight (four best; four worst) extended choice
 265 sets in Table 1. Then one possible pattern of results has the probability

$$b(e, \phi)b(\phi, f)b(\phi, g)b(\phi, h)w(e, \phi)w(\phi, g)w(f, \phi)w(h, \phi).$$

266 However, this pattern does not arise under the logic of the extended choice sets
 267 in Table 1. Therefore, the sum of the probabilities of the terms in (4) over all
 268 the patterns allowed by the logic of Table 1 is less than 1 – as they are a subset
 269 of all the possible independent patterns. Hence those terms do not correspond
 270 to a model of the original best-worst choice data.

Appendix B

B1. The Discrete Choice Experiments (DCEs)

272 The Laundry Detergent, Toothpaste, Delivered Pizza, and Solar Panel data
 273 sets involved DCEs (best-worst Case 3, choice between profiles) using an exper-
 274 imental design suitable for identifying particular forms of indirect utility func-
 275 tion. In particular, these studies used a recently proposed approach to DCE
 276 design that involves two stages: 1) construct a set of product profiles (i.e., at-
 277 tribute level combinations, product descriptions) using a suitable experimental
 278 design to allow identification of particular forms of indirect utility functions of
 279 interest (e.g., additive), and 2) use a B(alanced) I(ncomplete) B(lock) D(esign)
 280 (BIBD) to assign the profiles to choice sets of fixed size (Louviere et al., 2008).
 281 For each of the four product categories, there were 9 attributes, 3 of which had
 282 4 levels and the remaining 6 each had 2 levels. Specifically, the design used an
 283 orthogonal main effects plan to design (construct) 16 profiles (choice options);
 284 this was a purposive sample of a $4^3 \times 2^6 = 4096$ factorial that allowed estimation
 285 of attribute main effects, assuming that all interactions were non-significant. A
 286 BIBD then assigned the 16 profiles to 20 choice sets, each of which had four
 287 choice options. For each choice set, respondents first selected the best profile
 288 from the 4 available options; they then selected the worst option from the re-
 289 remaining 3 options. For the detergent, toothpaste, and pizza data sets there were,
 290 respectively, 218, 234, and 186 respondents who each completed four waves of
 291

292 data collection across a two-year period. The solar data set consisted of 298
293 respondents who completed a single wave of data collection. In each of these
294 designs, the option selected as best remained on the screen whilst the worst op-
295 tion was being chosen; but at this stage, the best option could not be changed.
296 Further analyses of these data appear in Islam (2014) and Islam and Louviere
297 (2015).

298 For Budget Saving and Budget Spending (best-worst Case 1, choice between
299 objects), a BIBD was used to assign the 9 attributes of the study to 12 choice
300 sets each having three choice options. For each choice set, respondents first
301 selected the best from the 3 available attributes, then selected the worst from
302 the remaining 2 attributes. The same set of 561 respondents completed the
303 budget saving and budget spending designs. In each of these designs, the option
304 selected as best was removed from the screen before the worst option was chosen.
305 Further analyses of these data appear in Louviere et al (2015, Chap. 2).¹¹

306 **B2. In-Sample and Out-Of-Sample Fit to Data Using NBW (Resp.,** 307 **ABW) Scores**

308 Table B1 (resp., B2) shows the RMSE for in-sample prediction of aggregate
309 (marginal) worst (resp., best-worst) choices in a form that mirrors Table 2 of
310 the main text. Table B1 uses *minus* the NBW (resp., *minus* the ABW) scores
311 as the parameters in a multinomial logit (MNL) model. Table B2 estimates the
312 utility difference in a maxdiff model for a pair of options (x, y) , $x \neq y$, by the
313 NBW score for x minus the NBW score for y (resp., by the ABW score for x
314 minus the NBW score for y).

¹¹This data set was collected at the Centre for the Study of Choice, University of Technology Sydney; we thank those involved in its collection for allowing us to analyse it.

315 Table B1. RMSE of in-sample prediction of aggregate (marginal) worst choices
 316 using the NBW score (resp., ABW score). The RMSE of a set of L predictions
 317 (number of choice sets times alternatives in each choice set) is calculated as
 318 $RMSE = \sqrt{\frac{\sum_i^L (Worst_{obs} - Worst_{pred})^2}{L}}$. W_i denotes wave i of the design, where
 319 $i = 1, 2, 3, 4$.

	Product	Measure	RMSE			
			W ₁	W ₂	W ₃	W ₄
Case 3	Detergent	NBW	.097	.091	.098	.092
		ABW	.043	.036	.038	.040
	Toothpaste	NBW	.126	.131	.129	.128
		ABW	.054	.053	.053	.054
	Pizza	NBW	.108	.103	.105	.102
		ABW	.032	.031	.035	.036
	Solar	NBW	.063			
		ABW	.027			
Case 1	Budget Saving	NBW	.027			
		ABW	.017			
	Budget Spending	NBW	.024			
		ABW	.013			

321 Table B2. RMSE of in-sample prediction of aggregate best-worst choices using
 322 the NBW score (resp., ABW score). The RMSE of a set of L predictions
 323 (number of choice sets times $m \times (m - 1)$ alternatives in each choice set) is
 324 calculated as $RMSE = \sqrt{\frac{\sum_i^L (BW_{obs} - BW_{pred})^2}{L}}$. W_i denotes wave i of the
 325 design, where $i = 1, 2, 3, 4$.

	Product	Measure	RMSE			
			W ₁	W ₂	W ₃	W ₄
Case 3	Detergent	NBW	.051	.048	.053	.053
		ABW	.037	.031	.034	.036
	Toothpaste	NBW	.064	.066	.066	.063
		ABW	.043	.045	.046	.045
	Pizza	NBW	.053	.053	.054	.053
		ABW	.031	.031	.034	.034
	Solar	NBW	.040			
		ABW	.028			
Case 1	Budget Saving	NBW	.031			
		ABW	.022			
	Budget Spending	NBW	.025			
		ABW	.017			

327 We also examined out-of-sample prediction error according to RMSE. We
328 used multiple time origins (i.e., waves) to produce one, two, and three step
329 ahead predictions. The first time origin is wave 1 and from this origin one step
330 ahead (six months) to three steps ahead (one and half years) predictions are
331 made. The time origin is advanced by one wave and a set of one and two step
332 ahead predictions are generated. Finally, the origin is advanced to wave 3 and a
333 final one step ahead prediction is generated. Therefore, from the three data sets
334 that each have four waves we can generate three one-step, two two-step, and
335 one three-step ahead predictions. Tables B3, B4 and B5 show the mean RMSE
336 averaged over the three one-step (resp., two two-step, and one three-step) ahead
337 predictions according to NBW and ABW scores for each data set, separately
338 for (marginal) best, (marginal) worst, and best-worst choices. The NBW and
339 ABW scores are used in Table B3 (resp., B4, B5) in exactly the way they are
340 used in Table 2 (resp., Table B1, B2).

341 Table B3. RMSE for out-of-sample prediction of aggregate (marginal) best
342 choices using the NBW score (resp., ABW score).

Product	Measure	RMSE		
		One step	Two step	Three step
Detergent	NBW	.098	.099	.096
	ABW	.049	.055	.060
Toothpaste	NBW	.113	.115	.111
	ABW	.061	.065	.067
Pizza	NBW	.103	.102	.107
	ABW	.048	.050	.053

344 Table B4. RMSE for out-of-sample prediction of aggregate (marginal) worst
345 choices using the NBW score (resp., ABW score).

Product	Measure	RMSE		
		One step	Two step	Three step
Detergent	NBW	.098	.097	.103
	ABW	.049	.053	.058
Toothpaste	NBW	.131	.134	.131
	ABW	.060	.065	.065
Pizza	NBW	.110	.110	.116
	ABW	.049	.047	.049

347 Table B5. RMSE for out-of-sample prediction of aggregate best-worst choices
 348 using the NBW score (resp., ABW score).

Product	Measure	RMSE		
		One step	Two step	Three step
Detergent	NBW	.052	.050	.052
	ABW	.036	.037	.040
Toothpaste	NBW	.067	.069	.072
	ABW	.045	.044	.043
Pizza	NBW	.054	.053	.053
	ABW	.035	.033	.034

349

References

- 351 Dumont, J., Giergiczny, M., & Hess, S. (2015). Individual level models
352 vs. sample level models: Contrasts and mutual benefits. *Transportmetrica A:
353 Transport Science*, 11, 465-483.
- 354 Flynn, T. N. & Marley, A. A. J. Best-worst scaling: practice and theory.
355 Invited chapter in S. Hess & A. Daly (Eds.) *Handbook of Choice Modelling*.
356 Edward Elgar Publishing, 2014, pp. 178-201.
- 357 Giergiczny, M., Chintakayala, P., Dekker, T., & Hess, S. (2013). Testing
358 the consistency (or lack thereof) between choices in best-worst surveys. Paper
359 presented at the *Third International Choice Modelling Conference*, Sydney, July
360 3-5.
- 361 Greenacre, L., Dunn, S., & Mocuano, A. (2015). Heterogeneity in the consistency
362 of best-worst scale responses. *Australian Marketing Journal*. 23, 227-234.
- 363 Hawkins, G. E., Marley, A. A. J., Heathcote, A., Flynn, T. N., Louviere,
364 J. J., & Brown, S. D. (2014). The best of times and the worst of times are
365 interchangeable. *Decision*, 1, 192-214.
- 366 Islam, T. (2014). Household level innovation diffusion model of photo-voltaic
367 (PV) solar cells from stated preference data. *Energy Policy*, 65, 340-350.
- 368 Islam, T., & Louviere, J. J. (2015). The stability of aggregate-level preferences
369 in longitudinal discrete choice experiments. In J. J. Louviere, T. N.
370 Flynn & A. A. J. Marley (2015). *Best-Worst Scaling: Theory, Methods and
371 Applications*. Cambridge University Press. Cambridge: UK. Chap. 13, pp.
372 265-277.
- 373 Lipovetsky, S., & Conklin, M. W. (2014). Best-worst scaling in analytical
374 closed-form solution. *The Journal of Choice Modelling*, 10, 60-68.
- 375 Louviere, J. J. (2015). Using alternative-specific DCE designs and best and
376 worst choices to model choices. In J. J. Louviere, T. N. Flynn & A. A. J. Marley
377 (2015). *Best-Worst Scaling: Theory, Methods and Applications*. Cambridge
378 University Press. Cambridge: UK. Chap. 15, pp. 297-315.
- 379 Louviere, J. J., Flynn, T. N. & Marley, A. A. J. (2015). *Best-Worst Scaling:
380 Theory, Methods and Applications*. Cambridge University Press. Cambridge:
381 UK.
- 382 Louviere, J. J., Street, D. J., Burgess, L., Wasi, N., Islam, T., & Marley, A.
383 A. J. (2008). Modelling the choices of single individuals by combining efficient
384 choice experiment designs with extra preference information. *Journal of Choice
385 Modelling*, 1, 128-163.
- 386 Marley, A. A. J., & Louviere, J. J. (2005). Some probabilistic models of
387 Best, Worst, and Best-Worst choices. *Journal of Mathematical Psychology*, 49,
388 464-480.
- 389 Marley, A. A. J., & Pihlens, D. (2012). Models of best-worst choice and
390 ranking among multiattribute options (profiles). *Journal of Mathematical Psychology*,
391 56, 24-34.
- 392 Rose, J. M. (2014). Interpreting discrete choice models based on best-worst
393 data: A matter of framing. *Transportation Research Board 93rd Annual Meeting*.
394 *Washington, DC. Jan. 12-16.*